

AIX Version 7.2

General programming concepts

IBM

Note

Before using this information and the product it supports, read the information in [“Notices” on page 803.](#)

This edition applies to AIX Version 7.2 and to all subsequent releases and modifications until otherwise indicated in new editions.

© **Copyright International Business Machines Corporation 2015, 2019.**

US Government Users Restricted Rights – Use, duplication or disclosure restricted by GSA ADP Schedule Contract with IBM Corp.

Contents

About this document.....	ix
Highlighting.....	ix
Case-sensitivity in AIX.....	ix
ISO 9000.....	ix
General programming concepts.....	1
What's new.....	1
Tools and utilities.....	1
Curses library.....	3
Initializing curses.....	5
Windows in the curses environment.....	6
Manipulating window data with curses.....	7
Controlling the cursor with curses.....	9
Manipulating characters with curses.....	10
Understanding terminals with curses	17
Working with color.....	23
Manipulating video attributes.....	23
Manipulating soft labels.....	26
Curses compatibility.....	26
List of additional curses subroutines.....	27
Debugging programs.....	28
adb debug program overview.....	29
Getting started with the adb debug program.....	29
Controlling program execution.....	30
Using adb expressions.....	34
Customizing the adb debug program.....	35
Computing numbers and displaying text.....	39
Displaying and manipulating the source file with the adb program.....	39
adb debug program reference information.....	46
Example adb program: adbsamp.....	52
Example adb program: adbsamp2.....	53
Example adb program: adbsamp3.....	53
Example of directory and i-node dumps in adb debugging.....	54
Example of data formatting in adb debugging.....	56
Example of tracing multiple functions in adb debugging.....	58
dbx symbolic debug program overview.....	60
Using the dbx debug program.....	60
Displaying and manipulating the source file with the dbx debug program.....	64
Examining program data.....	66
Debugging at the machine level with dbx.....	72
Customizing the dbx debugging environment.....	75
Developing for the dbx plug-in framework.....	77
List of dbx subcommands.....	94
Error-logging overview.....	96
Error-logging facility.....	97
Managing error logging.....	98
Error notification.....	101
Error logging tasks.....	104
Error logging and alerts.....	111
Error logging controls	112

File systems and logical volumes.....	113
File types.....	114
Working with JFS directories.....	116
Working with JFS2 directories.....	118
Working with JFS i-nodes.....	120
Working with JFS2 i-nodes	121
Allocating JFS file space.....	123
Allocating JFS2 file space.....	126
JFS file system layout.....	127
JFS2 file system layout.....	129
Writing programs that access large files.....	130
Linking for programmers.....	134
Using file descriptors.....	136
Creating and removing files.....	139
Working with file I/O.....	140
File status.....	147
File accessibility.....	148
Creating new file system types.....	149
Logical volume programming.....	152
J2_CFG_ASSIST ioctl operation.....	153
Floating-point exceptions.....	154
Input and output handling.....	163
Storage protect keys.....	167
Large program support.....	172
Programming on multiprocessor systems.....	176
Identifying processors.....	176
Controlling processor use.....	177
Using Dynamic Processor Deallocation.....	177
Dynamic memory guarding.....	181
Creating locking services.....	181
ProbeVue dynamic tracing facility.....	184
ProbeVue concepts.....	184
ProbeVue variables.....	190
Running ProbeVue.....	266
Vue functions.....	374
Process attributes.....	404
Multithreaded programming.....	407
Understanding threads and processes.....	408
Threadsafe and threaded libraries in AIX	412
Creating threads.....	413
Terminating threads.....	416
Synchronization overview.....	423
Using mutexes.....	424
Using condition variables.....	429
Using read/write locks.....	434
Joining threads.....	441
Scheduling threads.....	444
Contention scope and concurrency level.....	447
Synchronization scheduling.....	448
One-time initializations.....	451
Thread-specific data.....	452
Creating complex synchronization objects.....	456
Signal management.....	459
Process duplication and termination.....	463
Threads library options.....	464
Writing reentrant and threadsafe code.....	473
Developing multithreaded programs.....	478
Developing multithreaded programs to examine and modify pthread library objects.....	482

Developing multithreaded program debuggers.....	485
Benefits of threads.....	489
lex and yacc program information.....	490
Generating a lexical analyzer with the lex command.....	491
Using the lex program with the yacc program.....	492
Extended regular expressions in the lex command.....	493
Passing code to the generated lex program.....	496
Defining lex substitution strings.....	496
lex library.....	497
Actions taken by the lexical analyzer.....	497
lex program start conditions.....	501
Creating a parser with the yacc program.....	502
The yacc grammar file.....	502
Using the yacc grammar file.....	504
yacc grammar file declarations.....	505
yacc rules.....	507
yacc actions.....	509
yacc program error handling.....	510
Parser operation generated by the yacc command.....	511
Using ambiguous rules in the yacc program.....	513
Turning on debug mode for a parser generated by the yacc command.....	515
Example program for the lex and yacc programs.....	515
make command.....	519
Creating a description file.....	519
Internal rules for the make program.....	522
Defining and using macros in a description file.....	525
Creating a target file with the make command.....	529
Using the make command with source code control system files.....	530
Using the make command with non-source code control system (SCCS) files.....	531
Understanding how the make command uses environment variables.....	532
Using the make command in parallel run mode.....	532
m4 macro processor overview.....	532
Object data manager.....	541
ODM commands and subroutines.....	551
ODM example code and output.....	553
Simultaneous multithreading.....	556
Dynamic logical partitioning.....	558
DLPAR-safe and aware programs.....	558
Processor bindings.....	562
Integrating the DLPAR operation into the application.....	563
Actions taken by DLPAR scripts.....	563
Making kernel extensions DLPAR-aware.....	570
sed program information	574
Manipulating strings with sed.....	574
Using text in commands.....	579
Using string replacement.....	579
Shared libraries and shared memory.....	580
Shared objects and run time linking.....	581
Shared libraries and lazy loading.....	583
Named shared library areas.....	584
Creating a shared library.....	586
Program address space overview.....	588
Understanding memory mapping.....	590
Interprocess communication limits.....	594
Creating a mapped data file with the shmat subroutine.....	597
Creating a copy-on-write mapped data file with the shmat subroutine.....	598
Creating a shared memory segment with the shmat subroutine.....	599
Paging space programming requirements.....	600

List of memory manipulation services.....	601
List of memory mapping services.....	602
AIX vector programming.....	602
System memory allocation using the malloc subsystem.....	609
User-defined malloc replacement.....	619
Debug malloc tool.....	623
Malloc multiheap.....	629
Malloc buckets.....	630
Malloc trace.....	633
Malloc log.....	634
Malloc disclaim.....	636
Malloc detect.....	636
Configuring and using the malloc thread cache.....	637
Writing reentrant and threadsafe code.....	637
Packaging software for installation.....	642
Source code control system	681
SCCS flag and parameter conventions.....	682
Creating, editing, and updating an SCCS file.....	683
Controlling and tracking SCCS file changes.....	684
Detecting and repairing damaged SCCS files.....	685
List of additional SCCS commands.....	686
Subroutines, example programs, and libraries.....	686
128-bit long double floating-point data type.....	688
List of character manipulation subroutines.....	689
List of executable program creation subroutines.....	691
List of files and directories subroutines.....	692
List of numerical manipulation subroutines.....	694
List of long long integer numerical manipulation subroutines.....	696
List of 128-bit long double numerical manipulation subroutines.....	696
List of processes subroutines.....	697
List of multithreaded programming subroutines.....	700
List of programmer's workbench library subroutines.....	701
List of security and auditing subroutines.....	702
List of string manipulation subroutines.....	705
Example: Program for manipulating characters.....	706
Example: Searching and sorting program.....	709
List of operating system libraries.....	713
System Management Interface Tool (SMIT).....	714
SMIT screen types.....	715
SMIT object classes.....	717
SMIT aliases and fast paths.....	720
SMIT information command descriptors.....	721
SMIT command generation and execution.....	723
Adding tasks to the SMIT database.....	725
Debugging SMIT database extensions.....	726
Creating SMIT help information for a new task.....	726
sm_menu_opt (SMIT menu) object class.....	727
sm_name_hdr (SMIT selector header) object class.....	729
sm_cmd_opt (SMIT dialog/selector command option) object class.....	732
sm_cmd_hdr (SMIT dialog header) object class.....	736
SMIT example program.....	739
System resource controller.....	749
SRC objects.....	750
SRC communication types.....	755
Programming subsystem communication with the SRC.....	757
Defining your subsystem to the SRC.....	763
List of additional SRC subroutines.....	764
Trace facility.....	764

Start the trace facility.....	769
Tracing user application.....	772
Tracing data structures.....	772
Trace stream attributes.....	776
Trace event type definitions.....	776
Trace subroutines.....	777
tty subsystem.....	780
Line discipline module (ldterm).....	784
Converter modules.....	787
TTY drivers.....	788
Loader domains.....	789
Data management application programming interface	791
AIX transactional memory programming.....	795
Notices.....	803
Privacy policy considerations.....	804
Trademarks.....	805
Index.....	807

About this document

This topic collection provides application developers with complete information about writing applications for the AIX® operating system. Programmers can use this topic collection to gain knowledge of programming guidelines and resources. Topics include input and output handling, curses, file systems and directories, lex and yacc, logical volume programming, shared libraries, large program support, packaging, trace facility, and System Management Interface Tool (SMIT).

Highlighting

The following highlighting conventions are used in this document:

Bold	Identifies commands, subroutines, keywords, files, structures, directories, and other items whose names are predefined by the system. Also identifies graphical objects such as buttons, labels, and icons that the user selects.
<i>Italics</i>	Identifies parameters whose actual names or values are to be supplied by the user.
Monospace	Identifies examples of specific data values, examples of text similar to what you might see displayed, examples of portions of program code similar to what you might write as a programmer, messages from the system, or information you should actually type.

Case-sensitivity in AIX

Everything in the AIX® operating system is case-sensitive, which means that it distinguishes between uppercase and lowercase letters. For example, you can use the **ls** command to list files. If you type LS, the system responds that the command is not found. Likewise, **FILEA**, **FiLea**, and **filea** are three distinct file names, even if they reside in the same directory. To avoid causing undesirable actions to be performed, always ensure that you use the correct case.

ISO 9000

ISO 9000 registered quality systems were used in the development and manufacturing of this product.

General programming concepts

This topic collection provides application developers with complete information about writing applications for the AIX® operating system. Programmers can use this information to gain knowledge of programming guidelines and resources. Topics include input and output handling, curses, file systems and directories, lex and yacc, logical volume programming, shared libraries, large program support, packaging, trace facility, and System Management Interface Tool (SMIT).



The AIX operating system is designed to support The Open Group's Single UNIX Specification Version 3 (UNIX 03) for portability of operating systems based on the UNIX operating system. Many new interfaces, and some current ones, have been added or enhanced to meet this specification. To determine the correct way to develop a UNIX 03 portable application, see The Open Group's UNIX 03 specification on The UNIX System website (<http://www.unix.org>).

What's new in General programming concepts

Read about new or significantly changed information for the General programming concepts topic collection.

How to see what's new or changed

To help you see where technical changes have been made, the General programming concepts topic collection uses:

- The  image to mark where new or changed information begins.
- The  image to mark where new or changed information ends.

November 2019

The following information is a summary of the updates that are made to the General programming topic collection:

- Added information about preferring the writer-thread over a reader thread to get the pthread read/write locks in the [“Using read/write locks”](#) on page 434 topic.

November 2018

The following information is a summary of the updates that are made to the General programming topic collection:

- Added process ID information to the [“UFT probe manager”](#) on page 285 topic.
- Added information about the `get_cwd` interface in the [“Process attributes”](#) on page 404 topic.

Tools and utilities

This section provides an overview of the tools and utilities that you can use to develop C compiled language programs.

Many tools are provided to help you develop C compiled programs. The tools provide help with the following programming tasks:

Subroutines and shell commands are provided for use in a C compiled program.

- [Entering a program into the system](#)
- [Checking a program](#)
- [Compiling and linking a program](#)

- [Subroutines](#)
- [Shell commands](#)

Entering a program into the system

The system has a line editor called **ed** for use in entering a program into a file. The system also has the full-screen editor called **vi**, which displays one full screen of data at a time and allows interactive editing of a file.

Checking a program

Use the following commands to check the format of a program for consistency and accuracy:

Commands	Description
cb	Reformats a C language source program into a consistent format that uses indentation levels to show the structure of the program.
cflow	Generates a diagram of the logic flow of a C language source program.
cxref	Generates a list of all external references for each module of a C language source program, including where the reference is resolved (if it is resolved in the program).
lint	Checks for syntax and data type errors in a C language source program. The lint command might check these areas of a program more carefully than the C language compiler does, and displays many messages that point out possible problems.

Compiling and linking a program

To convert source code into a program that the system can run, you must process the source file with a compiler and a linkage editor.

A *compiler* is a program that reads text from a file and changes the programming language in that file to a form that the system understands. The *linkage editor* connects program modules and determines how to put the finished program into memory. To create this final form of the program, the system does the following:

1. If a file contains compiler source code, the compiler translates it into object code.
2. If a file contains assembler language, the assembler translates it into object code.
3. The linkage editor links the object files created in the previous step with any other object files specified in the compiler command.

Other programming languages available for use on the operating system include the C++, FORTRAN, COBOL, and Assembler and other compiler languages.

You can write parts of a program in different languages and have one main routine call and start the separate routines to execute. You can also use the compiler program to create object code and link the program.

Correcting errors in a program

You can use the following debugging tools that are provided with the base operating system:

- The **dbx** symbolic debug program allows you to debug programs written in C language, C++, FORTRAN, COBOL and Assembler languages.
- The **adb** debug program provides subcommands you can use to examine, debug, and repair executable binary files and to examine non-ASCII data files.
- [KDB Kernel Debugger and kdb command](#) can help you determine errors in code running in the kernel. You can use this debug program to debug device drivers and kernel extensions.
- The trace facility helps isolate system problems by monitoring selected system events.

When syntax errors or parameter-naming inconsistencies are discovered in a program file, you can use a text editor or string-searching and string-editing programs to locate and change strings in the file. String-searching and string-editing programs include the **grep**, **sed**, and **awk** commands. To make many changes in one or more program files, you can include the commands in a shell program and then run the shell program to locate and change the code in the files.

Building and Maintaining a Program

The following facilities help you control program changes and build a program from many source modules. These facilities can be particularly useful in software development environments in which many source modules are produced.

- The **make** command builds a program from source modules. Because the **make** command compiles only those modules changed since the last build, its use can reduce compilation time when many source modules must be processed.
- The Source Code Control System (SCCS) allows you to maintain separate versions of a program without storing separate, complete copies of each version. The use of SCCS can reduce storage requirements and help in tracking the development of a project that requires keeping many versions of large programs.

Subroutines

Subroutines from system libraries handle many complex or repetitive programming situations so that you can concentrate on unique programming situations.

Shell commands

You can include the functions of many of the shell commands in a C language program. Any shell command used in a program must be available on all systems that use the program.

You can then use the **fork** and **exec** subroutines in a program to run the command as a process in a part of the system that is separate from the program. The **system** subroutine also runs a shell command in a program, and the **popen** subroutine uses shell filters.

Related concepts

[Manipulating strings with sed](#)

The **sed** program performs its editing without interacting with the person requesting the editing.

[Generating a lexical analyzer with the lex command](#)

The **lex** command helps write a C language program that can receive and translate character-stream input into program actions.

[make command](#)

This topic provides information about simplifying the recompiling and relinking processes using the **make** command.

[Subroutines, example programs, and libraries](#)

This topic provides information about what subroutines are, how to use them, and where they are stored.

Curses library

The curses library provides a set of functions that enable you to manipulate a terminal's display regardless of the terminal type. The curses library supports color. However, multibyte characters are not supported. All references to characters in the curses documentation refer to single-byte characters.

Throughout this documentation, the curses library is referred to as *curses*.

The basis of curses programming is the window data structure. Using this structure, you can manipulate data on a terminal's display. You can instruct curses to treat the entire terminal display as one large window, or you can create multiple windows on the display. The windows can be different sizes and can overlap one another. A typical curses application has a single large window and one subwindow within it.

Each window on a terminal's display has its own window data structure. This structure keeps state information about the window, such as its size and where it is located on the display. Curses uses the window data structure to obtain the relevant information it needs to carry out your instructions.

Terminology

When programming with curses, you should be familiar with the following terms:

Term	Definition
current character	The character that the logical cursor is currently on.
current line	The line that the logical cursor is currently on.
curscr	A virtual default window provided by curses. The curscr (current screen) is an internal representation of what currently appears on the terminal's external display. Do not modify the curscr.
display	A physical display connected to a workstation.
logical cursor	The cursor location within each window. The window data structure keeps track of the location of its logical cursor.
pad	A pad is a window that is not restricted by the size of the screen
physical cursor	The cursor that appears on a display. The workstation uses this cursor to write to the display. There is only one physical cursor per display.
screen	The window that fills the entire display. The screen is synonymous with the stdscr.
stdscr	A virtual default window (standard screen) provided by curses that represents the entire display.
window	A pointer to a C data structure and the graphic representation of that data structure on the display. A window can be thought of as a two-dimensional array representing how all or part of the display looks at any point in time.

Naming conventions

A single curses subroutine can have more than one version. Curses subroutines with multiple versions follow distinct naming conventions that identify the separate versions. These conventions add a prefix to a standard curses subroutine and identify what arguments the subroutine requires or what actions take place when the subroutine is called. The different versions of curses subroutine names use the following prefixes:

Prefix	Description
w	Identifies a subroutine that requires a window argument
p	Identifies a subroutine that requires a pad argument
mv	Identifies a subroutine that first performs a move to the program-supplied coordinates

If a curses subroutine has multiple versions and does not include one of the preceding prefixes, the curses default window stdscr (standard screen) is used. The majority of subroutines that use the stdscr are macros created in the `/usr/include/curses.h` file using `#define` statements. The preprocessor replaces these statements at compilation time. As a result, these macros do not appear in the compiled assembler code, a trace, a debug program, or the curses source code.

If a curses subroutine has only a single version, it does not necessarily use stdscr. For example, the `printw` subroutine prints a string to the stdscr. The `wprintw` subroutine prints a string to a specific window by supplying the `window` argument. The `mvprintw` subroutine moves the specified coordinates to the stdscr and then performs the same function as the `printw` subroutine. Likewise, the `mvwprintw` subroutine moves the specified coordinates to the specified window and then performs the same function as the `wprintw` subroutine.

Structure of a curses Pprogram

In general, a curses program has the following progression:

1. Start curses.
2. Check for color support (optional).
3. Start color (optional).

4. Create one or more windows.
5. Manipulate windows.
6. Destroy one or more windows.
7. Stop curses.

Some steps are optional, so your program does not have to follow this progression exactly.

Return values

With a few exceptions, all curses subroutines return either the integer value ERR or the integer value OK. Subroutines that do not follow this convention are noted appropriately. Subroutines that return pointers always return a null pointer or an error.

Initializing curses

This section describes the commands used for initializing curses.

Use the following commands to initialize curses:

Command	Description
<u>endwin</u>	Terminates the curses subroutine libraries and their data structures
<u>initscr</u>	Initializes the curses subroutine library and its data structures
<u>isendwin</u>	Returns TRUE if the endwin subroutine has been called without any subsequent calls to the wrefresh subroutine
<u>newterm</u>	Sets up a new terminal
<u>setupterm</u>	Sets up the TERMINAL structure for use by curses

You must include the **curses.h** file at the beginning of any program that calls curses subroutines. To do this, use the following statement:

```
#include <curses.h>
```

Before you can call subroutines that manipulate windows or screens, you must call the **initscr** or **newterm** subroutine. These subroutines first save the terminal's settings and then call the **setupterm** subroutine to establish a curses terminal.

If you need to temporarily suspend curses, use a shell escape or subroutine. To resume after a temporary escape, call the **wrefresh** or **doupdate** subroutine. Before exiting a curses program, you must call the **endwin** subroutine. The **endwin** subroutine restores tty modes, moves the cursor to the lower-left corner of the screen, and resets the terminal into the proper nonvisual mode.

Most interactive, screen-oriented programs require character-at-a-time input without echoing the result to the screen. To establish your program with character-at-a-time input, call the **cbreak** and **noecho** subroutines after calling the **initscr** subroutine. When accepting this type of input, programs should also call the following subroutines:

- **nonl** subroutine.
- **intrflush** subroutine with the *Window* parameter set to the **stdscr** and the *Flag* parameter set to **FALSE**. The *Window* parameter is required but ignored. You can use **stdscr** as the value of the *Window* parameter, because **stdscr** is already created for you.
- **keypad** subroutine with the *Window* parameter set to the **stdscr** and the *Flag* parameter set to **TRUE**.

The **isendwin** subroutine is helpful if, for optimization reasons, you do not want to call the **wrefresh** subroutine needlessly. To determine if the **endwin** subroutine was called without any subsequent calls to the **wrefresh** subroutine, use the **isendwin** subroutine.

Windows in the curses environment

A curses program manipulates windows that appear on a terminal's display. A window can be as large as the entire display or as small as a single character in length and height.

Note: A pad is a window that is not restricted by the size of the screen.

Within a curses program, windows are variables declared as type `WINDOW`. The `WINDOW` data type is defined in the `/usr/include/curses.h` file as a C data structure. You create a window by allocating a portion of a machine's memory for a window structure. This structure describes the characteristics of the window. When a program changes the window data internally in memory, it must use the **wrefresh** subroutine (or equivalent subroutine) to update the external, physical screen to reflect the internal change in the appropriate window structure.

Default window structure

Curses provides a virtual default window structure called *stdscr*. The *stdscr* represents, in memory, the entire terminal display. The *stdscr* window structure is created automatically when the curses library is initialized and it describes the display. When the library is initialized, the *length* and *width* variables are set to the length and width of the physical display.

Programs that use the *stdscr* first manipulate the *stdscr*. They then call the **refresh** subroutine to refresh the external display so that it matches the *stdscr* window.

In addition to the *stdscr*, you can define your own windows. These windows are known as *user-defined windows* to distinguish them from the *stdscr*. Like the *stdscr*, user-defined windows exist in machine memory as structures. Except for the amount of memory available to a program, there is no limit to the number of windows you can create. A curses program can manipulate the default window, user-defined windows, or both.

Current window structure

Curses supports another virtual window called *curscr* (current screen). The *curscr* window is an internal representation of what currently appears on the terminal's external display.

When a program requires the external representation to match the internal representation, it must call a subroutine, such as the **wrefresh** subroutine, to update the physical display (or the **refresh** subroutine if the program is working with the *stdscr*).

The *curscr* is reserved for internal use by curses. Do not manipulate the *curscr*.

Subwindows

Curses also allows you to construct *subwindows*. Subwindows are rectangular portions within other windows. A subwindow is also of type `WINDOW`. The window that contains a subwindow is known as the subwindow's parent, and the subwindow is known as the containing window's child.

Changes to either the parent window or the child window within the area overlapped by the subwindow are made to both windows. After modifying a subwindow, call the **touchline** or **touchwin** subroutine on the parent window before refreshing it.

Subroutine	Description
<u>touchline</u>	Forces a range of lines to be refreshed at the next call to the wrefresh subroutine.
<u>touchwin</u>	Forces every character in a window's character array to be refreshed at the next call of the wrefresh subroutine. The touchwin subroutine does not save optimization information. This subroutine is useful with overlapping windows.

A refresh called on the parent also refreshes the children. A subwindow can also be a parent window. The process of layering windows inside of windows is called *nesting*.

Before you can delete a parent window, you must first delete all of its children using the **delwin** subroutine. Curses returns an error if you try to delete a window before first deleting all of its children.

Pads

A pad is a type of window that is not restricted by the terminal's display size or associated with a particular part of the display. Because a pad is usually larger than the physical display, only a portion of a pad is visible to the user at a given time.

Use pads if you have a large amount of related data that you want to keep all together in one window but you do not need to display all of the data at one time.

Windows within pads are known as *subpads*. Subpads are positioned within a pad at coordinates relative to the parent pad. This placement differs from subwindows, which are positioned using screen coordinates.

Unlike other windows, scrolling or echoing of input does not automatically refresh a pad. Like subwindows, when changing the image of a subpad, you must call either the **touchline** or **touchwin** subroutine on the parent pad before refreshing the parent.

You can use all the curses subroutines with pads except for the **newwin**, **subwin**, **wrefresh**, and **wnoutrefresh** subroutines. These subroutines are replaced with the **newpad**, **subpad**, **prefresh**, and **pnoutrefresh** subroutines.

Manipulating window data with curses

When curses is initialized, the `stdscr` is provided automatically. You can manipulate the `stdscr` using the curses subroutine library or you can create user-defined windows.

Creating windows

You can create your own window using the **newwin** subroutine.

Each time you call the **newwin** subroutine, curses allocates a new window structure in memory. This structure contains all the information associated with the new window. Curses does not put a limit on the number of windows you can create. The number of nested subwindows is limited to the amount of memory available, up to the value of `SHRT_MAX` as defined in the `/usr/include/limits.h` file.

You can change windows without regard to the order in which they were created. Updates to the terminal's display occur through calls to the **wrefresh** subroutine.

Subwindows

You must supply coordinates for the subwindow relative to the terminal's display. The subwindow, created using the **subwin** subroutine, must fit within the bounds of the parent window. Otherwise, a null value is returned.

Pads

Use the following subroutines to create pads:

Subroutine	Description
------------	-------------

newpad	Creates a pad data structure.
---------------	-------------------------------

subpad	Creates and returns a pointer to a subpad within a pad.
---------------	---

Removing windows, pads, and subwindows

The new subpad is positioned relative to its parent.

To remove a window, pad, or subwindow, use the **delwin** subroutine. Before you can delete a window or pad, you must have already deleted its children; otherwise, the **delwin** subroutine returns an error.

Changing the screen or window images

When curses subroutines change the appearance of a window, the internal representation of the window is updated, while the display remains unchanged until the next call to the **wrefresh** subroutine. The **wrefresh** subroutine uses the information in the window structure to update the display.

Refreshing windows

Whenever you write output to a window or pad structure, you must refresh the terminal's display to match the internal representation. A refresh does the following:

- Compares the contents of the `curscr` to the contents of the user-defined or `stdscr`
- Updates the `curscr` structure to match the user-defined or `stdscr`
- Redraws the portion of the physical display that changed

Use the following subroutines to refresh windows:

Subroutine	Description
<u>refresh</u> , or <u>wrefresh</u>	Updates the terminal and <code>curscr</code> to reflect changes made to a window.
<u>wnoutrefresh</u> or <u>doupdate</u>	Updates the designated windows and outputs them all at once to the terminal. These subroutines are useful for faster response when there are multiple updates.

The **refresh** and **wrefresh** subroutines first call the **wnoutrefresh** subroutine to copy the window being refreshed to the current screen. They then call the **doupdate** subroutine to update the display.

If you need to refresh multiple windows at the same time, use one of the two available methods. You can use a series of calls to the **wrefresh** subroutine that result in alternating calls to the **wnoutrefresh** and **doupdate** subroutines. You can also call the **wnoutrefresh** subroutine once for each window and then call the **doupdate** subroutine once. With the second method, only one burst of output is sent to the display.

Subroutines used for refreshing pads

The **prefresh** and **pnoutrefresh** subroutines are similar to the **wrefresh** and **wnoutrefresh** subroutines.

The **prefresh** subroutine updates both the current screen and the physical display, while the **pnoutrefresh** subroutine updates `curscr` to reflect changes made to a user-defined pad. Because pads instead of windows are involved, these subroutines require additional parameters to indicate which part of the pad and screen are involved.

Refreshing areas that have not changed

During a refresh, only those areas that have changed are redrawn on the display. You can refresh areas of the display that have not changed using the **touchwin** and **touchline** subroutines:

Subroutine	Description
<u>touchline</u>	Forces a range of lines to be refreshed at the next call to the wrefresh subroutine.
<u>touchwin</u>	Forces every character in a window's character array to be refreshed at the next call of the wrefresh subroutine. The touchwin subroutine does not save optimization information. This subroutine is useful with overlapping windows.

Combining the **touchwin** and **wrefresh** subroutines is helpful when dealing with subwindows or overlapping windows. To bring a window forward from behind another window, call the **touchwin** subroutine followed by the **wrefresh** subroutine.

Garbled displays

If text is sent to the terminal's display with a noncurses subroutine, such as the **echo** or **printf** subroutine, the external window can become garbled. In this case, the display changes, but the current screen is not updated to reflect these changes. Problems can arise when a refresh is called on the garbled screen because after a screen is garbled, there is no difference between the window being refreshed and the current screen structure. As a result, spaces on the display caused by garbled text are not changed.

A similar problem can also occur when a window is moved. The characters sent to the display with the noncurses subroutines do not move with the window internally.

If the screen becomes garbled, call the **wrefresh** subroutine on the `curscr` to update the display to reflect the current physical display.

Manipulating window content

After a window or subwindow is created, programs often must manipulate them in some way, by using the following subroutines:

Subroutine	Description
<u>box</u>	Draws a box in or around a window
<u>copywin</u>	Provides more precise control over the overlay and overwrite subroutine
<u>garbagedlines</u>	Indicates to curses that a screen line is discarded and should be thrown away before having anything written over
<u>mvwin</u>	Moves a window or subwindow to a new location
<u>overlay</u> or <u>overwrite</u>	Copies one window on top of another
<u>ripoffline</u>	Removes a line from the default screen

To use the **overlay** and **overwrite** subroutines, the two windows must overlap. The **overwrite** subroutine is destructive, whereas the **overlay** subroutine is not. When text is copied from one window to another using the **overwrite** subroutine, blank portions from the copied window overwrite any portions of the window copied to. The **overlay** subroutine is nondestructive because it does not copy blank portions from the copied window.

Similar to the **overlay** and **overwrite** subroutines, the **copywin** subroutine allows you to copy a portion of one window to another. Unlike **overlay** and **overwrite** subroutines, the windows do not have to overlap for you to use the **copywin** subroutine.

To remove a line from the stdscr, you can use the **ripoffline** subroutine. If you pass this subroutine a positive *line* argument, the specified number of lines is removed from the top of the stdscr. If you pass the subroutine a negative *line* argument, the lines are removed from the bottom of the stdscr.

To discard a specified range of lines before writing anything new, you can use the **garbagedlines** subroutine.

Support for filters

The [filter](#) subroutine is provided for curses applications that are filters. This subroutine causes curses to operate as if the stdscr was only a single line. When running with the **filter** subroutine, curses does not use any terminal capabilities that require knowledge of the line that curses is on.

Controlling the cursor with curses

This section explains the different types of cursors that exist in the curses library.

The following types of cursors exist in the curses library:

logical cursor

The cursor location within each window. A window's data structure keeps track of the location of its logical cursor. Each window has a logical cursor.

physical cursor

The display cursor. The workstation uses this cursor to write to the display. There is only one physical cursor per display.

You can only add to or erase characters at the logical cursor in a window. The following subroutines are provided for controlling the cursor:

[getbegyx](#)

Places the beginning coordinates of the window in integer variables *y* and *x*

[getmaxyx](#)

Places the size of the window in integer variables *y* and *x*

[getsyx](#)

Returns the current coordinates of the virtual screen cursor

getyx

Returns the position of the logical cursor associated with a specified window

leaveok

Controls physical cursor placement after a call to the **wrefresh** subroutine

move

Moves the logical cursor associated with the stdscr

mvcur

Moves the physical cursor

setsyx

Sets the virtual screen cursor to the specified coordinate

wmove

Moves the logical cursor associated with a user-defined window

After a call to the **refresh** or **wrefresh** subroutine, curses places the physical cursor at the last updated character position in the window. To leave the physical cursor where it is and not move it after a refresh, call the **leaveok** subroutine with the *Window* parameter set to the desired window and the *Flag* parameter set to **TRUE**.

Manipulating characters with curses

You can add characters to a curses window using a keyboard or a curses application. This section describes how you can add, remove, or change characters that appear in a curses window.

Character size

Some character sets define multi-column characters that occupy more than one column position when displayed on the screen.

Writing a character whose width is greater than the width of the destination window produces an error.

Adding characters to the screen image

The curses library provides a number of subroutines that write text changes to a window and mark the area to be updated at the next call to the **wrefresh** subroutine.

waddch subroutines

The **waddch** subroutines overwrite the character at the current logical cursor location with a specified character. After overwriting, the logical cursor is moved one space to the right. If the **waddch** subroutines are called at the right margin, these subroutines also add an automatic newline character. Additionally, if you call one of these subroutines at the bottom of a scrolling region and the **scrollok** subroutine is enabled, the region is scrolled up one line. For example, if you added a new line at the bottom line of a window, the window would scroll up one line.

If the character to add is a tab, newline, or backspace character, curses moves the cursor appropriately in the window to reflect the addition. Tabs are set at every eighth column. If the character is a newline, curses first uses the **wclrtoeol** subroutine to erase the current line from the logical cursor position to the end of the line before moving the cursor. The **waddch** subroutine family is made up of the following:

Subroutine	Description
addch macro	Adds a character to the stdscr
mvaddch macro	Moves a character to the specified location before adding it to the stdscr
mvwaddch macro	Moves a character to the specified location before adding it to the user-defined window
waddch subroutine	Adds a character to the user-defined window

By using the **winch** and **waddch** subroutine families together, you can copy text and video attributes from one place to another. Using the **winch** subroutine family, you can retrieve a character and its video

attributes. You can then use one of the **waddch** subroutines to add the character and its attributes to another location.

You can also use the **waddch** subroutines to add control characters to a window. Control characters are drawn in the ^X notation.

Note: Calling the **winch** subroutine on a position in the window containing a control character does not return the character. Instead, it returns one character of the control character representation.

Outputting single, noncontrol characters

When outputting single, noncontrol characters, there can be significant performance gain to using the **wechochar** subroutines. These subroutines are functionally equivalent to a call to the corresponding **waddchr** subroutine followed by the corresponding **wrefresh** subroutine. The **wechochar** subroutines include the **wechochar** subroutine, the **echochar** macro, and the **pechochar** subroutine.

Some character sets may contain nonspacing characters. (Nonspacing characters are those, other than the '\0' character, for which the **wcwidth** subroutine returns a width of zero.) The application may write nonspacing characters to a window. Every nonspacing character in a window is associated with a spacing character and modifies the spacing character. Nonspacing characters in a window cannot be addressed separately. A nonspacing character is implicitly addressed whenever a Curses operation affects the spacing character with which the nonspacing character is associated.

Nonspacing characters do not support attributes. For interfaces that use wide characters and attributes, the attributes are ignored if the wide character is a nonspacing character. Multi-column characters have a single set of attributes for all columns. The association of nonspacing characters with spacing characters can be controlled by the application using the wide character interfaces. The wide character string functions provide codeset-dependent association.

The typical effects of a nonspacing character associated with a spacing character called *c*, are as follows:

- The nonspacing character may modify the appearance of *c*. (For instance, there may be nonspacing characters that add diacritical marks to characters. However, there may also be spacing characters with built-in diacritical marks.)
- The nonspacing characters may bridge *c* to the character following *c*. Examples of this usage are the formation of ligatures and the conversion of characters into compound display forms, words, or ideograms.

Implementations may limit the number of nonspacing characters that can be associated with a spacing character, provided any limit is at least 5.

Complex characters

A complex character is a set of associated characters, which may include a spacing character and may also include any nonspacing characters associated with it. A spacing complex character is a complex character that includes one spacing character and any nonspacing characters associated with it. An example of a code set that has complex characters is ISO/IEC 10646-1:1993.

A complex character can be written to the screen. If the complex character does not include a spacing character, any nonspacing characters are associated with the spacing complex character that exists at the specified screen position. When the application reads information back from the screen, it obtains spacing complex characters.

The `cchar_t` data type represents a complex character and its rendition. When a `cchar_t` represents a nonspacing complex character (that is, when there is no spacing character within the complex character), then its rendition is not used. When it is written to the screen, it uses the rendition specified by the spacing character already displayed.

An object of type `cchar_t` can be initialized using the **setchar** subroutine, and its contents can be extracted using the **getchar** subroutine. The behavior of functions that take a `cchar_t` value that was not initialized in this way are obtained from a curses function that has a `cchar_t` output argument.

Special characters

Some functions process special characters. In functions that do not move the cursor based on the information placed in the window, these special characters would only be used within a string in order to affect the placement of subsequent characters. The cursor movement specified below does not persist in the visible cursor beyond the end of the operation. In functions that do not move the cursor, these special characters can be used to affect the placement of subsequent characters and to achieve movement of the physical cursor.

Character	Description
Backspace	Unless the cursor was already in column 0, Backspace moves the cursor one column toward the start of the current line, and any characters after the Backspace are added or inserted starting there.
Carriage return	Unless the cursor was already in column 0, Carriage return moves the cursor to the start of the current line. Any characters after the Carriage return are added or inserted starting there.
newline	In an add operation, curses adds the background character into successive columns until reaching the end of the line. Scrolling occurs, and any characters after the newline character are added, starting at the beginning of the new line. In an insert operation, newline erases the remainder of the current line with the background character (effectively a wclrtoeol subroutine), and moves the cursor to the start of a new line. When scrolling is enabled, advancing the cursor to a new line may cause scrolling. Any characters after the newline character are inserted at the beginning of the new line. The filter function may inhibit this processing.
Tab	Tab characters in text move subsequent characters to the next horizontal tab stop. By default, tab stops are in columns 0, 8, 16, and so on. In an insert or add operation, curses inserts or adds, respectively, the background character into successive columns until reaching the next tab stop. If there are no more tab stops in the current line, wrapping and scrolling occur.

Control characters

The curses functions that perform special-character processing conceptually convert control characters to the (' ^ ') character followed by a second character (which is an uppercase letter if it is alphabetic) and write this string to the window in place of the control character. The functions that retrieve text from the window will not retrieve the original control character.

Line graphics:

You can use the following variables to add line-drawing characters to the screen with the **waddch** subroutine. When defined for the terminal, the variable will have the **A_ALTCHARSET** bit turned on. Otherwise, the default character listed in the following table is stored in the variable.

Variable name	Default character	Glyph description
ACS_ULCORNER	+	upper left corner
ACS_LLCORNER	+	lower left corner
ACS_URCORNER	+	upper right corner
ACS_LRCORNER	+	lower right corner
ACS_RTEE	+	right tee
ACS_LTEE	+	left tee
ACS_BTEE	+	bottom tee

Variable name	Default character	Glyph description
ACS_TTEE	+	top tee
ACS_HLINE	—	horizontal line
ACS_VLINE		vertical line
ACS_PLUS	+	plus
ACS_S1	-	scan line 1
ACS_S9	_	scan line 9
ACS_DIAMOND	+	diamond
ACS_CKBOARD	:	checkerboard (stipple)
ACS_DEGREE	,	degree symbol
ACS_PLMINUS	#	plus/minus
ACS_BULLET	o	bullet
ACS_LARROW	<	arrow pointing left
ACS_RARROW	>	arrow pointing right
ACS_DARROW	v	arrow pointing down
ACS_UARROW	^	arrow pointing up
ACS_BOARD	#	board of squares
ACS_LANTERN	#	lantern symbol
ACS_BLOCK	#	solid square block

waddstr subroutines

The **waddstr** subroutines add a null-terminated character string to a window, starting with the current character. If you are adding a single character, use the **waddch** subroutine. Otherwise, use the **waddstr** subroutine. The following are part of the **waddstr** subroutine family:

Subroutine	Description
addstr macro	Adds a character string to the stdscr
mvaddstr macro	Moves the logical cursor to a specified location before adding a character string to the stdscr
waddstr subroutine	Adds a character string to a user-defined window
wmvaddstr macro	Moves the logical cursor to a specified location before adding a character string to a user-defined window

winsch subroutines

The **winsch** subroutines insert a specified character before the current character in a window. All characters to the right of the inserted character are moved one space to the right. As a result, the rightmost character on the line may be lost. The positions of the logical and physical cursors do not change after the move. The **winsch** subroutines include the following:

Subroutine	Description
insch macro	Inserts a character in the stdscr
mvinsch macro	Moves the logical cursor to a specified location in the stdscr before inserting a character

Subroutine	Description
mvwansch macro	Moves the logical cursor to a specified location in a user-defined window before inserting a character
wansch subroutine	Inserts a character in a user-defined window

winsertln subroutines

The **winsertln** subroutines insert a blank line above the current line in a window. The **insertln** subroutine inserts a line in the stdscr. The bottom line of the window is lost. The **winsertln** subroutine performs the same action in a user-defined window.

wprintw subroutines

The **wprintw** subroutines replace a series of characters (starting with the current character) with formatted output. The format is the same as for the **printf** command. The **printw** family is made up of the following:

Subroutine	Description
mvprintw macro	Moves the logical cursor to a specified location in the stdscr before replacing any characters
mvwprintw macro	Moves the logical cursor to a specified location in a user-defined window before replacing any characters
printw macro	Replaces a series of characters in the stdscr
wprintw subroutine	Replaces a series of characters in a user-defined window

The **wprintw** subroutines make calls to the **waddch** subroutine to replace characters.

unctrl macro

The **unctrl** macro returns a printable representation of the specified control character, displayed in the ^X notation. The **unctrl** macro returns print characters as is.

Enabling text scrolling

Use the following subroutines to enable scrolling:

Subroutine	Description
idlok	Allows curses to use the hardware insert/delete line feature
scrollok	Enables a window to scroll when the cursor is moved off the right edge of the last line of a window
setscreg or wsetscreg	Sets a software scrolling region within a window

Scrolling occurs when a program or user moves a cursor off a window's bottom edge. For scrolling to occur, you must first use the **scrollok** subroutine to enable scrolling for a window. A window is scrolled if scrolling is enabled and if any of the following occurs:

- The cursor is moved off the edge of a window.
- A newline character is encountered on the last line.
- A character is inserted in the last position of the last line.

When a window is scrolled, curses will update both the window and the display. However, to get the physical scrolling effect on the terminal, you must call the **idlok** subroutine with the *Flag* parameter set to **TRUE**.

If scrolling is disabled, the cursor remains on the bottom line at the location where the character was entered.

When scrolling is enabled for a window, you can use the **setscrreg** subroutines to create a software scrolling region inside the window. You pass the **setscrreg** subroutines values for the top line and bottom line of the region. If **setscrreg** is enabled for the region and scrolling is enabled for the window, any attempt to move off the specified bottom line causes all the lines in the region to scroll up one line. You can use the **setscrreg** macro to define a scrolling region in the stdscr. Otherwise, you use the **wsetscrreg** subroutine to define scrolling regions in user-defined windows.

Note: Unlike the **idlok** subroutine, the **setscrreg** subroutines have no bearing on the use of the physical scrolling region capability that the terminal may have.

Deleting characters

You can delete text by replacing it with blank spaces or by removing characters from a character array and sliding the rest of the characters on the line one space to the left.

werase subroutines

The **erases** macro copies blank space to every position in the stdscr. The **werase** subroutine puts a blank space at every position in a user-defined window. To delete a single character in a window, use the **wdelch** subroutine.

wclear subroutines

Use the following subroutines to clear the screen:

Subroutine	Description
clear , or wclear	Clears the screen and sets a clear flag for the next refresh.
clearok	Determines whether curses clears a window on the next call to the refresh or wrefresh subroutine.

The **wclear** subroutines are similar to the **werase** subroutines. However, in addition to putting a blank space at every position of a window, the **wclear** subroutines also call the **wclearok** subroutine. As a result, the screen is cleared on the next call to the **wrefresh** subroutine.

The **wclear** subroutine family contains the **wclear** subroutine, the **clear** macro, and the **clearok** subroutine. The **clear** macro puts a blank at every position in the stdscr.

wclrtoeol subroutines

The **clrtoeol** macro operates in the stdscr, while the **wclrtoeol** subroutine performs the same action within a user-defined window.

wclrtoeol subroutines

The **clrtoeol** macro operates in the stdscr, while the **wclrtoeol** performs the same action in a user-defined window.

wdelch subroutines

Use the following subroutines to delete characters from the screen:

Subroutine	Description
delch macro	Deletes the current character from the stdscr
mvdelch macro	Moves the logical cursor before deleting a character from the stdscr
mvwdelch macro	Moves the logical cursor before deleting a character from a user-defined window
wdelch subroutine	Deletes the current character in a user-defined window

The **wdelch** subroutines delete the current character and move all the characters to the right of the current character on the current line one position to the left. The last character in the line is filled with a blank. The **delch** subroutine family consists of the following subroutine and macros:

wdeleteln subroutines

The **deleteln** subroutines delete the current line and move all lines below the current line up one line. This action clears the window's bottom line.

Getting characters

Your program can retrieve characters from the keyboard or from the display. The **wgetch** subroutines retrieve characters from the keyboard. The **winch** subroutines retrieve characters from the display.

wgetch subroutines

The **wgetch** subroutines read characters from the keyboard attached to the terminal associated with the window. Before getting a character, these subroutines call the **wrefresh** subroutines if anything in the window has changed: for example, if the cursor has moved or text has changed.

The **wgetch** subroutine family is made up of the following:

Subroutine	Description
getch macro	Gets a character from the stdscr
mvgetch macro	Moves the cursor before getting a character from the stdscr
mvwgetch macro	Moves the cursor before getting a character from a user-defined window
wgetch subroutine	Gets a character from a user-defined window

To place a character previously obtained by a call to the **wgetch** subroutine back in the input queue, use the **ungetch** subroutine. The character is retrieved by the next call to the **wgetch** subroutine.

Terminal modes

The output of the **wgetch** subroutines is, in part, determined by the mode of the terminal. The following list describes the action of the **wgetch** subroutines in each type of terminal mode:

Subroutine	Description
DELAY mode	Stops reading until the system passes text through the program. If CBREAK mode is also set, the program stops after one character. If CBREAK mode is not set (NOCBREAK mode), the wgetch subroutine stops reading after the first newline character. If ECHO is set, the character is also echoed to the window.
HALF-DELAY mode	Stops reading until a character is typed or a specified timeout is reached. If ECHO mode is set, the character is also echoed to the window.
NODELAY mode	Returns a value of ERR if there is no input waiting.

Note: When you use the **getch**, **mvgetch**, **mvwgetch**, or **wgetch** subroutines, do not set both the NOCBREAK mode and the ECHO mode at the same time. Setting both modes can cause undesirable results depending on the state of the tty driver when each character is typed.

Function keys

Function keys are defined in the **curses.h** file. Function keys can be returned by the **wgetch** subroutine if the keypad is enabled. A terminal may not support all of the function keys. To see if a terminal supports a particular key, check its **terminfo** database definition.

Getting function keys

If your program enables the keyboard with the **keypad** subroutine, and the user presses a function key, the token for that function key is returned instead of raw characters. The **/usr/include/curses.h** file defines the possible function keys. Each define statement begins with a **KEY_** prefix, and the keys are defined as integers beginning with the value 03510.

If a character is received that could be the beginning of a function key (such as an Escape character), **curses** sets a timer (a structure of type **timeval** that is defined in **/usr/include/sys/time.h**). If the remainder of the sequence is not received before the timer expires, the character is passed through.

Otherwise, the function key's value is returned. For this reason, after a user presses the Esc key there is a delay before the escape is returned to the program. Avoid using the Esc key where possible when you call a single-character subroutine such as the **wgetch** subroutine. This timer can be overridden or extended by the use of the **ESCDELAY** environment variable.

The **ESCDELAY** environment variable sets the length of time to wait before timing out and treating the ESC keystroke as the Escape character rather than combining it with other characters in the buffer to create a key sequence. The **ESCDELAY** value is measured in fifths of a millisecond. If the **ESCDELAY** variable is 0, the system immediately composes the Escape response without waiting for more information from the buffer. You may choose any value from 0 to 99,999. The default setting for the **ESCDELAY** variable is 500 (1/10th of a second).

To prevent the **wgetch** subroutine from setting a timer, call the **notimeout** subroutine. If **notimeout** is set to TRUE, **curses** does not distinguish between function keys and characters when retrieving data.

keyname subroutine

The **keyname** subroutine returns a pointer to a character string containing a symbolic name for the *Key* argument. The *Key* argument can be any key returned from the **wgetch**, **getch**, **mvgetch**, or **mvwgetch** subroutines.

winch subroutines

The **winch** subroutines retrieve the character at the current position. If any attributes are set for the position, the attribute values are ORed into the value returned. You can use the **winch** subroutines to extract only the character or its attributes. To do this, use the predefined constants **A_CHARTEXT** and **A_ATTRIBUTES** with the logical & (ampersand) operator. These constants are defined in the **curses.h** file. The following are the **winch** subroutines:

Subroutine	Description
inch macro	Gets the current character from the stdscr
mvinch macro	Moves the logical cursor before calling the inch subroutine on the stdscr
mvwinch macro	Moves the logical cursor before calling the winch subroutine in the user-defined window
winch subroutine	Gets the current character from a user-defined window

wscanw subroutines

The **wscanw** subroutines read character data, interpret it according to a conversion specification, and store the converted results into memory. The **wscanw** subroutines use the **wgetstr** subroutines to read the character data. The following are the **wscanw** subroutines:

Subroutine	Description
mvscanw macro	Moves the logical cursor before scanning the stdscr
mvwscanw macro	Moves the logical cursor in the user-defined window before scanning
scanw macro	Scans the stdscr
wscanw subroutine	Scans a user-defined window

Understanding terminals with curses

The capabilities of your program are limited, in part, by the capabilities of the terminal on which it runs.

This section provides information about initializing terminals and identifying their capabilities.

Manipulating multiple terminals

With **curses**, you can use one or more terminals for input and output. The terminal subroutines enable you to establish new terminals, to switch input and output processing, and to retrieve terminal capabilities.

You can start curses on a single default screen using the **initscr** subroutine. If your application sends output to more than one terminal, use the **newterm** subroutine. Call the **newterm** subroutine for each terminal. Also use the **newterm** subroutine if your application wants an indication of error conditions so that it can continue to run in a line-oriented mode if the terminal cannot support a screen-oriented program.

When it completes, a program must call the **endwin** subroutine for each terminal it used. If you call the **newterm** subroutine more than once for the same terminal, the first terminal referred to must be the last one for which you call the **endwin** subroutine.

The **set_term** subroutine switches input and output processing between different terminals.

Determining terminal capabilities

Curses supplies the following subroutines to help you determine the capabilities of a terminal:

Subroutine	Description
has_ic	Determines whether a terminal has the insert-character capability
has_il	Determines whether a terminal has the insert-line capability
longname	Returns the verbose name of the terminal

The **longname** subroutine returns a pointer to a static area containing a verbose description of the current terminal. This static area is defined only after a call to the **initscr** or **newterm** subroutine. If you intend to use the **longname** subroutine with multiple terminals, each call to the **newterm** subroutine overwrites this area. Calls to the **set_term** subroutine do not restore the value. Instead, save this area between calls to the **newterm** subroutine.

The **has_ic** subroutine returns TRUE if the terminal has insert and delete character capabilities.

The **has_il** subroutine returns TRUE if the terminal has insert and delete line capabilities or can simulate the capabilities using scrolling regions. Use the **has_il** subroutine to check whether it is appropriate to turn on physical scrolling using the **scrollok** or **idlok** subroutines.

Setting terminal input and output modes

The subroutines that control input and output determine how your application retrieves and displays data to users.

Input Modes

Special input characters include the flow-control characters, the interrupt character, the erase character, and the kill character. The following mutually-exclusive curses modes let the application control the effect of the input characters:

Cooked mode

This mode achieves normal line-at-a-time processing with all special characters handled outside the application, achieving the same effect as canonical-mode input processing. The state of the ISIG and IXON flags is not changed upon entering this mode by calling `nocbreak()` and are set upon entering this mode by calling `noraw()`.

The implementation supports erase and kill characters from any supported locale, regardless of the width of the character.

cbreak mode

Characters typed by the user are immediately available to the application and curses does not perform special processing on either the erase character or the kill character. An application can select cbreak mode to do its own line editing but to let the abort character be used to abort the task. This mode achieves the same effect as noncanonical mode, Case B input processing (with MIN set to 1 and ICRNL cleared). The state of the ISIG and IXON flags is not changed upon entering this mode.

Half-delay mode

The effect is the same as **cbreak**, except that input functions wait until a character is available or an interval defined by the application elapses, whichever comes first. This mode achieves the same

effect as noncanonical mode, Case C input processing (with TIME set to the value specified by the application). The state of the ISIG and IXON flags is not changed upon entering this mode.

Raw mode

Raw mode gives maximum control to the application over terminal input. The application sees each character as it is typed. This achieves the same effect as noncanonical mode, Case D input processing. The ISIG and IXON flags are cleared upon entering this mode.

The terminal interface settings are recorded when the process calls the **initscr** or **newterm** subroutines to initialize curses and restores these settings when the **endwin** subroutine is called. The initial input mode for curses operations is unspecified unless the implementation supports enhanced curses compliance, in which the initial input mode is **cbreak** mode.

The behavior of the BREAK key depends on other bits in the display driver that are not set by curses.

Delay mode

The following mutually exclusive delay modes specify how quickly certain curses functions return to the application when there is no terminal input waiting when the function is called:

Delay	Description
No Delay	The function fails.
Delay	The application waits until the implementation passes text through to the application. If cbreak mode or Raw Mode is set, this is after one character. Otherwise, this is after the first newline character, end-of-line character, or end-of-file character.

The effect of No Delay mode on function-key processing is unspecified.

Echo mode processing

Echo mode determines whether curses echoes typed characters to the screen. The effect of echo mode is analogous to the effect of the echo flag in the local mode field of the **termios** structure associated with the terminal device connected to the window. However, curses always clears the echo flag when invoked, to inhibit the operating system from performing echoing. The method of echoing characters is not identical to the operating system's method of echoing characters, because curses performs additional processing of terminal input.

If in echo mode, curses performs its own echoing. Any visible input character is stored in the current or specified window by the input function that the application called, at that window's cursor position, as though the **addch** subroutine was called, with all consequent effects such as cursor movement and wrapping.

If not in echo mode, any echoing of input must be performed by the application. Applications often perform their own echoing in a controlled area of the screen, or do not echo at all, so they disable echo mode.

It may not be possible to turn off echo processing for synchronous and network asynchronous terminals because echo processing is done directly by the terminals. Applications running on such terminals should be aware that any characters typed will display on the screen at the point where the cursor is positioned.

The following are a part of the echo processing family of subroutines:

Subroutine	Description
<u>cbreak</u> or <u>nocbreak</u>	Puts the terminal into or takes it out of CBREAK mode
<u>delay_output</u>	Sets the output delay in milliseconds
<u>echo</u> or <u>noecho</u>	Controls echoing of typed characters to the screen
<u>halfdelay</u>	Returns ERR if no input was typed after blocking for a specified amount of time

Subroutine	Description
<u>nl</u> or <u>nonl</u>	Determines whether curses translates a new line into a carriage return and line feed on output, and translates a return into a new line on input
<u>raw</u> or <u>noraw</u>	Places the terminal into or out of mode

The **cbreak** subroutine performs a subset of the functions performed by the **raw** subroutine. In cbreak mode, characters typed by the user are immediately available to the program, and erase or kill character processing is not done. Unlike RAW mode, interrupt and flow characters are acted upon. Otherwise, the tty driver buffers the characters typed until a new line or carriage return is typed.

Note: CBREAK mode disables translation by the tty driver.

The **delay_output** subroutine sets the output delay to the specified number of milliseconds. Do not use this subroutine excessively because it uses padding characters instead of a processor pause.

The **echo** subroutine puts the terminal into echo mode. In echo mode, curses writes characters typed by the user to the terminal at the physical cursor position. The **noecho** subroutine takes the terminal out of echo mode.

The **nl** and **nonl** subroutines, respectively, control whether curses translates new lines into carriage returns and line feeds on output, and whether curses translates carriage returns into new lines on input. Initially, these translations do occur. By disabling these translations, the curses subroutine library has more control over the line-feed capability, resulting in faster cursor motion.

The **nocbreak** subroutine takes the terminal out of cbreak mode.

The **raw** subroutine puts the terminal into raw mode. In raw mode, characters typed by the user are immediately available to the program. Additionally, the interrupt, quit, suspend, and flow-control characters are passed uninterpreted instead of generating a signal as they do in cbreak mode. The **noraw** subroutine takes the terminal out of raw mode.

Using the terminfo and termcap files

When curses is initialized, it checks the **TERM** environment variable to identify the terminal type. Then, curses looks for a definition explaining the capabilities of the terminal. This information is usually kept in a local directory specified by the **TERMINFO** environment variable or in the **/usr/share/lib/terminfo** directory. All curses programs first check to see if the **TERMINFO** environment variable is defined. If this variable is not defined, the **/usr/share/lib/terminfo** directory is checked.

For example, if the **TERM** variable is set to **vt100** and the **TERMINFO** variable is set to the **/usr/mark/myterms** file, curses checks for the **/usr/mark/myterms/v/vt100** file. If this file does not exist, curses checks the **/usr/share/lib/terminfo/v/vt100** file.

Additionally, the **LINES** and **COLUMNS** environment variables can be set to override the terminal description.

Writing programs that use the terminfo subroutines

Use the **terminfo** subroutines when your program must deal directly with the **terminfo** database. For example, use these subroutines to program function keys. In all other cases, curses subroutines are more suitable and their use is recommended.

Initializing terminals

Your program should begin by calling the **setupterm** subroutine. Normally, this subroutine is called indirectly by a call to the **initscr** or **newterm** subroutine. The **setupterm** subroutine reads the terminal-dependent variables defined in the **terminfo** database. The **terminfo** database includes boolean, numeric, and string variables. All of these **terminfo** variables use the values defined for the specified terminal. After reading the database, the **setupterm** subroutine initializes the **cur_term** variable with the terminal definition. When working with multiple terminals, you can use the **set_curterm** subroutine to set the **cur_term** variable to a specific terminal.

Another subroutine, **restartterm**, is similar to the **setupterm** subroutine. However, it is called after memory is restored to a previous state. For example, you would call the **restartterm** subroutine after a call to the **scr_restore** subroutine. The **restartterm** subroutine assumes that the input and output options are the same as when memory was saved, but that the terminal type and baud rate may differ.

The **del_curterm** subroutine frees the space containing the capability information for a specified terminal.

Header files

Include the **curses.h** and **term.h** files in your program in the following order:

```
#include <curses.h>
#include <term.h>
```

These files contain the definitions for the strings, numbers, and flags in the **terminfo** database.

Handling terminal capabilities

Pass all parameterized strings through the **tparm** subroutine to instantiate them. Use the **tputs** or **putp** subroutine to print all **terminfo** strings and the output of the **tparm** subroutine.

Subroutine Description

putp	Provides a shortcut to the tputs subroutine
tparm	Instantiates a string with parameters
tputs	Applies padding information to the given string and outputs it

Use the following subroutines to obtain and pass terminal capabilities:

Subroutine	Description
tigetflag	Returns the value of a specified boolean capability. If the capability is not boolean, a -1 is returned.
tigetnum	Returns the value of a specified numeric capability. If the capability is not numeric, a -2 is returned.
tigetstr	Returns the value of a specified string capability. If the capability specified is not a string, the tigetstr subroutine returns the value of (char *) -1 .

Exiting the program

When your program exits, restore the tty modes to their original state. To do this, call the **reset_shell_mode** subroutine. If your program uses cursor addressing, it should output the **enter_ca_mode** string at startup and the **exit_ca_mode** string when it exits.

Programs that use shell escapes should call the **reset_shell_mode** subroutine and output the **exit_ca_mode** string before calling the shell. After returning from the shell, the program should output the **enter_ca_mode** string and call the **reset_prog_mode** subroutine. This process differs from standard curses operations, which call the **endwin** subroutine on exit.

Low-level screen subroutines

Use the following subroutines for low-level screen manipulations:

Subroutine	Description
ripoffline	Strips a single line from the stdscr
scr_dump	Dumps the contents of the virtual screen to a specified file
scr_init	Initializes the curses data structures from a specified file
scr_restore	Restores the virtual screen to the contents of a previously dumped file

termcap subroutines

If your program uses the **termcap** file for terminal information, the **termcap** subroutines are included as a conversion aid. The parameters are the same for the **termcap** subroutines. Curses emulates the subroutines using the **terminfo** database. The following **termcap** subroutines are supplied:

Subroutine Description

<u>tgetent</u>	Emulates the <u>setupterm</u> subroutine.
<u>tgetflag</u>	Returns the boolean entry for a termcap identifier.
<u>tgetnum</u>	Returns the numeric entry for a termcap identifier.
<u>tgetstr</u>	Returns the string entry for a termcap identifier.
<u>tgoto</u>	Duplicates the <u>tparam</u> subroutine. The output from the <u>tgoto</u> subroutine should be passed to the <u>tputs</u> subroutine.

Converting termcap descriptions to terminfo descriptions

The **captoinfo** command converts **termcap** descriptions to **terminfo** descriptions. The following example illustrates how the **captoinfo** command works:

```
captoinfo /usr/lib/libtermcap/termcap.src
```

This command converts the **/usr/lib/libtermcap/termcap.src** file to **terminfo** source. The **captoinfo** command writes the output to standard output and preserves comments and other information in the file.

Manipulating TTYs

The following functions save and restore the state of terminal modes:

Functions Description

<u>savetty</u>	Saves the state of the tty modes.
<u>resetty</u>	Restores the state of the tty modes to what they were the last time the <u>savetty</u> subroutine was called.

Synchronous and networked asynchronous terminals

Synchronous, networked synchronous (NWA) or non-standard directly connected asynchronous terminals are often used in a mainframe environment and communicate to the host in block mode. That is, the user types characters at the terminal, then presses a special key to initiate transmission of the characters to the host.

Note: Although it may be possible to send arbitrarily sized blocks to the host, it is not possible or desirable to cause a character to be transmitted with only a single keystroke. Doing so could cause severe problems to an application that makes use of single-character input.

Output

The curses interface can be used for all operations pertaining to output to the terminal, with the possible exception that on some terminals, the **refresh** routine may have to redraw the entire screen contents in order to perform any update.

If it is additionally necessary to clear the screen before each such operation, the result could be undesirable.

Input

Because of the nature of operation of synchronous (block-mode) and NWA terminals, it might not be possible to support all or any of the curses input functions. In particular, note the following points:

- Single-character input might not be possible. It may be necessary to press a special key to cause all characters typed at the terminal to be transmitted to the host.
- It is sometimes not possible to disable echo. Character echo may be performed directly by the terminal. On terminals that behave in this way, any curses application that performs input should be aware that any characters typed will appear on the screen at the point where the cursor is positioned. This does not necessarily correspond to the position of the cursor in the window.

Working with color

If a terminal supports color, you can use the color manipulation subroutines to include color in your curses program.

Before manipulating colors, test whether a terminal supports color. To do this, you can use either the **has_colors** subroutine or the **can_change_color** subroutine. The **can_change_color** subroutine also checks to see if a program can change the terminal's color definitions. Neither of these subroutines requires an argument.

Subroutine	Description
<u>can_change_color</u>	Checks to see if the terminal supports colors and changing of the color definition
<u>has_colors</u>	Checks that the terminal supports colors
start_color	Initializes the eight basic colors and two global variables, COLORS and COLOR_PAIRS

After you have determined that the terminal supports color, call the **start_color** subroutine before you call other color subroutines. It is a good practice to call this subroutine immediately after the **initscr** subroutine and after a successful color test. The **COLORS** global variable defines the maximum number of colors that the terminal supports. The **COLOR_PAIRS** global variable defines the maximum number of color pairs that the terminal supports.

Manipulating video attributes

Your program can manipulate a number of video attributes.

Video attributes, bit masks, and default colors

Curses enables you to control the following attributes:

A_ALTCHARSET

Alternate character set.

A_BLINK

Blinking.

A_BOLD

Extra bright or bold.

A_DIM

Half-bright.

A_NORMAL

Normal attributes.

A_REVERSE

Reverse video.

A_STANDOUT

Terminal's best highlighting mode.

A_UNDERLINE

Underline.

COLOR_PAIR (Number)

Displays the color pair represented by *Number*. You must have already initialized the color pair using the **init_pair** subroutine.

These attributes are defined in the **curses.h** file. You can pass attributes to the **wattron**, **wattroff**, and **wattrset** subroutines, or you can OR them with the characters passed to the **waddch** subroutine. The C logical OR operator is a | (pipe symbol). The following bit masks are also provided:

A_ATTRIBUTES

Extracts attributes

A_CHARTEXT

Extracts a character

A_COLOR

Extracts color-pair field information

A_NORMAL

Turns all video attributes off

The following macros are provided for working with color pairs: **COLOR_PAIR(Number)** and **PAIR_NUMBER(Attribute)**. The **COLOR_PAIR(Number)** macro and the **A_COLOR** mask are used by the **PAIR_NUMBER(Attribute)** macro to extract the color-pair number found in the attributes specified by the *Attribute* parameter.

If your program uses color, the **curses.h** file defines a number of macros that identify the following default colors:

Color	Integer Value
COLOR_BLACK	0
COLOR_BLUE	1
COLOR_GREEN	2
COLOR_CYAN	3
COLOR_RED	4
COLOR_MAGENTA	5
COLOR_YELLOW	6
COLOR_WHITE	7

Curses assumes that the default background color for all terminals is 0 (**COLOR_BLACK**).

Setting video attributes

The current window attributes are applied to all characters written into the window with the **addch** subroutines. These attributes remain as a property of the characters. The characters retain these attributes during terminal operations.

attroff or wattroff

Turns off attributes

attron or wattron

Turns on attributes

attrset or wattrset

Sets the current attributes of a window

standout, wstandout, standend, or wstandend

Puts a window into and out of the terminal's best highlight mode

vidputs or vidattr

Outputs a string that puts the terminal in a video-attribute mode

The **attrset** subroutine sets the current attributes of the default screen. The **wattrset** subroutine sets the current attributes of the user-defined window.

Use the **attron** and **attroff** subroutines to turn on and off the specified attributes in the stdscr without affecting any other attributes. The **wattron** and **wattroff** subroutines perform the same actions in user-defined windows.

The **standout** subroutine is the same as a call to the **attron** subroutine with the **A_STANDOUT** attribute. It puts the stdscr into the terminal's best highlight mode. The **wstandout** subroutine is the same as a call to the **wattron(Window, A_STANDOUT)** subroutine. It puts the user-defined window into the terminal's best highlight mode. The **standend** subroutine is the same as a call to the **attrset(0)** subroutine. It turns

off all attributes for `stdscr`. The **wstandend** subroutine is the same as a call to the **wattrset(Window, 0)** subroutine. It turns off all attributes for the specified window.

The **vidputs** subroutine outputs a string that puts the terminal in the specified attribute mode. Characters are output through the **putc** subroutine. The **vidattr** subroutine is the same as the **vidputs** subroutine except that characters are output through the **putchar** subroutine.

Working with color pairs

The **COLOR_PAIR (Number)** macro is defined in the **curses.h** file so you can manipulate color attributes as you would any other attributes. You must initialize a color pair with the **init_pair** subroutine before you use it. The **init_pair** subroutine has the following parameters: *Pair*, *Foreground*, and *Background*. The *Pair* parameter must be between 1 and **COLOR_PAIRS -1**. The *Foreground* and *Background* parameters must be between 0 and **COLORS -1**. For example, to initialize color pair 1 to a foreground of black with a background of cyan, you would use the following:

```
init_pair(1, COLOR_BLACK, COLOR_CYAN);
```

You could then set the attributes for the window as follows:

```
wattrset(win, COLOR_PAIR(1));
```

If you then write the string `Let's add color` to the terminal, the string displays as black characters on a cyan background.

Extracting attributes

You can use the results from the call to the **winch** subroutine to extract attribute information, including the color-pair number. The following example uses the value returned by a call to the **winch** subroutine with the C logical AND operator (&) and the **A_ATTRIBUTES** bit mask to extract the attributes assigned to the current position in the window. The results from this operation are used with the **PAIR_NUMBER** macro to extract the color-pair number, and the number 1 is printed on the screen.

```
win = newwin(10, 10, 0, 0);
init_pair(1, COLOR_RED, COLOR_YELLOW);
wattrset(win, COLOR_PAIR(1));
waddstr(win, "apple");

number = PAIR_NUMBER((mvwinch(win, 0, 0) & A_ATTRIBUTES));
wprintw(win, "%d\n", number);
wrefresh(win);
```

Lights and whistles

The curses library provides the following alarm subroutines to signal the user:

beep

Sounds an audible alarm on the terminal

flash

Displays a visible alarm on the terminal

Setting curses options

All curses options are initially turned off, so it is not necessary to turn them off before calling the **endwin** subroutine. The following subroutines allow you to set various options with curses:

curs_set

Sets the cursor visibility to invisible, normal, or very visible.

idlok

Specifies whether curses can use the hardware insert and delete line features of terminals so equipped.

intrflush

Specifies whether an interrupt key (interrupt, quit, or suspend) flushes all output in the tty driver. This option's default is inherited from the tty driver.

keypad

Specifies whether curses retrieves the information from the terminal's keypad. If enabled, the user can press a function key (such as an arrow key) and the **wgetch** subroutine returns a single value representing that function key. If disabled, curses will not treat the function keys specially and your program must interpret the escape sequences. For a list of these function keys, see the **wgetch** subroutine.

typeahead

Instructs curses to check for type ahead in an alternative file descriptor.

See the **wgetch** subroutines and [Understanding Terminals with Curses](#) for descriptions of additional curses options.

Manipulating soft labels

Curses provides subroutines for manipulating soft function-key labels.

These labels appear at the bottom of the screen and give applications, such as editors, a more user-friendly look. To use soft labels, you must call the **slk_init** subroutine before calling the **initscr** or **newterm** subroutines.

Subroutine	Description
<u>slk_clear</u>	Clears soft labels from the screen.
<u>slk_init</u>	Initializes soft function key labels.
<u>slk_label</u>	Returns the current label.
<u>slk_noutrefresh</u>	Refreshes soft labels. This subroutine is functionally equivalent to the wnoutrefresh subroutine.
<u>slk_refresh</u>	Refreshes soft labels. This subroutine is functionally equivalent to the refresh subroutine.
<u>slk_restore</u>	Restores the soft labels to the screen after a call to the slk_clear subroutine.
<u>slk_set</u>	Sets a soft label.
<u>slk_touch</u>	Updates soft labels on the next call to the slk_noutrefresh subroutine.

To manage soft labels, curses reduces the size of the **stdscr** by one line. It reserves this line for use by the soft-label functions. This reservation means that the **LINES** environment variable is also reduced. Many terminals support built-in soft labels. If built-in soft labels are supported, curses uses them. Otherwise, curses simulates the soft-labels with software.

Because many terminals that support soft labels have 8 labels, curses follows the same standard. A label string is restricted to 8 characters. Curses arranges labels in one of two patterns: 3-2-3 (3 left, 2 center, 3 right) or 4-4 (4 left, 4 right).

To specify a string for a particular label, call the **slk_set** subroutine. This subroutine also instructs curses to left-justify, right-justify, or center the string on the label. To obtain a label name before it was justified by the **slk_set** subroutine, use the **slk_label** subroutine. The **slk_clear** and **slk_restore** subroutines clear and restore soft labels respectively. Normally, to update soft labels, your program should call the **slk_noutrefresh** subroutine for each label and then use a single call to the **slk_refresh** subroutine to perform the actual output. To output all the soft labels on the next call to the **slk_noutrefresh** subroutine, use the **slk_touch** subroutine.

Curses compatibility

This section describes the compatibility issues.

The following compatibility issues need to be considered:

- In AIX 4.3, curses is not compatible with AT&T System V Release 3.2 curses.

- Applications compiled, rebound, or relinked may need source code changes for compatibility with the AIX Version 4 of curses. The curses library does not have or use AIX extended curses functions.
- Applications requiring multibyte support might still compile and link with extended curses. Use of the extended curses library is, however, discouraged except for applications that require multibyte support.

List of additional curses subroutines

The following sections describe additional curses subroutines:

The following sections describe additional curses subroutines:

- [Manipulating Windows](#)
- [Manipulating Characters](#)
- [Manipulating Terminals](#)
- [Manipulating Color](#)
- [Miscellaneous Utilities](#)

Manipulating windows

Use the following subroutines to manipulate windows:

Subroutine	Description
scr_dump	Writes the current contents of the virtual screen to the specified file
scr_init	Uses the contents of a specified file to initialize the curses data structures
scr_restore	Sets the virtual screen to the contents of the specified file

Manipulating characters

Use the following subroutines to manipulate characters:

Subroutine	Description
echochar , wechochar , or pechochar	Functionally equivalent to a call to the addch (or waddch) subroutine followed by a call to the refresh (or wrefresh) subroutine.
flushinp	Flushes any type-ahead characters typed by the user but not yet read by the program.
insertln or winsertln	Inserts a blank line in a window.
keyname	Returns a pointer to a character string containing a symbolic name for the <i>Key</i> parameter.
meta	Determines whether 8-bit character return for the wgetch subroutine is allowed.
nodelay	Causes a call to the wgetch subroutine to be a nonblocking call. If no input is ready, the wgetch subroutine returns ERR .
scroll	Scrolls a window up one line.
unctrl	Returns the printable representation of a character. Control characters are punctuated with a ^ (caret).
vwprintw	Performs the same operation as the wprintw subroutine, but takes a variable list of arguments.
vwscanw	Performs the same operation as the wscanw subroutine, but takes a variable list of arguments.

Manipulating terminals

Use the following subroutines to manipulate terminals:

Subroutine	Description
<u>def_prog_mode</u>	Identifies the current terminal mode as the in-curses mode
<u>def_shell_mode</u>	Saves the current terminal mode as the not-in-curses mode
<u>del_curterm</u>	Frees the space pointed to by the <i>oterm</i> variable
<u>notimeout</u>	Prevents the wgetch subroutine from setting a timer when interpreting an input escape sequence
<u>pechochar</u>	Equivalent to a call to the waddch subroutine followed by a call to the prefresh subroutine.
<u>reset_prog_mode</u>	Restores the terminal into the in-curses program mode.
<u>reset_shell_mode</u>	Restores the terminal to shell mode (out-of-curses mode). The endwin subroutine does this automatically.
<u>restartterm</u>	Sets up a TERMINAL structure for use by curses. This subroutine is similar to the setupterm subroutine. Call the restartterm subroutine after restoring memory to a previous state. For example, call this subroutine after a call to the scr_restore subroutine.

Manipulating color

Use the following subroutines to manipulate colors:

Subroutine	Description
<u>color_content</u>	Returns the composition of a color
<u>init_color</u>	Changes a color to the desired composition
<u>init_pair</u>	Initializes a color pair to the specified foreground and background colors
<u>pair_content</u>	Returns the foreground and background colors for a specified color-pair number

Miscellaneous utilities

The following miscellaneous utilities are available:

Utilities	Description
<u>baudrate</u>	Queries the current terminal and returns its output speed
<u>erasechar</u>	Returns the erase character chosen by the user
<u>killchar</u>	Returns the line-kill character chosen by the user

Debugging programs

There are several debug programs available for debugging your programs: the **adb**, **dbx**, **dex**, **softdb**, and kernel debug programs. The **adb** program enables you to debug executable binary files and examine non-ASCII data files.

The **dbx** program enables source-level debugging of C, C++, and FORTRAN language programs, as well as assembler-language debugging of executable programs at the machine level. The **(dex)** provides an X interface for the **dbx** debug program, providing windows for viewing the source, context, and variables of the application program. The **softdb** debug program works much like the **dex** debug program, but **softdb** is used with AIX Software Development Environment Workbench. The kernel debug program is used to help determine errors in code running in the kernel.

The following articles provide information on the **adb** and **dbx** debug programs:

Related information

[adb Command](#)

adb debug program overview

The **adb** command provides a general purpose debug program. You can use this command to examine object and core files and provide a controlled environment for running a program.

While the **adb** command is running, it takes standard input and writes to standard output. The command does not recognize the Quit or Interrupt keys. If these keys are used, the **adb** command waits for a new command.

Getting started with the adb debug program

This section explains how to start the **adb** debugging program from a variety of files, use the **adb** prompt, use shell commands from within the **adb** program, and stop the **adb** program.

Starting adb with a program file

You can start the **adb** debug program without a file name. In this case, the **adb** program searches for the default **a.out** file in your current working directory and prepares it for debugging. Thus, the command:

```
adb
```

is the same as entering:

```
adb a.out
```

The **adb** program starts with the **a.out** file and waits for a command. If the **a.out** file does not exist, the **adb** program starts without a file and does not display an error message.

Starting adb with a Core image file

You can use the **adb** debug program to examine the core image files of programs that caused irrecoverable system errors. Core image files maintain a record of the contents of the CPU registers, stack, and memory areas of your program at the time of the error. Therefore, core image files provide a way to determine the cause of an error.

To examine a core image file with its corresponding program, you must give the name of both the core and the program file. The command line has the form:

```
adb ProgramFile CoreFile
```

where *ProgramFile* is the file name of the program that caused the error, and *CoreFile* is the file name of the core image file generated by the system. The **adb** program then uses information from both files to provide responses to your commands.

If you do not give the filename of the core image file, the **adb** program searches for the default core file, named **core**, in your current working directory. If such a file is found, the **adb** program determines whether the core file belongs to the *ProgramFile*. If so, the **adb** program uses it. Otherwise, the **adb** program discards the core file by giving an appropriate error message.

Note: The **adb** command cannot be used to examine 64-bit objects and AIX 4.3 core format. **adb** still works with pre-AIX 4.3 core format. On AIX 4.3, user can make kernel to generate pre-AIX 4.3 style core dumps using smitty.

Starting adb with a data file

The **adb** program provides a way to look at the contents of the file in a variety of formats and structures. You can use the **adb** program to examine data files by giving the name of the data file in place of the program or core file. For example, to examine a data file named **outdata**, enter:

```
adb outdata
```

The **adb** program opens a file called **outdata** and lets you examine its contents. This method of examining files is useful if the file contains non-ASCII data. The **adb** command may display a warning

when you give the name of a non-ASCII data file in place of a program file. This usually happens when the content of the data file is similar to a program file. Like core files, data files cannot be executed.

Starting **adb** with the write option

If you open a program or data file with the **-w** flag of the **adb** command, you can make changes and corrections to the file. For example, the command:

```
adb -w sample
```

opens the program file `sample` for writing. You can then use **adb** commands to examine and modify this file. The **-w** flag causes the **adb** program to create a given file if it does not already exist. The option also lets you write directly to memory after running the given program.

Using a prompt

After you have started the **adb** program you can redefine your prompt with the **\$P** subcommand.

To change the **[adb:scat]>>** prompt to **Enter a debug command-->**, enter:

```
$P"Enter a debug command-->"
```

The quotes are not necessary when redefining the new prompt from the **adb** command line.

Using shell commands from within the **adb** program

You can run shell commands without leaving the **adb** program by using the **adb** escape command (!) (exclamation point). The escape command has the form:

! Command

In this format *Command* is the shell command you want to run. You must provide any required arguments with the command. The **adb** program passes this command to the system shell that calls it. When the command is finished, the shell returns control to the **adb** program. For example, to display the date, enter the following command:

```
! date
```

The system displays the date and restores control to the **adb** program.

Exiting the **adb** debug program

You can stop the **adb** program and return to the system shell by using the **\$q** or **\$Q** subcommands. You can also stop the **adb** program by typing the Ctrl-D key sequence. You cannot stop the **adb** program by pressing the Interrupt or Quit keys. These keys cause **adb** to wait for a new command.

Controlling program execution

This section explains the commands and subcommands necessary to prepare programs for debugging; execute programs; set, display, and delete breakpoints; continue programs; single-step through a program; stop programs; and kill programs.

Preparing programs for debugging with the **adb** program

Compile the program using the **cc** command to a file such as **adbsamp2** by entering the following:

```
cc adbsamp2.c -o adbsamp2
```

To start the debug session, enter:

```
adb adbsamp2
```

The C language does not generate statement labels for programs. Therefore, you cannot refer to individual C language statements when using the debug program. To use execution commands effectively, you must be familiar with the instructions that the C compiler generates and how those instructions relate

to individual C language statements. One useful technique is to create an assembler language listing of your C program before using the **adb** program. Then, refer to the listing as you use the debug program. To create an assembler language listing, use the **-S** or **-qList** flag of the **cc** command.

For example, to create an assembler language listing of the example program, **adbsamp2.c**, use the following command:

```
cc -S adbsamp2.c -o adbsamp2
```

This command creates the **adbsamp2.s** file, that contains the assembler language listing for the program, and compiles the program to the executable file, **adbsamp2**.

Running program

You can execute a program by using the **:r** or **:R** subcommand. The commands have the form:

```
[ Address ][,Count ] :r [Arguments ]
```

OR

```
[ Address ][,Count ] :R [Arguments ]
```

In this format, the *Address* parameter gives the address at which to start running the program; the *Count* parameter is the number of breakpoints to skip before one is taken; and the *Arguments* parameter provides the command-line arguments, such as file names and options, to pass to the program.

If you do not supply an *Address* value, the **adb** program uses the start of the program. To run the program from the beginning enter:

```
:r
```

If you supply a *Count* value, the **adb** program ignores all breakpoints until the given number has been encountered. For example, to skip the first five named breakpoints, use the command:

```
,5:r
```

If you provide arguments, separate them by at least one space each. The arguments are passed to the program in the same way the system shell passes command-line arguments to a program. You can use the shell redirection symbols.

The **:R** subcommand passes the command arguments through the shell before starting program operation. You can use shell pattern-matching characters in the arguments to refer to multiple files or other input values. The shell expands arguments containing pattern-matching characters before passing them to the program. This feature is useful if the program expects multiple file names. For example, the following command passes the argument **[a-z]*** to the shell where it is expanded to a list of the corresponding file names before being passed to the program:

```
:R [a-z]*.s
```

The **:r** and **:R** subcommands remove the contents of all registers and destroy the current stack before starting the program. This operation halts any previous copy of the program that may be running.

Setting breakpoints

To set a breakpoint in a program, use the **:b** subcommand. Breakpoints stop operation when the program reaches the specified address. Control then returns to the **adb** debug program. The command has the form:

```
[Address] [,Count ] :b [Command]
```

In this format, the *Address* parameter must be a valid instruction address; the *Count* parameter is a count of the number of times you want the breakpoint to be skipped before it causes the program to stop; and the *Command* parameter is the **adb** command you want to execute each time that the instruction

is executed (regardless of whether the breakpoint stops the program). If the specified command sets . (period) to a value of 0, the breakpoint causes a stop.

Set breakpoints to stop program execution at a specific place in the program, such as the beginning of a function, so that you can look at the contents of registers and memory. For example, when debugging the example **adbsamp2** program, the following command sets a breakpoint at the start of the function named **f**:

```
.f :b
```

The breakpoint is taken just as control enters the function and before the function's stack frame is created.

A breakpoint with a count is used within a function that is called several times during the operation of a program, or within the instructions that correspond to a **for** or **while** statement. Such a breakpoint allows the program to continue to run until the given function or instructions have been executed the specified number of times. For example, the following command sets a breakpoint for the second time that the **f** function is called in the **adbsamp2** program:

```
.f,2 :b
```

The breakpoint does not stop the function until the second time the function is run.

Displaying breakpoints

Use the **\$b** subcommand to display the location and count of each currently defined breakpoint. This command displays a list of the breakpoints by address and any count or commands specified for the breakpoints. For example, the following sets two breakpoints in the **adbsamp2** file and then uses the **\$b** subcommand to display those breakpoints:

```
.f+4:b
.f+8:b$v
$b
breakpoints
count  brkpt          command
1      .f+8           $v
1      .f+4
```

When the program runs, it stops at the first breakpoint that it finds, such as **.f+4**. If you use the **:c** subcommand to continue execution, the program stops again at the next breakpoint and starts the **\$v** subcommand. The command and response sequence looks like the following example:

```
:r
adbsamp2:running
breakpoint      .f+4:          st          r3,32(r1)
:c
adbsamp2:running
variables
b = 268435456
d = 236
e = 268435512
m = 264
breakpoint      .f+8          1          r15,32(r1)
```

Deleting breakpoints

To use the **:d** subcommand to delete a breakpoint from a program, enter:

Address :d

In this format, the *Address* parameter gives the address of the breakpoint to delete.

For example, when debugging the example **adbsamp2** program, entering the following command deletes the breakpoint at the start of the **f** function:

```
.f:d
```

Continuing program execution

To use the **:c** subcommand to continue the execution of a program after it has been stopped by a breakpoint enter:

```
[Address] [,Count] :c [Signal]
```

In this format, the *Address* parameter gives the address of the instruction at which to continue operation; the *Count* parameter gives the number of breakpoints to ignore; and the *Signal* parameter is the number of the signal to send to the program.

If you do not supply an *Address* parameter, the program starts at the next instruction after the breakpoint. If you supply a *Count* parameter, the **adb** debug program ignores the first *Count* breakpoints.

If the program is stopped using the Interrupt or Quit key, this signal is automatically passed to the program upon restarting. To prevent this signal from being passed, enter the command in the form:

```
[Address] [,Count] :c 0
```

The command argument **0** prevents a signal from being sent to the subprocess.

Single-stepping a program

Use the **:s** subcommand to run a program in single steps or one instruction at a time. This command issues an instruction and returns control to the **adb** debug program. The command has the form:

```
[Address] [,Count] :s [Signal]
```

In this format, the *Address* parameter gives the address of the instruction you want to execute, and the *Count* parameter is the number of times you want to repeat the command. If there is no current subprocess, the *ObjectFile* parameter is run as a subprocess. In this case, no signal can be sent and the remainder of the line is treated as arguments to the subprocess. If you do not supply a value for the *Address* parameter, the **adb** program uses the current address. If you supply the *Count* parameter, the **adb** program continues to issue each successive instruction until the *Count* parameter instructions have been run. Breakpoints are ignored while single-stepping. For example, the following command issues the first five instructions in the **main** function:

```
.main,5:s
```

Stopping a program with the interrupt and quit keys

Use either the Interrupt or Quit key to stop running a program at any time. Pressing either of these keys stops the current program and returns control to the **adb** program. These keys are useful with programs that have infinite loops or other program errors.

When you press the Interrupt or Quit key to stop a program, the **adb** program automatically saves the signal. If you start the program again using the **:c** command, the **adb** program automatically passes the signal to the program. This feature is useful when testing a program that uses these signals as part of its processing. To continue running the program without sending signals, use the command:

```
:c 0
```

The command argument **0** (zero) prevents a signal from being sent to the program.

Stopping a program

To stop a program you are debugging, use the **:k** subcommand. This command stops the process created for the program and returns control to the **adb** debug program. The command clears the current contents of the system unit registers and stack and begins the program again. The following example shows the use of the **:k** subcommand to clear the current process from the **adb** program:

```
:k
```

```
560:   killed
```

Using adb expressions

This section describes the use of **adb** expressions.

Using integers in expressions

When creating an expression, you can use integers in three forms: decimal, octal, and hexadecimal. Decimal integers must begin with a non-zero decimal digit. Octal numbers must begin with a 0 (zero) and have octal digits only (0-7). Hexadecimal numbers must begin with the prefix 0x and can contain decimal digits and the letters a through f (in both uppercase and lowercase). The following are examples of valid numbers:

Decimal	Octal	Hexadecimal
34	042	0x22
4090	07772	0xffa

Using symbols in expressions

Symbols are the names of global variables and functions defined within the program being debugged. Symbols are equal to the address of the given variable or function. They are stored in the program symbol table and are available if the symbol table has not been stripped from the program file.

In expressions, you can spell the symbol exactly as it is in the source program or as it has been stored in the symbol table. Symbols in the symbol table are no more than 8 characters long.

When you use the **?** subcommand, the **adb** program uses the symbols found in the symbol table of the program file to create symbolic addresses. Thus, the **?** subcommand sometimes gives a function name when displaying data. This does not happen if the **?** subcommand is used for text (instructions) and the **/** command is used for data.

Local variables can only be addressed if the C language source program is compiled with the **-g** flag.

If the C language source program is not compiled using the **-g** flag the local variable cannot be addressed. The following command displays the value of the local variable **b** in a function sample:

```
.sample.b / x - value of local variable.  
.sample.b = x - Address of local variable.
```

Using operators in expressions

You can combine integers, symbols, variables, and register names with the following operators:

Unary operators:

~ (tilde)	Bitwise complementation
- (dash)	Integer negation
* (asterisk)	Returns contents of location

Binary operators:

+ (plus)	Addition
- (minus)	Subtraction
* (asterisk)	Multiplication
% (percent)	Integer division
& (ampersand)	Bitwise conjunction
] (right bracket)	Bitwise disjunction
^ (caret)	Modulo

Binary operators:

(number sign)

Round up to the next multiple

The **adb** debug program uses 32-bit arithmetic. Values that exceed 2,147,483,647 (decimal) are displayed as negative values. The following example shows the results of assigning two different values to the variable *n*, and then displaying the value in both decimal and hexadecimal:

```
2147483647>n<
n=D
 2147483647<
n=X
 7fffffff
2147483648>n<
n=D
-2147483648<
n=X
 80000000
```

Unary operators have higher precedence than binary operators. All binary operators have the same precedence and are evaluated in order from left to right. Thus, the **adb** program evaluates the following binary expressions as shown:

```
2*3+4=d
 10
4+2*3=d
 18
```

You can change the precedence of the operations in an expression by using parentheses. The following example shows how the previous expression is changed by using parentheses:

```
4+(2*3)=d
 10
```

The unary operator, ***** (asterisk), treats the given address as a pointer into the data segment. An expression using this operator is equal to the value pointed to by that pointer. For example, the expression:

```
*0x1234
```

is equal to the value at the data address 0x1234, whereas the example:

```
0x1234
```

is equal to 0x1234.

Customizing the adb debug program

This section describes how you can customize the **adb** debug program.

Combining commands on a single line

You can give more than one command on a line by separating the commands with a ; (semicolon). The commands are performed one at a time, starting at the left. Changes to the current address and format carry over to the next command. If an error occurs, the remaining commands are ignored. For example, the following sequence displays both the **adb** variables and then the active subroutines at one point in the **adbsamp2** program:

```
$v;$c
variables
b = 10000000
d = ec
e = 10000038
m = 108
t = 2f8.
```

```
f(0,0) .main+26.  
main(0,0,0) start+fa
```

Creating adb scripts

You can direct the **adb** debug program to read commands from a text file instead of from the keyboard by redirecting the standard input file when you start the **adb** program. To redirect standard input, use the input redirection symbol, < (less than), and supply a file name. For example, use the following command to read commands from the file script:

```
adb sample <script
```

The file must contain valid **adb** subcommands. Use the **adb** program script files when the same set of commands can be used for several different object files. Scripts can display the contents of core files after a program error. The following example shows a file containing commands that display information about a program error. When that file is used as input to the **adb** program using the following command to debug the **adbsamp2** file, the specified output is produced.

```
120$w  
4095$s.  
f:b:  
r  
=1n"=====  
adb Variables =====  
$v  
=1n"=====  
Address Map =====  
$m  
=1n"=====  
C Stack Backtrace =====  
$C  
=1n"=====  
C External Variables =====  
$e  
=1n"=====  
Registers =====  
$r  
0$s  
=1n"=====  
Data Segment =====<  
b,10/8xna
```

```
$ adb adbsamp2 <script
```

```
adbsamp2: running  
breakpoint .f: b .f+24  
=====  
adb Variables =====  
variables  
0 = TBD  
1 = TBD  
2 = TBD  
9 = TBD  
b = 10000000  
d = ec  
e = 10000038  
m = 108  
t = 2f8  
=====  
Address Map =====  
[0]? map .adbsamp2.  
b1 = 10000000 e1 = 100002f8 f1 = 0  
b2 = 200002f8 e2 = 200003e4 f2 = 2f8  
[0]/ map .-.  
b1 = 0 e1 = 0 f1 = 0  
b2 = 0 e2 = 0 f2 = 0  
=====  
C Stack Backtrace =====.  
f(0,0) .main+26.  
main(0,0,0) start+fa  
=====  
C External Variables =====Full word.  
errno: 0.  
environ: 3fffe6bc.  
NLinit: 10000238.  
main: 100001ea.  
exit: 1000028c.  
fcnt: 0
```

```
.loop .count: 1.  
f: 100001b4.  
NLgetfile: 10000280.  
write: 100002e0.
```

```

NLinit. .X: 10000238 .
NLgetfile. .X: 10000280 .
cleanup: 100002bc.
exit: 100002c8 .
exit . .X: 1000028c . .
cleanup . .X: 100002bc

```

```

===== Registers =====
mq 20003a24 .errno+3634
cs 100000 gt
ics 1000004
pc 100001b4 .f
r15 10000210 .main+26
r14 20000388 .main
r13 200003ec .loop .count
r12 3fffe3d0
r11 3fffe44c
r10 0
r9 20004bcc
r8 200041d8 .errno+3de8
r7 0
r6 200030bc .errno+2ccc
r5 1
r4 200003ec .loop .count
r3 f4240
r2 1
r1 3fffe678
r0 20000380 .f.
f: b .f+24

```

```

===== Data Segment =====
10000000: 103 5313 3800 0 0 2f8 0 ec
10000010: 0 10 1000 38 0 0 0 1f0
10000020: 0 0 0 0 1000 0 2000 2f8
10000030: 0 0 0 0 4 6000 0 6000
10000040: 6e10 61d0 9430 a67 6730 6820 c82e 8
10000050: 8df0 94 cd0e 60 6520 a424 a432 c84e
10000060: 8 8df0 77 cd0e 64 6270 8df0 86
10000070: cd0e 60 6520 a424 a432 6470 8df0 6a
10000080: cd0e 64 c82e 19 8df0 78 cd0e 60
10000090: 6520 a424 a432 c84e 19 8df0 5b cd0e
100000a0: 64 cd2e 5c 7022 d408 64 911 c82e
100000b0: 2e 8df0 63 cd0e 60 6520 a424 a432
100000c0: c84e 2e 8df0 46 cd0e 64 15 6280
100000d0: 8df0 60 cd0e 68 c82e 3f 8df0 4e
100000e0: cd0e 60 6520 a424 a432 c84e 3f 8df0
100000f0: 31 cd0e 64 c820 14 8df0 2b cd0e
10000100:

```

Setting output width

Use the **\$w** subcommand to set the maximum width (in characters) of each line of output created by the **adb** program. The command has the form:

*Width***\$w**

In this format, the *Width* parameter is an integer that specifies the width in characters of the display. You can give any width convenient for your display device. When the **adb** program is first invoked, the default width is 80 characters.

This command can be used when redirecting output to a line printer or special output device. For example, the following command sets the display width to 120 characters, a common maximum width for line printers:

```
120$w
```

Setting the maximum offset

The **adb** debug program normally displays memory and file addresses as the sum of a symbol and an offset. This format helps to associate the instructions and data on the display with a particular function or variable. When the **adb** program starts up, it sets the maximum offset to 255, so that symbolic addresses

are assigned only to instructions or data that occur less than 256 bytes from the start of the function or variable. Instructions or data beyond that point are given numeric addresses.

In many programs, the size of a function or variable is actually larger than 255 bytes. For this reason the **adb** program lets you change the maximum offset to accommodate larger programs. You can change the maximum offset by using the **\$s** subcommand.

The subcommand has the form:

*Offset***\$s**

In this format, the *Offset* parameter is an integer that specifies the new offset. For example, the following command increases the maximum possible offset to 4095:

```
4095$s
```

All instructions and data that are less than 4096 bytes away are given symbolic addresses. You can disable all symbolic addressing by setting the maximum offset to zero. All addresses are given numeric values instead.

Setting default input format

To alter the default format for numbers used in commands, use the **\$d** or **\$o** (octal) subcommands. The default format tells the **adb** debug program how to interpret numbers that do not begin with 0 (octal) or 0x (hexadecimal), and how to display numbers when no specific format is given. Use these commands to work with a combination of decimal, octal, and hexadecimal numbers.

The **\$o** subcommand sets the radix to 8 and thus sets the default format for numbers used in commands to octal. After you enter that subcommand, the **adb** program displays all numbers in octal format except those specified in some other format.

The format for the **\$d** subcommand is the *Radix***\$d** command, where the *Radix* parameter is the new value of the radix. If the *Radix* parameter is not specified, the **\$d** subcommand sets the radix to a default value of 16. When you first start the **adb** program, the default format is hexadecimal. If you change the default format, you can restore it as necessary by entering the **\$d** subcommand by itself:

```
$d
```

To set the default format to decimal, use the following command:

```
0xa$d
```

Changing the disassembly mode

Use the **\$i** and **\$n** subcommands to force the **adb** debug program to disassemble instructions using the specified instruction set and mnemonics. The **\$i** subcommand specifies the instruction set to be used for disassembly. The **\$n** subcommand specifies the mnemonics to be used in disassembly.

If no value is entered, these commands display the current settings.

The **\$i** subcommand accepts the following values:

com

Specifies the instruction set for the common intersection mode of the PowerPC® and POWER® family.

pwr

Specifies the instruction set and mnemonics for the POWER implementation of the POWER Architecture.

pwrx

Specifies the instruction set and mnemonics for the POWER2 implementation of the POWER family.

ppc

Specifies the instruction set and mnemonics for the PowerPC.

601

Specifies the instruction set and mnemonics for the PowerPC 601 RISC Microprocessor.

603

Specifies the instruction set and mnemonics for the PowerPC 603 RISC Microprocessor.

604

Specifies the instruction set and mnemonics for the PowerPC 604 RISC Microprocessor.

ANY

Specifies any valid instruction. For instruction sets that overlap, the mnemonics will default to PowerPC mnemonics.

The **\$n** subcommand accepts the following values:

pwr

Specifies the instruction set and mnemonics for the POWER implementation of the POWER Architecture.

ppc

Specifies the mnemonics for the PowerPC architecture.

Computing numbers and displaying text

You can perform arithmetic calculations while in the **adb** debug program by using the **=** (equal sign) subcommand. This command directs the **adb** program to display the value of an expression in a specified format.

The command converts numbers in one base to another, double-checks the arithmetic performed by a program, and displays complex addresses in simpler form. For example, the following command displays the hexadecimal number 0x2a as the decimal number 42:

```
0x2a=d
    42
```

Similarly, the following command displays 0x2a as the ASCII character * (asterisk):

```
0x2a=c
    *
```

Expressions in a command can have any combination of symbols and operators. For example, the following command computes a value using the contents of the **r0** and **r1** registers and the **adb** variable **b**.

```
<r0-12*<r1+<b+5=X
    8fa86f95
```

You can also compute the value of external symbols to check the hexadecimal value of an external symbol address, by entering:

```
main+5=X
    2000038d
```

The **=** (equal sign) subcommand can also display literal strings. Use this feature in the **adb** program scripts to display comments about the script as it performs its commands. For example, the following subcommand creates three lines of spaces and then prints the message **C Stack Backtrace**:

```
=3n"C Stack Backtrace"
```

Displaying and manipulating the source file with the adb program

This section describes how to use the **adb** program to display and manipulate the source file.

Displaying instructions and data

The **adb** program provides several subcommands for displaying the instructions and data of a given program and the data of a given data file. The subcommands and their formats are:

Display address

Address [, Count] = Format

Display instruction

Address [, Count] ? Format

Display value of variable

Address [, Count] / Format

In this format, the symbols and variables have the following meaning:

Address

Gives the location of the instruction or data item.

Count

Gives the number of items to be displayed.

Format

Defines how to display the items.

=

Displays the address of an item.

?

Displays the instructions in a text segment.

/

Displays the value of variables.

Forming addresses

In the **adb** program addresses are 32-bit values that indicate a specific memory address. They can, however, be represented in the following forms:

Absolute address

The 32-bit value is represented by an 8-digit hexadecimal number, or its equivalent in one of the other number-base systems.

Symbol name

The location of a symbol defined in the program can be represented by the name of that symbol in the program.

Entry points

The entry point to a routine is represented by the name of the routine preceded by a . (period). For example, to refer to the address of the start of the `main` routine, use the following notation:

```
.main
```

Displacements

Other points in the program can be referred to by using displacements from entry points in the program. For example, the following notation references the instruction that is 4 bytes past the entry point for the symbol `main`:

```
.main+4
```

Displaying an address

Use the = (equal sign) subcommand to display an address in a given format. This command displays instruction and data addresses in a simpler form and can display the results of arithmetic expressions. For example, entering:

```
main=an
```

displays the address of the symbol `main`:

```
10000370:
```

The following example shows a command that displays (in decimal) the sum of the internal variable **b** and the hexadecimal value 0x2000, together with its output:

```
<b+0x2000=D
268443648
```

If a count is given, the same value is repeated that number of times. The following example shows a command that displays the value of `main` twice and the output that it produces:

```
main,2=x
370 370
```

If no address is given, the current address is used. After running the above command once (setting the current address to `main`), the following command repeats that function:

```
,2=x
370 370
```

If you do not specify a format, the **adb** debug program uses the last format that was used with this command. For example, in the following sequence of commands, both `main` and `one` are displayed in hexadecimal:

```
main=x
370
one=
33c
```

Displaying the C stack backtrace

To trace the path of all active functions, use the **\$c** subcommand. This subcommand lists the names of all functions that have been called and have not yet returned control. It also lists the address from which each function was called and the arguments passed to each function. For example, the following command sequence sets a breakpoint at the function address `.f+2` in the **adbsamp2** program. The breakpoint calls the **\$c** subcommand. The program is started, runs to the breakpoint, and then displays a backtrace of the called C language functions:

```
.f+2:b$c
:r
adbsamp2:running
.f(0,0) .main+26
.main(0,0,0) start+fa
breakpoint          f+2:          tgte          r2,r2
```

By default, the **\$c** subcommand displays all calls. To display fewer calls, supply a count of the number of calls to display. For example, the following command displays only one of the active functions at the preceding breakpoint:

```
,1$c
```

```
.f(0,0) .main+26
```

Choosing data formats

A *format* is a letter or character that defines how data is to be displayed. The following are the most commonly used formats:

- a** The current symbolic address
- b** One byte in octal (displays data associated with instructions, or the high or low byte of a register)
- c** One byte as a character (char variables)

- d** Halfword in decimal (short variables)
- D** Fullword in decimal (long variables)
- i** Machine instructions in mnemonic format
- n** A new line
- o** Halfword in octal (short variables)
- O** Fullword in octal (long variables)
- r** A blank space
- s** A null-terminated character string (null-terminated arrays of char variables)
- t** A horizontal tab
- u** Halfword as an unsigned integer (short variables)
- x** Halfword in hexadecimal (short variables)
- X** Fullword in hexadecimal (long variables)

For example, the following commands produce the indicated output when using the **adbsamp** example program:

```

main=o
  1560
main=O
  4000001560
main=d
  880
main=D
  536871792
main=x
  370
main=X
  20000370
main=u
  880

```

A format can be used by itself or combined with other formats to present a combination of data in different forms. You can combine the **a**, **n**, **r**, and **t** formats with other formats to make the display more readable.

Changing the memory map

You can change the values of a memory map by using the **?m** and **/m** subcommands. These commands assign specified values to the corresponding map entries. The commands have the form:

```

[ ,count] ?m b1 e1 f1
[ ,count] /m b1 e1 f2

```

The following example shows the results of these commands on the memory map displayed with the **\$m** subcommand in the previous example:

```
,0?m 10000100 10000470 0
/m 100 100 100
$m
[0] : ?map : 'adbsamp3'
b1 = 0x10000100, e1 = 10000470, f1 = 0
b2 = 0x20000600, e2 = 0x2002c8a4, f2 = 0x600

[1] : ?map : 'shr.o' in library '/usr/ccs/lib/libc.a'
b1 = 0xd00d6200, e1 = 0xd01397bf, f1 = 0xd00defbc
b2 = 0x20000600, e2 = 0x2002beb8, f2 = 0x4a36c

[-] : /map : '-'
b1 = 100, e1 = 100, f1 = 100
b2 = 0, e2 = 0, f2 = 0
```

To change the data segment values, add an ***** (asterisk) after the **/** or **?**.

```
,0?*m 20000270 20000374 270
/*m 200 200 200
$m
[0] : ?map : 'adbsamp3'
b1 = 0x10000100, e1 = 10000470, f1 = 0
b2 = 0x20000270, e2 = 0x20000374, f2 = 0x270

[1] : ?map : 'shr.o' in library '/usr/ccs/lib/libc.a'
b1 = 0xd00d6200, e1 = 0xd01397bf, f1 = 0xd00defbc
b2 = 0x20000600, e2 = 0x2002beb8, f2 = 0x4a36c

[-] : /map : '-'
b1 = 100, e1 = 100, f1 = 100
b2 = 0, e2 = 0, f2 = 0
```

Patching binary files

You can make corrections or changes to any file, including executable binary files, by starting the **adb** program with the **-w** flag and by using the **w** and **W** () subcommands.

Locating values in a file

Locate specific values in a file by using the **l** and **L** subcommands. The subcommands have the form:

?l Value

OR

/l Value

The search starts at the current address and looks for the expression indicated by *Value*. The **l** subcommand searches for 2-byte values. The **L** subcommand searches for 4-byte values.

The **?l** subcommand starts the search at the current address and continues until the first match or the end of the file. If the value is found, the current address is set to that value's address. For example, the following command searches for the first occurrence of the **f** symbol in the **adbsamp2** file:

```
?l .f.
write+a2
```

The value is found at **.write+a2** and the current address is set to that address.

Writing to a file

Write to a file by using the **w** and **W** subcommands. The subcommands have the form:

[*Address*] **?w Value**

In this format, the *Address* parameter is the address of the value you want to change, and the *Value* parameter is the new value. The **w** subcommand writes 2-byte values. The **W** subcommand writes 4-byte values. For example, the following commands change the word "This" to "The":

```
?l .Th.  
?w .The.
```

The **W** subcommand changes all four characters.

Making changes to memory

Make changes to memory whenever a program has run. If you have used an **:r** subcommand with a breakpoint to start program operation, subsequent **w** subcommands cause the **adb** program to write to the program in memory rather than to the file. This command is used to make changes to a program's data as it runs, such as temporarily changing the value of program flags or variables.

Using adb variables

The **adb** debug program automatically creates a set of its own variables when it starts. These variables are set to the addresses and sizes of various parts of the program file as defined in the following table:

0	Last value printed
1	Last displacement part of an instruction source
2	Previous value of the 1 variable
9	Count on the last \$< or \$<< command
b	Base address of the data segment
d	Size of the data segment
e	Entry address of the program
m	"Magic" number
s	Size of the stack segment
t	Size of the text segment

The **adb** debug program reads the program file to find the values for these variables. If the file does not seem to be a program file, then the **adb** program leaves the values undefined.

To display the values that the **adb** debug program assigns to these variables, use the **\$v** subcommand. This subcommand lists the variable names followed by their values in the current format. The subcommand displays any variable whose value is not 0 (zero). If a variable also has a non-zero segment value, the variable's value is displayed as an address. Otherwise, it is displayed as a number. The following example shows the use of this command to display the variable values for the sample program **adbsamp**:

```
$v
```

```
Variables
```

```
0 = undefined
```

```
1 = undefined
```

```
2 = undefined
```

```
9 = undefined
```

```
b = 10000000
```

```
d = 130
```

```
e = 10000038
```

```
m = 108
```

```
t = 298
```

Specify the current value of an **adb** variable in an expression by preceding the variable name with < (less than sign). The following example displays the current value of the **b** base variable:

```
<b=X
```

```
10000000
```

Create your own variables or change the value of an existing variable by assigning a value to a variable name with > (greater than sign). The assignment has the form:

Expression > *VariableName*

where the *Expression* parameter is the value to be assigned to the variable and the *VariableName* parameter is the variable to receive the value. The *VariableName* parameter must be a single letter. For example, the assignment:

```
0x2000>b
```

assigns the hexadecimal value 0x2000 to the **b** variable. Display the contents of **b** again to show that the assignment occurred:

```
<b=X
```

```
2000
```

Finding the current address

The **adb** program has two special variables that keep track of the last address used in a command and the last address typed with a command. The **.** (period) variable, also called the current address, contains the last address used in a command. The **"** (double quotation mark) variable contains the last address typed with a command. The **.** and **"** variables usually contain the same address except when implied commands, such as the newline and **^** (caret) characters, are used. These characters automatically increase and decrease the **.** variable but leave the **)** (right parenthesis) variable unchanged.

Both the **.** and the **"** variables can be used in any expression. The < (less than sign) is not required. For example, the following commands display these variables at the start of debugging with the **adbsamp** program:

```
. =  
  0.  
 =  
  0
```

Displaying external variables

Use the **\$e** subcommand to display the values of all external variables in the **adb** program. External variables are the variables in your program that have global scope or have been defined outside of any function, and include variables defined in library routines used by your program, as well as all external variables of shared libraries.

The **\$e** subcommand is useful to get a list of the names for all available variables or a summary of their values. The command displays one name on each line with the variable's value (if any) on the same line. If the *Count* parameter is specified, only the external variables associated with that file are printed.

The following example illustrates the setting of a breakpoint in the **adbsamp2** sample program that calls the **\$e** subcommand, and the output that results when the program runs (be sure to delete any previous breakpoints that you may have set):

```
.f+2:b,0$e
:r
adbsamp2: running
_errno: 0
_environ: 3fffe6bc
__Nlinit: 10000238
_main: 100001ea
_exit: 1000028c
_fcnt: 0
_loop_count: 1
_f: 100001b4
_NLgetfile: 10000280
_write: 100002e0
__Nlinit__X: 10000238
_NLgetfile__X: 10000280
__cleanup: 100002bc
__exit: 100002c8
_exit__X: 1000028c
__cleanup__X: 100002bc
breakpoint .f+2: st r2,1c(r1)
```

Displaying the address maps

The **adb** program prepares a set of maps for the text and data segments in your program and uses these maps to access items that you request for display. Use the **\$m** subcommand to display the contents of the address maps. The subcommand displays the maps for all segments in the program and uses information taken from either the program and core files or directly from memory.

The **\$m** subcommand displays information similar to the following:

```
$m
[0] : ?map : 'adbsamp3'
      b1 = 0x10000200, e1 = 0x10001839, f1 = 0x10000200
      b2 = 0x2002c604, e2 = 0x2002c8a4, f2 = 0x600

[1] : ?map : 'shr.o' in library 'lib/libc.a'
      b1 = 0xd00d6200, e1 = 0xd013976f, f1 = 0xd00defbc
      b2 = 0x20000600, e2 = 0x2002bcb8, f2 = 0x4a36c

[-] : /map : '-'
      b1 = 0x00000000, e1 = 0x00000000, f1 = 0x00000000
      b2 = 0x00000000, e2 = 0x00000000, f2 = 0x00000000
```

The display defines address-mapping parameters for the text (b1, e1, and f1) and data (b2, e2, and f2) segments for the two files being used by the **adb** debug program. This example shows values for the **adbsamp3** sample program only. The second set of map values are for the core file being used. Since none was in use, the example shows the file name as - (dash).

The value displayed inside the square brackets can be used as the *Count* parameter in the **?e** and **?m** subcommands.

adb debug program reference information

The **adb** debug program uses addresses, expressions, operators, subcommands, and variables to organize and manipulate data.

adb debug program addresses

The address in a file associated with a written address is determined by a mapping associated with that file. Each mapping is represented by two triples (*B1*, *E1*, *F1*) and (*B2*, *E2*, *F2*). The *FileAddress* parameter that corresponds to a written *Address* parameter is calculated as follows:

$$B1 \leftarrow \text{Address} \langle E1 \rangle \text{FileAddress} = \text{Address} + F1 - B1$$

OR

$B2 \leq \text{Address} < E2 \Rightarrow \text{FileAddress} = \text{Address} + F2 - B2$

If the requested *Address* parameter is neither between *B1* and *E1* nor between *B2* and *E2*, the *Address* parameter is not valid. In some cases, such as programs with separated I and D space, the two segments for a file may overlap. If a ? (question mark) or / (slash) subcommand is followed by an * (asterisk), only the second triple is used.

The initial setting of both mappings is suitable for normal **a.out** and **core** files. If either file is not of the kind expected, the *B1* parameter for that file is set to a value of 0, the *E1* parameter is set to the maximum file size, and the *F1* parameter is set to a value of 0. In this way, the whole file can be examined with no address translation.

adb debug program expressions

The following expressions are supported by the **adb** debug program:

. (period)

Specifies the last address used by a subcommand. The last address is also known as the current address.

+ (plus)

Increases the value of . (period) by the current increment.

^ (caret)

Decreases the value of . (period) by the current increment.

" (double quotes)

Specifies the last address typed by a command.

Integer

Specifies an octal number if this parameter begins with **0o**, a hexadecimal number if preceded by **0x** or **#**, or a decimal number if preceded by **0t**. Otherwise, this expression specifies a number interpreted in the current radix. Initially, the radix is 16.

` Cccc `

Specifies the ASCII value of up to 4 characters. A \ (backslash) can be used to escape an ' (apostrophe).

< Name

Reads the current value of the *Name* parameter. The *Name* parameter is either a variable name or a register name. The **adb** command maintains a number of variables named by single letters or digits. If the *Name* parameter is a register name, the value of the register is obtained from the system header in the *CoreFile* parameter. Use the **\$r** subcommand to see the valid register names.

Symbol

Specifies a sequence of uppercase or lowercase letters, underscores, or digits, though the sequence cannot start with a digit. The value of the *Symbol* parameter is taken from the symbol table in the *ObjectFile* parameter. An initial _ (underscore) is prefixed to the *Symbol* parameter, if needed.

. Symbol

Specifies the entry point of the function named by the *Symbol* parameter.

Routine.Name

Specifies the address of the *Name* parameter in the specified C language routine. Both the *Routine* and *Name* parameters are *Symbol* parameters. If the *Name* parameter is omitted, the value is the address of the most recently activated C stack frame corresponding to the *Routine* parameter.

(Expression)

Specifies the value of the expression.

adb debug program operators

Integers, symbols, variables, and register names can be combined with the following operators:

Unary operators

*** Expression**

Returns contents of the location addressed by the *Expression* parameter in the *CoreFile* parameter.

@ Expression

Returns contents of the location addressed by the *Expression* parameter in the *ObjectFile* parameter.

- Expression

Performs integer negation.

~ Expression

Performs bit-wise complement.

Expression

Performs logical negation.

Binary operators**Expression1+Expression2**

Performs integer addition.

Expression1-Expression2

Performs integer subtraction.

Expression1*Expression2

Performs integer multiplication.

Expression1%Expression2

Performs integer division.

Expression1&Expression2

Performs bit-wise conjunction.

Expression1|Expression2

Performs bit-wise disjunction.

Expression1#Expression2

Rounds up the *Expression1* parameter to the next multiple of the *Expression2* parameter.

Binary operators are left-associative and are less binding than unary operators.

adb debug program subcommands

You can display the contents of a text or data segment with the **?** (question mark) or the **/** (slash) subcommand. The **=** (equal sign) subcommand displays a given address in the specified format. The **?** and **/** subcommands can be followed by an ***** (asterisk).

?Format

Displays the contents of the *ObjectFile* parameter starting at the *Address* parameter. The value of . (period) increases by the sum of the increment for each format letter.

/Format

Displays the contents of the *CoreFile* parameter starting at the *Address* parameter. The value of . (period) increases by the sum of the increment for each format letter.

=Format

Displays the value of the *Address* parameter. The **i** and **s** format letters are not meaningful for this command.

The *Format* parameter consists of one or more characters that specify print style. Each format character may be preceded by a decimal integer that is a repeat count for the format character. While stepping through a format, the . (period) increments by the amount given for each format letter. If no format is given, the last format is used.

The available format letters are as follows:

a

Prints the value of . (period) in symbolic form. Symbols are checked to ensure that they have an appropriate type.

b

Prints the addressed byte in the current radix, unsigned.

- c**
Prints the addressed character.
- C**
Prints the addressed character using the following escape conventions:
- Prints control characters as ~ (tilde) followed by the corresponding printing character.
 - Prints nonprintable characters as ~ (tilde) <Number>, where *Number* specifies the hexadecimal value of the character. The ~ character prints as ~ ~ (tilde tilde).
- d**
Prints in decimal.
- D**
Prints long decimal.
- f**
Prints the 32-bit value as a floating-point number.
- F**
Prints double floating point.
- i Number**
Prints as instructions. *Number* is the number of bytes occupied by the instruction.
- n**
Prints a new line.
- o**
Prints 2 bytes in octal.
- O**
Prints 4 bytes in octal.
- p**
Prints the addressed value in symbolic form using the same rules for symbol lookup as the **a** format letter.
- q**
Prints 2 bytes in the current radix, unsigned.
- Q**
Prints 4 unsigned bytes in the current radix.
- r**
Prints a space.
- s Number**
Prints the addressed character until a zero character is reached.
- S Number**
Prints a string using the ~ (tilde) escape convention. The *Number* variable specifies the length of the string including its zero terminator.
- t**
Tabs to the next appropriate tab stop when preceded by an integer. For example, the **8t** format command moves to the next 8-space tab stop.
- u**
Prints as an unsigned decimal number.
- U**
Prints a long unsigned decimal.
- x**
Prints 2 bytes in hexadecimal.
- X**
Prints 4 bytes in hexadecimal.
- Y**
Prints 4 bytes in date format.

- / Local or global data symbol.
- ? Local or global text symbol.
- = Local or global absolute symbol.
- "..." Prints the enclosed string.
- ^ Decreases the . (period) by the current increment. Nothing prints.
- + Increases the . (period) by a value of 1. Nothing prints.
- Decreases the . (period) decrements by a value of 1. Nothing prints.

newline

Repeats the previous command incremented with a *Count* of 1.

[?/]!Value Mask

Words starting at the . (period) are masked with the *Mask* value and compared with the *Value* parameter until a match is found. If **L** is used, the match is for 4 bytes at a time instead of 2 bytes. If no match is found, then . (period) is unchanged; otherwise, . (period) is set to the matched location. If the *Mask* parameter is omitted, a value of -1 is used.

[?/]wValue...

Writes the 2-byte *Value* parameter into the addressed location. If the command is **W**, write 4 bytes. If the command is **V**, write 1 byte. Alignment restrictions may apply when using the **w** or **W** command.

[,Count][?/]m B1 E1 F1[?/]

Records new values for the *B1*, *E1*, and *F1* parameters. If less than three expressions are given, the remaining map parameters are left unchanged. If the ? (question mark) or / (slash) is followed by an * (asterisk), the second segment (*B2*, *E2*, *F2*) of the mapping is changed. If the list is terminated by ? or /, the file (*ObjectFile* or *CoreFile*, respectively) is used for subsequent requests. (For example, the **/m?** command causes / to refer to the *ObjectFile* file. If the *Count* parameter is specified, the **adb** command changes the maps associated with that file or library only. The **\$m** command shows the count that corresponds to a particular file. If the *Count* parameter is not specified, a default value of 0 is used.

>Name

Assigns a . (period) to the variable or register specified by the *Name* parameter.

!

Calls a shell to read the line following ! (exclamation mark).

\$ Modifier

Miscellaneous commands. The available values for *Modifier* are:

<File

Reads commands from the specified file and returns to standard input. If a count is given as 0, the command will be ignored. The value of the count is placed in the **adb 9** variable before the first command in the *File* parameter is executed.

<<File

Reads commands from the specified file and returns to standard input. The **<<File** command can be used in a file without causing the file to be closed. If a count is given as 0, the command is ignored. The value of the count is placed in the **adb 9** variable before the first command in *File* is executed. The **adb 9** variable is saved during the execution of the **<<File** command and restored when **<<File** completes. There is a limit to the number of **<<File** commands that can be open at once.

>File

Sends output to the specified file. If the *File* parameter is omitted, output returns to standard output. The *File* parameter is created if it does not exist.

b

Prints all breakpoints and their associated counts and commands.

c

Stacks back trace. If the *Address* parameter is given, it is taken as the address of the current frame (instead of using the frame pointer register). If the format letter **C** is used, the names and values of all automatic and static variables are printed for each active function. If the *Count* parameter is given, only the number of frames specified by the *Count* parameter are printed.

d

Sets the current radix to the *Address* value or a value of 16 if no address is specified.

e

Prints the names and values of external variables. If a count is specified, only the external variables associated with that file are printed.

f

Prints the floating-point registers in hexadecimal.

i instruction set

Selects the instruction set to be used for disassembly.

I

Changes the default directory as specified by the **-I** flag to the *Name* parameter value.

m

Prints the address map.

n mnem_set

Selects the mnemonics to be used for disassembly.

o

Sets the current radix to a value of 8.

q

Exits the **adb** command.

r

Prints the general registers and the instruction addressed by **iar** and sets the . (period) to **iar**. The *Number***\$r** parameter prints the register specified by the *Number* variable. The *Number,Count***\$r** parameter prints registers *Number+Count-1,...,Number*.

s

Sets the limit for symbol matches to the *Address* value. The default is a value of 255.

v

Prints all non-zero variables in octal.

w

Sets the output page width for the *Address* parameter. The default is 80.

P Name

Uses the *Name* value as a prompt string.

?

Prints the process ID, the signal that caused stoppage or termination, and the registers of **\$r**.

: Modifier

Manages a subprocess. Available modifiers are:

bCommand

Sets the breakpoint at the *Address* parameter. The breakpoint runs the *Count* parameter -1 times before causing a stop. Each time the breakpoint is encountered, the specified command runs. If this command sets . (period) to a value of 0, the breakpoint causes a stop.

cSignal

Continues the subprocess with the specified signal. If the *Address* parameter is given, the subprocess is continued at this address. If no signal is specified, the signal that caused the subprocess to stop is sent. Breakpoint skipping is the same as for the **r** modifier.

d

Deletes the breakpoint at the *Address* parameter.

k

Stops the current subprocess, if one is running.

r

Runs the *ObjectFile* parameter as a subprocess. If the *Address* parameter is given explicitly, the program is entered at this point. Otherwise, the program is entered at its standard entry point. The *Count* parameter specifies how many breakpoints are to be ignored before stopping. Arguments to the subprocess can be supplied on the same line as the command. An argument starting with **<** or **>** establishes standard input or output for the command. On entry to the subprocess, all signals are turned on.

sSignal

Continues the subprocess in single steps up to the number specified in the *Count* parameter. If there is no current subprocess, the *ObjectFile* parameter is run as a subprocess. In this case no signal can be sent. The remainder of the line is treated as arguments to the subprocess.

adb debug program variables

The **adb** command provides a number of variables. When the **adb** program is started, the following variables are set from the system header in the specified core file. If the *CoreFile* parameter does not appear to be a **core** file, these values are set from the *ObjectFile* parameter:

0

Last value printed

1

Last displacement part of an instruction source

2

Previous value of the **1** variable

9

Count on the last **\$<** or **\$<<** subcommand

b

Base address of the data segment

d

Size of the data segment

e

Entry address of the program

m

"Magic" number

s

Size of the stack segment

t

Size of the text segment

Example adb program: adbsamp

The following sample program is used in this example:

```
/* Program Listing for adbsamp.c */
char str1[ ] = "This is a character string";
int one = 1;
int number = 456;
long lnum = 1234;
```

```

float fpt = 1.25;
char str2[ ] = "This is the second character string";
main()
{
    one = 2;
    printf("First String = %s\n",str1);
    printf("one = %d\n",one);
    printf("Number = %d\n",lnum);
    printf("Floating point Number = %g\n",fpt);
    printf("Second String = %s\n",str2);
}

```

Compile the program using the **cc** command to the **adbsamp** file as follows:

```
cc -g adbsamp.c -o adbsamp
```

To start the debug session, enter:

```
adb adbsamp
```

Example adb program: adbsamp2

The following sample program is used in this example:

```

/*program listing for adbsamp2.c*/
int    fcnt,loop_count;

f(a,b)
int a,b;
{
    a = a+b;
    fcnt++;
    return(a);
}
main()
{
    loop_count = 0;
    while(loop_count <= 100)
    {
        loop_count = f(loop_count,1);
        printf("%s%d\n","Loop count is: ", loop_count);
        printf("%s%d\n","fcnt count is: ",fcnt);
    }
}

```

Compile the program using the **cc** command to the **adbsamp2** file with the following command:

```
cc -g adbsamp2.c -o adbsamp2
```

To start the debug session, enter:

```
adb adbsamp2
```

Example adb program: adbsamp3

The following sample program **adbsamp3.c** contains an infinite recursion of subfunction calls.

If you run this program to completion, it causes a memory fault error and quits.

```

int    fcnt,gcnt,hcnt;
h(x,y)
int x,y;
{
    int hi;
    register int hr;
    hi = x+1;
    hr = x-y+1;
    hcnt++;
    hj:
    f(hr,hi);
}

```

```

}
g(p,q)
int p,q;
{
    int gi;
    register int gr;
    gi = q-p;
    gr = q-p+1;
    gcnt++;
    gj:
    h(gr,gi);
}
f(a,b)
int a,b;
{
    int fi;
    register int fr;
    fi = a+2*b;
    fr = a+b;
    fcnt++;
    fj:
    g(fr,fi);
}
main()
{
    f(1,1);
}

```

Compile the program using the **cc** command to create the **adbsamp3** file with the following command:

```
cc -g adbsamp3.c -o adbsamp3
```

To start the debug session, enter:

```
adb adbsamp3
```

Example of directory and i-node dumps in adb debugging

This example shows how to create **adb** scripts to display the contents of a directory and the i-node map of a file system. In the example, the directory is named **dir** and contains a variety of files.

The file system is associated with the `/dev/hd3` device file (`/tmp`), which has the necessary permissions to be read by the user.

To display a directory, create an appropriate script. A directory normally contains one or more entries. Each entry consists of an unsigned i-node number (i-number) and a 14-character file name. You can display this information by including a command in your script file. The **adb** debug program expects the object file to be an **xcoff** format file. This is not the case with a directory. The **adb** program indicates that the directory, because it is not an **xcoff** format file, has a text length of 0. Use the **m** command to indicate to the **adb** program that this directory has a text length of greater than 0. Therefore, display entries in your **adb** session by entering:

```
,0?m 360 0
```

For example, the following command displays the first 20 entries separating the i-node number and file name with a tab:

```
0,20?ut14cn
```

You can change the second number, 20, to specify the number of entries in the directory. If you place the following command at the beginning of the script, the **adb** program displays the strings as headings for each column of numbers:

```
= "i number" 8t "Name"
```

Once you have created the script file, redirect it as input when you start the **adb** program with the name of your directory. For example, the following command starts the **adb** program on the `geo` directory using command input from the `ddump` script file:


```
adb geo - <ddump
```

The minus sign (-) prevents the **adb** program from opening a core file. The **adb** program reads the commands from the script file.

To display the i-node table of a file system, create a new script and then start the **adb** program with the file name of the device associated with the file system. The i-node table of a file system has a complex structure. Each entry contains:

- A word value for status flags
- A byte value for number links
- 2-byte values for the user and group IDs
- A byte and word value for the size
- 8-word values for the location on disk of the file's blocks
- 2-word values for the creation and modification dates

The following is an example directory dump output:

```
0:      inumber Name
        26      .
        2       ..
        27     .estate
        28     adbsamp
        29     adbsamp.c
        30     calc.lex
        31     calc.yacc
        32     cbtest
        68     .profile
        66     .profile.bak
        46     adbsamp2.c
        52     adbsamp2
        35     adbsamp.s
        34     adbsamp2.s
        48     forktst1.c
        49     forktst2.c
        50     forktst3.c
        51     lpp&us1.name
        33     adbsamp3.c
        241    sample
        198    adbsamp3
        55     msgqtst.c
        56     newsig.c
```

The i-node table starts at the address 02000. You can display the first entry by putting the following command in your script file:

```
02000, -1?on3bnbrdn8un2Y2na
```

The command specifies several new-line characters for the output display to make it easier to read.

To use the script file with the i-node table of the **/dev/hd3** file, enter the following command:

```
adb /dev/hd3 - <script
```

Each entry in the display has the form:

```
02000: 073145
        0163 0164 0141
        0162 10356
        28770 8236 25956 27766 25455 8236 25956 25206
        1976 Feb 5 08:34:56 1975 Dec 28 10:55:15
```

Example of data formatting in adb debugging

To display the current address after each machine instruction, enter:

```
main , 5 ? ia
```

This produces output such as the following when used with the example program **adbsamp**:

```
.main:
.main:      mflr 0
.main+4:    st r0, 0x8(r1)
.main+8:    stu rs, (r1)
.main+c:    lil r4, 0x1
.main+10:   oril r3, r4, 0x0
.main+14:
```

To make it clearer that the current address does not belong to the instruction that appears on the same line, add the new-line format character (n) to the command:

```
.main , 5 ? ian
```

In addition, you can put a number before a formatting character to indicate the number of times to repeat that format.

To print a listing of instructions and include addresses after every fourth instruction, use the following command:

```
.main,3?4ian
```

This instruction produces the following output when used with the example program **adbsamp**:

```
.main:
      mflr 0
      st r0, 0x8(r1)
      stu r1, -56(r1)
      lil r4, 0x1

.main+10:
      oril r3, r4, 0x0
      bl .f
      l r0, 0x40(r1)
      ai r1, r1, 0x38

.main+20:
      mtlr r0
      br
      Invalid opcode
      Invalid opcode

.main+30:
```

Be careful where you put the number.

The following command, though similar to the previous command, does not produce the same output:

```
main,3?i4an
```

```
.main:
.main:      mflr 0
.main+4:    .main+4:      .main+4:      .main+4:
      st r0, 0x8(r1)
.main+8:    .main+8:      .main+8:      .main+8:
      stu r1, (r1)
.main+c:    .main+c:      .main+c:      .main+c:
```

You can combine format requests to provide elaborate displays. For example, entering the following command displays instruction mnemonics followed by their hexadecimal equivalent:

```
.main,-1?i^xn
```

In this example, the display starts at the address `main`. The negative count (-1) causes an indefinite call of the command, so that the display continues until an error condition (such as end-of-file) occurs. In the format, `i` displays the mnemonic instruction at that location, the `^` (caret) moves the current address back to the beginning of the instruction, and `x` re-displays the instruction as a hexadecimal number. Finally, `n` sends a newline character to the terminal. The output is similar to the following, only longer:

```
.main:
.main:      mflr 0
            7c0802a6
            st r0, 0x8(r1)
            9001008
            st r1, -56(r1)
            9421ffc8
            lil r4, 0x1
            38800001
            oril r3, r4, 0x0
            60830000
            bl - .f
            4bffff71
            l r0, 0x40(r1)
            80010040
            ai r1, r1, 0x38
            30210038
            mtlr r0
            7c0803a6
```

The following example shows how to combine formats in the `? or /` subcommand to display different types of values when stored together in the same program. It uses the **adbsamp** program. For the commands to have variables with which to work, you must first set a breakpoint to stop the program, and then run the program until it finds the breakpoint. Use the `:b` command to set a breakpoint:

```
.main+4:b
```

Use the `$b` command to show that the breakpoint is set:

```
$b
breakpoints
count bkpt      command
1      .main+4
```

Run the program until it finds the breakpoint by entering:

```
:r
adbsamp: running
breakpoint .main+4:      st r0, 0x8(r1)
```

You can now display conditions of the program when it stopped. To display the value of each individual variable, give its name and corresponding format in a `/` (slash) command. For example, the following command displays the contents of `str1` as a string:

```
str1/s
str1:
str1:      This is a character string
```

The following command displays the contents of `number` as a decimal integer:

```
number/D
number:
number:      456
```

You can choose to view a variable in a variety of formats. For example, you can display the long variable **lnum** as a 4-byte decimal, octal, and hexadecimal number by entering the commands:

```
lnum/D
lnum:
lnum:      1234

lnum/O
lnum:
```

```

lnum:      2322

lnum/X
lnum:
lnum:      4d2

```

You can also examine variables in other formats. For example, the following command displays some variables as eight hexadecimal values on a line and continues for five lines:

```

str1,5/8x
str1:
str1:  5468 6973 2069 7320 6120 6368 6172 6163
      7465 7220 7374 7269 6e67 0 0 0 0

number: 0          1c8 0 0 0 4d2 0 0
        3fa0 0 0 0 5468 6973 2069 7320
        7468 6520 7365 636f 6e64 2063 6861 7261

```

Since the data contains a combination of numeric and string values, display each value as both a number and a character to see where the actual strings are located. You can do this with one command:

```

str1,5/4x4^8Cn
str1:
str1:  5468 6973 2069 7320 This is
      6120 6368 6172 6163 a charac
      7465 7220 7374 7269 ter stri
      6e67 0 0 0 ng~@~@~@~@~@
      0 1c8 0 0 ~@~@~A~<c8>~@~@~@~@

```

In this case, the command displays four values in hexadecimal, then displays the same values as eight ASCII characters. The ^ (caret) is used four times just before displaying the characters to set the current address back to the starting address for that line.

To make the display easier to read, you can insert a tab between the values and characters and give an address for each line:

```

str1,5/4x4^8t8Cna
str1:
str1:  5468 6973 2069 7320 This is
str1+8: 6120 6368 6172 6163 a charac
str1+10: 7465 7220 7374 7269 ter stri
str1+18: 6e67 0 0 1 ng~@~@~@~@~A

number:
number: 0          1c8 0 0 ~@~@~A~<c8>~@~@~@~@
fpt:

```

Example of tracing multiple functions in adb debugging

The following example shows how to execute a program under **adb** control and carry out the basic debugging operations described in the following sections.

Note: The example program used in this section, **adbsamp3**, contains an infinite recursion of subfunction calls. If you run this program to completion, it causes a memory fault error and quits.

The source program for this example is stored in a file named **adbsamp3.c**. Compile this program to an executable file named **adbsamp3** using the **cc** command:

```
cc adbsamp3.c -o adbsamp3
```

Starting the adb program

To start the session and open the program file, use the following command (no core file is used):

```
adb adbsamp3
```

Setting breakpoints

First, set breakpoints at the beginning of each function using the **:b** subcommand:

```
.f:b
.g:b
.h:b
```

Displaying a set of instructions

Next, display the first five instructions in the **f** function:

```
.f,5?ia
.f:
.f:          mflr  r0
.f+4:       st   r0, 0x8(r1)
.f+8:       stu  r1, -64(r1)
.f+c:       st   r3, 0x58(r1)
.f+10:      st   r4, 0x5c(r1)
.f+14:
```

Display five instructions in function **g** without their addresses:

```
.g,5?i
.g:         mflr  r0
           st   r0, 0x8(r1)
           stu  r1, -64(r1)
           st   r3, 0x58(r1)
           st   r4, 0x5c(r1)
```

Starting the adsamp3 program

Start the program by entering the following command:

```
:r
adbsamp3: running
breakpoint .f:          mflr  r0
```

The **adb** program runs the sample program until it reaches the first breakpoint where it stops.

Removing a breakpoint

Since running the program to this point causes no errors, you can remove the first breakpoint:

```
.f:d
```

Continuing the program

Use the **:c** subcommand to continue the program:

```
:c
adbsamp3: running
breakpoint .g:          mflr  r0
```

The **adb** program restarts the **adbsamp3** program at the next instruction. The program operation continues until the next breakpoint, where it stops.

Tracing the path of execution

Trace the path of execution by entering:

Displaying a variable value

Display the contents of the **fcnt** integer variable by entering the command:

```
fcnt/D
fcnt:
fcnt:    1
```

Skipping breakpoints

Next, continue running the program and skip the first 10 breakpoints by entering:

```
,10:c
adbsamp3: running
breakpoint .g:      mflr r0
```

The **adb** program starts the **adbsamp3** program and displays the running message again. It does not stop the program until exactly 10 breakpoints have been encountered. To ensure that these breakpoints have been skipped, display the backtrace again:

```
$c
.g(0,0) .f+2a
.f(2,11) .h+28
.h(10,f) .g+2a
.g(11,20) .f+2a
.f(2,f) .h+28
.h(e,d) .g+2a
.g(f,1c) .f+2a
.f(2,d) .h+28
.h(c,b) .g+2a
.g(d,18) .f+2a
.f(2,b) .h+28
.h(a,9) .g+2a
.g(b,14) .f+2a
.f(2,9) .h+28
.h(8,7) .g+2a
.g(9,10) .f+2a
.f(2,7) .h+28
.h(6,5) .g+2a
.g(7,c) .f+2ae
.f(2,5) .h+28
.h(4,3) .g+2a
.g(5,8) .f+2a
.f(2,3) .h+28
.h(2,1) .g+2a
.g(2,3) .f+2a
.f(1,1) .main+e
.main(0,0,0) start+fa
```

dbx symbolic debug program overview

The **dbx** symbolic debug program allows you to debug an application program at two levels: the source-level and the assembler language-level. Source level debugging allows you to debug your C, C++, or FORTRAN language program.

Assembler language level debugging allows you to debug executable programs at the machine level. The commands used for machine level debugging are similar to those used for source-level debugging.

Using the **dbx** debug program, you can step through the application program you want to debug one line at a time or set breakpoints in the object program that will stop the debug program. You can also search through and display portions of the source files for an application program.

The following sections contain information on how to perform a variety of tasks with the **dbx** debug program:

Using the dbx debug program

This section contains information on how to use the **dbx** debug program.

Starting the dbx debug program

The **dbx** program can be started with a variety of flags. The three most common ways to start a debug session with the **dbx** program are:

- Running the **dbx** command on a specified object file
- Using the **-r** flag to run the **dbx** command on a program that ends abnormally
- Using the **-a** flag to run the **dbx** command on a process that is already in progress

When the **dbx** command is started, it checks for a **.dbxinit** file in the user's current directory and in the user's **\$HOME** directory. If a **.dbxinit** file exists, its subcommands run at the beginning of the debug

session. If a **.dbxinit** file exists in both the home and current directories, then both are read in that order. Because the current directory **.dbxinit** file is read last, its subcommands can supercede those in the home directory.

If no object file is specified, then the **dbx** program asks for the name of the object file to be examined. The default is **a.out**. If the **core** file exists in the current directory or a *CoreFile* parameter is specified, then the **dbx** program reports the location where the program faulted. Variables, registers, and memory held in the core image may be examined until execution of the object file begins. At that point the **dbx** debug program prompts for commands.

Debugging a core image with missing dependent modules

Beginning with AIX 5.3, the **dbx** program has the ability to examine a core image even if one or more dependent modules are inaccessible. During initialization, notification messages are displayed for each missing dependent module.

In normal operation, the **dbx** program relies on the information contained in the dependent modules' symbol tables and text sections. Because some of that information is missing, a **dbx** session with missing dependent modules has the following limitations:

- All attempts by the user to read the contents of memory residing in the text sections of the missing dependent modules result in an error message. The error message is like the error that occurs when data cannot be accessed because it does not reside in the core file.
- The user cannot obtain information concerning any symbols that would have been read from the symbol tables of the missing dependent modules. The behavior of the **dbx** program is similar to the case where a dependent module's symbol table was stripped.
- Stack frames corresponding to routines within the missing dependent modules are displayed simply as:

```
. ()
```

In addition, the instruction address within the unknown routine and the name of the corresponding missing dependent module are displayed.

The user always has the option of directing the **dbx** program to accessible dependent modules using the **-p** flag.

Debugging a core image with mismatched dependent modules

Beginning with AIX 5.3, the **dbx** program detects if any dependent modules referenced in the core file are different than at core file creation. During initialization, notification messages are displayed for each mismatched dependent module.

The user should be aware that any information displayed by the **dbx** program that is based on the contents of a mismatched dependent module might be unreliable. In an effort to alert the user to information that should not be trusted, the **dbx** program sends notification messages whenever questionable information is displayed.

To disable this function and force the **dbx** program to treat mismatched dependent modules as missing dependent modules, the user can export the *DBX_MISMATCH_MODULE* environment variable with a value of **DISCARD**. With this variable exported, the **dbx** program still notifies the user of the mismatch, but proceeds as if mismatched dependent modules were inaccessible.

The user always has the option of directing the **dbx** program to matching dependent modules using the **-p** flag.

Running shell commands from **dbx**

You can run shell commands without exiting from the debug program using the **sh** subcommand.

If **sh** is entered without any commands specified, the shell is entered for use until it is exited, at which time control returns to the **dbx** program.

Command line editing in **dbx**

The **dbx** command provides command line editing features similar to those provided by **Korn Shell**. **vi** mode provides **vi**-like editing features, while **emacs** mode gives you controls similar to **emacs**.

You can turn these features on by using **dbx** subcommand **set -o** or **set edit**. So, to turn on **vi**-style command line editing, you would type the subcommand `set edit vi` or `set -o vi`.

You can also use the **EDITOR** environment variable to set the editing mode.

The **dbx** command saves commands entered to history file **.dbxhistory**. If the **DBXHISTFILE** environment variable is not set, then the history file used is **\$HOME/.dbxhistory**.

By default, the **dbx** command saves the text of the last 128 commands entered. The **DBXHISTSIZE** environment variable can be used to increase this limit.

Using program control

The **dbx** debug program allows you to set breakpoints (stopping places) in the program. After entering the **dbx** program you can specify which lines or addresses are to be breakpoints and then run the program you want to debug with the **dbx** program. The program halts and reports when it reaches a breakpoint. You can then use **dbx** commands to examine the state of your program.

An alternative to setting breakpoints is to run your program one line or instruction at a time, a procedure known as single-stepping.

Setting and deleting breakpoints

Use the **stop** subcommand to set breakpoints in the **dbx** program. The **stop** subcommand halts the application program when certain conditions are fulfilled:

- The *Variable* is changed when the *Variable* parameter is specified.
- The *Condition* is true when the **if** *Condition* flag is used.
- The *Procedure* is called when the **in** *Procedure* flag is used.
- The *SourceLine* line number is reached when the **at** *SourceLine* flag is used.

Note: The *SourceLine* variable can be specified as an integer or as a file name string followed by a **:** (colon) and an integer.

After any of these commands, the **dbx** program responds with a message reporting the event ID associated with your breakpoint along with an interpretation of your command. You can associate **dbx** subcommands to the specified event ID by using the **addcmd** subcommand. These associated **dbx** subcommands are issued when the breakpoint, tracepoint, or watchpoint corresponding to this event is hit. Use the **delcmd** subcommand to delete the associated **dbx** subcommands from the specified event ID.

Running a program

The **run** subcommand starts your program. It tells the **dbx** program to begin running the object file, reading any arguments just as if they were typed on the shell command line. The **rerun** subcommand has the same form as **run**; the difference is that if no arguments are passed, the argument list from the previous execution is used. After your program begins, it continues until one of the following events occurs:

- The program reaches a breakpoint.
- A signal occurs that is not ignored, such as **INTERRUPT** or **QUIT**.
- A multiprocess event occurs while multiprocess debugging is enabled.
- The program performs a **load**, **unload**, or **loadbind** subroutine.

Note: The **dbx** program ignores this condition if the **\$ignoreload** debug variable is set. This is the default. For more information see the **set** subcommand.

- The program completes.

In each case, the **dbx** debug program receives control and displays a message explaining why the program stopped.

There are several ways to continue the program once it stops:

Command	Description
cont	Continues the program from where it stopped.
detach	Continues the program from where it stopped, exiting the debug program. This is useful after you have patched the program and want to continue without the debug program.
return	Continues execution until a return to <i>Procedure</i> is encountered, or until the current procedure returns if <i>Procedure</i> is not specified.
skip	Continues execution until the end of the program or until <i>Number</i> + 1 breakpoints execute.
step	Runs one or a specified <i>Number</i> of source lines.
next	Runs up to the next source line, or runs a specified <i>Number</i> of source lines.

A common method of debugging is to step through your program one line at a time. The **step** and **next** subcommands serve that purpose. The distinction between these two commands is apparent only when the next source line to be run involves a call to a subprogram. In this case, the **step** subcommand stops in the subprogram; the **next** subcommand runs until the subprogram has finished and then stops at the next instruction after the call.

The **\$stepignore** debug variable can be used to modify the behavior of the **step** subcommand.

There is no event number associated with these stops because there is no permanent event associated with stopping a program.

If your program has multiple threads, they all run normally during the **cont**, **next**, **nexti**, and **step** subcommands. These commands act on the running thread (the thread which stopped execution by hitting a breakpoint), so even if another thread runs the code which is being stepped, the **cont**, **next**, **nexti**, or **step** operation continues until the running thread has also executed that code.

If you want these subcommands to execute the running thread only, you can set the **dbx** debug program variable **\$hold_next**; this causes the **dbx** debug program to hold all other user threads during **cont**, **next**, **nexti**, and **step** subcommands.

Note: If you use this feature, remember that a held thread will not be able to release any locks which it has acquired; another thread which requires one of these locks could deadlock your program.

Separating dbx output from program output

Use the **screen** subcommand for debugging programs that are screen-oriented, such as text editors or graphics programs. This subcommand opens an Xwindow for **dbx** command interaction. The program continues to operate in the window in which it originated. If **screen** is not used, **dbx** program output is intermixed with the screen-oriented program output.

Tracing execution

The **trace** subcommand tells the **dbx** program to print information about the state of the program being debugged while that program is running. The **trace** subcommand can slow a program considerably, depending on how much work the **dbx** program has to do. There are five forms of program tracing:

- You can single-step the program, printing out each source line that is executed. The **\$stepignore** debug variable can be used to modify the behavior of the **trace** subcommand. See the **set** subcommand for more information.
- You can restrict the printing of source lines to when the specified procedure is active. You can also specify an optional condition to control when trace information is produced.
- You can display a message each time a procedure is called or returned.
- You can print the specified source line when the program reaches that line.
- You can print the value of an expression when the program reaches the specified source line.

Deleting trace events is the same as deleting stop events. When the **trace** subcommand is executed, the event ID associated is displayed along with the internal representation of the event.

Displaying and manipulating the source file with the dbx debug program

This section describes the process of displaying and manipulating the source file with the dbx debug program.

You can use the **dbx** debug program to search through and display portions of the source files for a program.

You do not need a current source listing for the search. The **dbx** debug program keeps track of the current file, current procedure, and current line. If a core file exists, the current line and current file are set initially to the line and file containing the source statement where the process ended. This is only true if the process stopped in a location compiled for debugging.

Changing the source directory path

By default, the **dbx** debug program searches for the source file of the program being debugged in the following directories:

- Directory where the source file was located when it was compiled. This directory is searched only if the compiler placed the source path in the object.
- Current directory.
- Directory where the program is currently located.

You can change the list of directories to be searched by using the **-I** option on the **dbx** invocation line or issuing the **use** subcommand within the **dbx** program. For example, if you moved the source file to a new location since compilation time, you might want to use one of these commands to specify the old location, the new location, and some temporary location.

Displaying the current file

The **list** subcommand allows you to list source lines.

The \$ (dollar sign) and @ (at sign) symbols represent *SourceLineExpression* and are useful with the **list**, **stop**, and **trace** subcommands. The \$ symbol represents the next line to be run. The @ symbol represents the next line to be listed.

The **move** subcommand changes the next line number to be listed.

Changing the current file or procedure

Use the **func** and **file** subcommands to change the current file, current procedure, and current line within the **dbx** program without having to run any part of your program.

Search through the current file for text that matches regular expressions. If a match is found, the current line is set to the line containing the matching text. The syntax of the search subcommand is:

/ RegularExpression [/]

Searches forward in the current source file for the given expression.

? RegularExpression [?]

Searches backward in the current source file for the given expression.

If you repeat the search without arguments, the **dbx** command searches again for the previous regular expression. The search wraps around the end or beginning of the file.

You can also invoke an external text editor for your source file using the **edit** subcommand. You can override the default editor (**vi**) by setting the **EDITOR** environment variable to your desired editor before starting the **dbx** program.

The **dbx** program resumes control of the process when the editing session is completed.

Debugging programs involving multiple threads

Programs involving multiple user threads call the subroutine **pthread_create**. When a process calls this subroutine, the operating system creates a new thread of execution within the process. When debugging a multithreaded program, it is necessary to work with individual threads instead of with processes. The **dbx**

program only works with user threads: in the **dbx** documentation, the word *thread* is usually used alone to mean *user thread*. The **dbx** program assigns a unique thread number to each thread in the process being debugged, and also supports the concept of a running and current thread:

Running thread

The user thread that was responsible for stopping the program by hitting a breakpoint. Subcommands that single-step the program work with the running thread.

Current thread

The user thread that you are examining. Subcommands that display information work in the context of the current thread.

By default, the running thread and current thread are the same. You can select a different current thread by using the **thread** subcommand. When the **thread** subcommand displays threads, the current thread line is preceded by a **>**. If the running thread is not the same as the current thread, its line is preceded by a *****.

Debugging programs involving multiple processes

Programs involving multiple processes call the **fork** and **exec** subroutines. When a program forks, the operating system creates another process that has the same image as the original. The original process is called the parent process, the created process is called the child process.

When a process performs an **exec** subroutine, a new program takes over the original process. Under normal circumstances, the debug program debugs only the parent process. However, the **dbx** program can follow the execution and debug the new processes when you issue the **multproc** subcommand. The **multproc** subcommand enables multiprocess debugging.

When multiprocess debugging is enabled and a fork occurs, the parent and child processes are halted. A separate virtual terminal Xwindow is opened for a new version of the **dbx** program to control running of the child process:

```
(dbx) multproc on
(dbx) multproc
multi-process debugging is enabled
(dbx) run
```

When the fork occurs, execution is stopped in the parent, and the **dbx** program displays the state of the program:

```
application forked, child pid = 422, process stopped, awaiting input
stopped due to fork with multiprocessing enabled in fork at 0x1000025a (fork+0xe)
(dbx)
```

Another virtual terminal Xwindow is then opened to debug the child process:

```
debugging child, pid=422, process stopped, awaiting input
stopped due to fork with multiprocessing enabled in fork at 0x10000250
10000250 (fork+0x4) )80010010 1 r0,0x10(r1)
(dbx)
```

At this point, two distinct debugging sessions are running. The debugging session for the child process retains all the breakpoints from the parent process, but only the parent process can be rerun.

When a program performs an **exec** subroutine in multiprocess debugging mode, the program overwrites itself, and the original symbol information becomes obsolete. All breakpoints are deleted when the **exec** subroutine runs; the new program is stopped and identified for the debugging to be meaningful. The **dbx** program attaches itself to the new program image, makes a subroutine to determine the name of the new program, reports the name, and then prompts for input. The prompt is similar to the following:

```
(dbx) multproc
Multi-process debugging is enabled
(dbx) run
Attaching to program from exec . . .
Determining program name . . .
Successfully attached to /home/user/execprog . . .
```

```
Reading symbolic information . . .
(dbx)
```

If a multithreaded program forks, the new child process will have only one thread. The process should call the **exec** subroutine. Otherwise, the original symbol information is retained, and thread-related subcommands (such as **thread**) display the objects of the parent process, which are obsolete. If an **exec** subroutine is called, the original symbol information is reinitialized, and the thread-related subcommands display the objects in the new child process.

It is possible to follow the child process of a fork without a new Xwindow being opened by using the **child** flag of the **multproc** subcommand. When a forked process is created, **dbx** follows the child process. The **parent** flag of the **multproc** subcommand causes **dbx** to stop when a program forks, but then follows the parent. Both the **child** and **parent** flags follow an **execed** process. These flags are very useful for debugging programs when Xwindows is not running.

Examining program data

This section explains how to examine, test, and modify program data.

Handling signals

The **dbx** debug program can either trap or ignore signals before they are sent to your program. Each time your program is to receive a signal, the **dbx** program is notified. If the signal is to be ignored, it is passed to your program; otherwise, the **dbx** program stops the program and notifies you that a signal has been trapped. The **dbx** program cannot ignore the **SIGTRAP** signal if it comes from a process outside of the debug process. In a multithreaded program, a signal can be sent to a particular thread via the **pthread_kill** subroutine. By default, the **dbx** program stops and notifies you that a signal has been trapped. If you request a signal be passed on to your program using the **ignore** subcommand, the **dbx** program ignores the signal and passes it on to the thread. Use the **catch** and **ignore** subcommands to change the default handling.

In the following example, a program uses **SIGGRANT** and **SIGREQUEST** to handle allocation of resources. In order for the **dbx** program to continue each time one of these signals is received, enter:

```
(dbx) ignore GRANT
(dbx) ignore SIGREQUEST
(dbx) ignore
CONT CLD ALARM KILL GRANT REQUEST
```

The **dbx** debug program can block signals to your program if you set the **\$sigblock** variable. By default, signals received through the **dbx** program are sent to the source program or the object file specified by the **dbx ObjectFile** parameter. If the **\$sigblock** variable is set using the **set** subcommand, signals received by the **dbx** program are not passed to the source program. If you want a signal to be sent to the program, use the **cont** subcommand and supply the signal as an operand.

You can use this feature to interrupt execution of a program running under the **dbx** debug program. Program status can be examined before continuing execution as usual. If the **\$sigblock** variable is not set, interrupting execution causes a **SIGINT** signal to be sent to the program. This causes execution, when continued, to branch to a signal handler if one exists.

The following example program illustrates how execution using the **dbx** debug program changes when the **\$sigblock** variable is set:

```
#include <signal.h>
#include <stdio.h>
void inthand( ) {
    printf("\nSIGINT received\n");
    exit(0);
}

main( )
{
    signal(SIGINT, inthand);
    while (1) {
        printf(".");
        fflush(stdout);
    }
}
```

```

        }
        sleep(1);
    }
}

```

The following sample session with the **dbx** program uses the preceding program as the source file. In the first run of the program, the **\$sigblock** variable is not set. During rerun, the **\$sigblock** variable is set. Comments are placed between angle brackets to the right:

```

dbx version 3.1.
Type 'help' for help.
reading symbolic information ...
(dbx) run
.....^C <User pressed Ctrl-C here!>
interrupt in sleep at 0xd00180bc
0xd00180bc (sleep+0x40) 80410014      1      r2,0x14(r1)
(dbx) cont

```

SIGINT received

```

execution completed
(dbx) set $sigblock
(dbx) rerun
[ looper ]
.....^C <User pressed Ctrl-C here!>
interrupt in sleep at 0xd00180bc
0xd00180bc (sleep+0x40) 80410014      1      r2,0x14(r1)
(dbx) cont
....^C <Program did not receive signal, execution continued>

```

```

interrupt in sleep at 0xd00180bc
0xd00180bc (sleep+0x40) 80410014      1      r2,0x14(r1)
(dbx) cont 2 <End program with a signal 2>

```

SIGINT received

```

execution completed
(dbx)

```

Calling procedures

You can call your program procedures from the **dbx** program to test different arguments. You can also call diagnostic routines that format data to aid in debugging. Use the **call** subcommand or the **print** subcommand to call a procedure.

Displaying a stack trace

To list the procedure calls preceding a program halt, use the **where** command.

In the following example, the executable object file, **hello**, consists of two source files and three procedures, including the standard procedure **main**. The program stopped at a breakpoint in procedure **sub2**.

```

(dbx) run
[1] stopped in sub2 at line 4 in file "hellosub.c"
(dbx) where
sub2(s = "hello", n = 52), line 4 in "hellosub.c"
sub(s = "hello", a = -1, k = delete), line 31 in "hello.c"
main(), line 19 in "hello.c"

```

The stack trace shows the calls in reverse order. Starting at the bottom, the following events occurred:

1. Shell called **main**.
2. **main** called sub procedure at line 19 with values **s = "hello"**, **a = -1**, and **k = delete**.
3. sub called **sub2** procedure at line 31 with values **s = "hello"** and **n = 52**.
4. The program stopped in **sub2** procedure at line 4.

A portion of the stack trace from frame number 0 to frame number 1 can be displayed using `where 0 1`.

```
(dbx)run
[1] stopped in sub2 at line 4 in file "hellosub.c"
(dbx) where 0 1
sub2(s = "hello", n = 52), line 4 in "hellosub.c"
sub(s = "hello", a = -1, k = delete), line 31 in "hello.c"
```

Note: Set the debug program variable `$noargs` to turn off the display of arguments passed to procedures. Set the debug program variable `$stack_details` to display the frame number and the register set for each active function or procedure.

You can also display portions of the stack with the `up`, `down`, and `frame` subcommands.

Displaying and modifying variables

To display an expression, use the `print` subcommand. To print the names and values of variables, use the `dump` subcommand. If the given procedure is a period, then all active variables are printed. If the `PATTERN` parameter is specified, then instead of displaying only the specified symbol, all the symbols that match the `PATTERN` are printed. To modify the value of a variable, use the `assign` subcommand.

In the following example, a C program has an automatic integer variable `x` with value 7, and `s` and `n` parameters in the `sub2` procedure:

```
(dbx) print x, n
7 52
(dbx) assign x = 3*x
(dbx) print x
21
(dbx) dump
sub2(s = "hello", n = 52)
x = 21
```

Displaying thread-related information

To display information on user threads, mutexes, conditions, and attribute objects, use the `thread`, `mutex`, `condition`, and `attribute` subcommands. You can also use the `print` subcommand on these objects. In the following example, the running thread is thread 1. The user sets the current thread to be thread 2, lists the threads, prints information on thread 1, and finally prints information on several thread-related objects.

```
(dbx) thread current 2
(dbx) thread
  thread state-k   wchan state-u   k-tid mode held scope function
*$t1    run                running 12755  u  no  pro  main
>$t2    run                running 12501  k  no  sys  thread_1
```

```
(dbx) print $t1
(thread_id = 0x1, state = run, state_u = 0x0, tid = 0x31d3, mode = 0x1, held = 0x0, priority = 0x3c,
 policy = other, scout = 0x1, cursig = 0x5, attributes = 0x200050f8)
```

```
(dbx) print $a1,$c1,$m2
(attr_id = 0x1, type = 0x1, state = 0x1, stacksize = 0x0, detachedstate = 0x0, process_shared = 0x0,
 contentionscope = 0x0, priority = 0x0, sched = 0x0, inherit = 0x0, protocol = 0x0, prio_ceiling = 0x0)
(cv_id = 0x1, lock = 0x0, semaphore_queue = 0x200032a0, attributes = 0x20003628)
(mutex_id = 0x2, islock = 0x0, owner = (nil), flags = 0x1, attributes = 0x200035c8)
```

Scoping of names

Names resolve first using the static scope of the current function. The dynamic scope is used if the name is not defined in the first scope. If static and dynamic searches do not yield a result, an arbitrary symbol is chosen and the message using `QualifiedName` is printed. You can override the name resolution procedure by qualifying an identifier with a block name (such as `Module.Variable`). Source files are treated as modules named by the file name without the suffix. For example, the `x` variable, which is declared in the `sub` procedure inside the `hello.c` file, has the fully qualified name `hello.sub.x`. The program itself has a period for a name.

The **which** and **whereis** subcommands can be helpful in determining which symbol is found when multiple symbols with the same name exist.

Using operators and modifiers in expressions

The **dbx** program can display a wide range of expressions. Specify expressions with a common subset of C syntax, with some FORTRAN extensions.

* (asterisk) or ^ (caret)

Denotes indirection or pointer dereferencing.

[] (brackets) or () (parentheses)

Denotes subscript array expressions.

. (period)

Use this field reference operator with pointers and structures. This makes the C operator -> (arrow) unnecessary, although it is allowed.

& (ampersand)

Gets the address of a variable.

.. (two periods)

Separates the upper and lower bounds when specifying a subsection of an array. For example: **n[1..4]**.

The following types of operations are valid in expressions:

Algebraic

=, -, *, / (floating division), **div** (integral division), **mod**, **exp** (exponentiation)

Bitwise

~, **I**, **bitand**, **xor**, ~, <<, >>

Logical

or, and, not, **!!**, **&&**

Comparison

<, >, <=, >=, <> or !=, = or ==

Other

sizeof

Logical and comparison expressions are allowed as conditions in **stop** and **trace** subcommands.

Checking of expression types

The **dbx** debug program checks expression types. You can override the expression type by using a renaming or casting operator. There are three forms of type renaming:

- *Typename (Expression)*
- *Expression \ Typename*
- *(Typename) Expression*

Note: When you cast to or from a structure, union, or class, the casting is left-justified. However, when casting from a class to a base class, C++ syntax rules are followed.

For example, to rename the x variable where x is an integer with a value of 97, enter:

```
(dbx) print char (x), x \ char, (char) x, x,  
'a' 'a' 'a' 97
```

The following examples show how you can use the *(Typename) Expression* form of type renaming:

```
print (float) i
```

```
print ((struct qq *) void_pointer)->first_element
```

The following restrictions apply to C-style typecasting for the **dbx** debug program:

- The Fortran types (integer*1, integer*2, integer*4, logical*1, logical*2, logical*4, and so on) are not supported as cast operators.
- If an active variable has the same name as one of the base types or user-defined types, the type cannot be used as a cast operator for C-style typecasting.

The **whatis** subcommand prints the declaration of an identifier, which you can then qualify with block names.

Use the **\$\$TagName** construct to print the declaration of an enumeration, structure, or union tag.

The type of the **assign** subcommand expression must match the variable type you assigned. If the types do not match, an error message is displayed. Change the expression type using a type renaming. Disable type checking by setting a special **dbx** debug program **\$unsafeassign** variable.

Folding variables to lowercase and uppercase

By default, the **dbx** program folds symbols based on the current language. If the current language is C, C++, or undefined, the symbols are not folded. If the current language is Fortran, the symbols are folded to lowercase. The current language is undefined if the program is in a section of code that has not been compiled with the **debug** flag. You can override default handling with the **case** subcommand.

Using the **case** subcommand without arguments displays the current case mode.

The Fortran compiler converts all program symbols to lowercase; the C compiler does not. However, some Fortran compilers might not always generate lowercase symbols. For example, given a procedure named **proc1** inside a module named **mod2**, the XLF Fortran compiler generates the **__mod2_MOD_proc1** symbol, which is mixed case. In such situations, you must change the **case** in the **dbx** program to **mixed** case.

Changing print output with special debug program variables

Use the **set** subcommand to set the following special **dbx** debug program variables to get different results from the **print** subcommand:

\$hexints

Prints integer expressions in hexadecimal.

\$hexchars

Prints character expressions in hexadecimal.

\$hexstrings

Prints the address of the character string, not the string itself.

\$octints

Prints integer expressions in octal.

\$expandunions

Prints fields within a union.

\$pretty

Displays complex C and C++ types in **pretty** format.

\$print_dynamic

Prints the dynamic type of the C++ objects.

\$show_vft

Prints Virtual Function Table while printing C++ objects.

Set and unset the debug program variables to get the desired results. For example:

```
(dbx) whatis x; whatis i; whatis s
int x;
char i;
char *s;
(dbx) print x, i, s
375 'c' "hello"
(dbx) set $hexstrings; set $hexints; set $hexchars
(dbx) print x, i, s
0x177 0x63 0x3fffe460
(dbx) unset $hexchars; set $octints
```



```

(dbx) print x, i
0567 'c'
(dbx) whatis p
struct info p;
(dbx) whatis struct info
struct info {
    int x;
    double position[3];
    unsigned char c;
    struct vector force;
};
(dbx) whatis struct vector
struct vector {
    int a;
    int b;
    int c;
};
(dbx) print p
(x = 4, position = (1.3262493258532527e-315, 0.0, 0.0),
c = '\0', force = (a = 0, b = 9, c = 1))(dbx) set $pretty="on"
(dbx) print p
{
    x = 4
    position[0] = 1.3262493258532527e-315
    position[1] = 0.0
    position[2] = 0.0
    c = '\0'
    force = {
        a = 0
        b = 9
        c = 1
    }
}
(dbx) set $pretty="verbose"
(dbx) print p
x = 4
position[0] = 1.3262493258532527e-315
position[1] = 0.0
position[2] = 0.0
c = '\0'
force.a = 0
force.b = 9
force.c = 1

```

When `show_vft` is not set and an object is printed using `print` sub command, the Virtual Function Table (VFT) is not printed. If it is set VFT is displayed. For example:

```

(dbx) p *d
B1:(int_in_b1 = 91)
B2:(int_in_b2 = 92)
(int_in_d = 93)
(dbx) p *b2
(int_in_b2 = 20)
(dbx)set $show_vft
(dbx) p *d
B1:(B1::f1(), int_in_b1 = 91)
B2:(D::f2(), int_in_b2 = 92)
(int_in_d = 93)
(dbx) p *b2
(B2::f2(), int_in_b2 = 20)
(dbx)

```

When `print_dynamic` is not set the object is displayed as per the static type's (what is defined in the source code) template. Otherwise it will be displayed as per the dynamic type's (what an object was before any casts were made to it) template. For example:

```

(dbx) r
[1] stopped in main at line 57
57 A *obj1 = new A();
(dbx) n
stopped in main at line 58
58 A *obj2 = new B();
(dbx) n
stopped in main at line 59
59 cout<<" a = "<<obj2->a<<" b = "<<obj2->b<<endl;
(dbx) p *obj2
(a = 1, b = 2)
(dbx)set $print_dynamic

```

```
(dbx) print *obj2
      A:(a = 1, b = 2)
      (c = 3, d = 4)
(dbx)
```

Debugging at the machine level with dbx

You can use the **dbx** debug program to examine programs at the assembly language level. You can display and modify memory addresses, display assembler instructions, single-step instructions, set breakpoints and trace events at memory addresses, and display the registers.

In the commands and examples that follow, an address is an expression that evaluates to a memory address. The most common forms of addresses are integers and expressions that take the address of an identifier with the **&** (ampersand) operator. You can also specify an address as an expression enclosed in parentheses in machine-level commands. Addresses can be composed of other addresses and the operators **+** (plus), **-** (minus), and indirection (unary *****).

The following sections contain more information on debugging at the machine level with the **dbx** program.

Using machine registers

Use the **registers** subcommand to see the values of the machine registers. Registers are divided into three groups: general-purpose, floating-point, and system-control.

General-purpose registers

General-purpose registers are denoted by **\$rNumber**, where *Number* represents the number of the register.

Note: The register value may be set to a hexadecimal value of 0xdeadbeef. This is an initialization value assigned to all general-purpose registers at process initialization.

Floating-point registers

Floating-point registers are denoted by **\$frNumber**, where *Number* represents the register number. Floating-point registers are not displayed by default. You can unset the **\$nofregs** debug program variable to enable the floating-point register display (**unset \$nofregs**). You can also reference floating-point registers by type when using these registers with the **print** and **assign** subcommands. **\$frNumber** defaults to the double type. **\$frNumberh** references the floating-point registers as type **_Decimal32**. **\$frNumberd** references the floating-point registers as type **_Decimal64**. The following are examples of the different types of floating-point registers:

```
(dbx) print $fr0
1.10000002
(dbx) print $fr0h
1.100001
(dbx) print $fr0d
1.10000062
(dbx) assign $fr0 = 9.876
(dbx) assign $fr0h = 9.876df
(dbx) assign $fr0d = 9.876dd
```

Vector registers

Vector registers are denoted by **\$vrNumber**, where *Number* represents the number of the register. Vector registers are not displayed by default, and are only present on processors supporting a Vector Processing Unit.

You can unset the **\$novregs** debug program variable to enable the vector register display with **unset \$novregs**. You can also reference vector registers by type when using them with the **print** and **assign**

subcommands. *\$vrNumber* defaults to a vector type of int. *\$vrNumberf* references the vector as type float. *\$vrNumbers* references the vector as type short. *\$vrNumberc* references the vector as type char.

Following are examples of the different types of vector registers:

```
(dbx) print $vr20
((1066192077, 1074161254, 1078355558, 1082340147))
(dbx) print $vr20f
((1.10000002, 2.0999999, 3.0999999, 4.0999999))
(dbx) print $vr20s
((16268, 52429, 16390, 26214, 16454, 26214, 16515, 13107))
(dbx) assign $vr20f[3] = 9.876
(dbx) print $vr20f ((1.10000002, 2.0999999, 3.0999999, 9.8760004))
```

System-control registers

Supported system-control registers are denoted by:

- The Instruction Address register, **\$iar** or **\$pc**
- The Condition Status register, **\$cr**
- The Multiplier Quotient register, **\$mq**
- The Machine State register, **\$msr**
- The Link register, **\$link**
- The Count register, **\$ctr**
- The Fixed Point Exception register, **\$xer**
- The Transaction ID register, **\$tid**
- The Floating-Point Status register, **\$fpSCR**

Examining memory addresses

Use the following command format to print the contents of memory starting at the first address and continuing up to the second address, or until the number of items specified by the *Count* variable are displayed. The *Mode* specifies how memory is to print.

Address, Address / [Mode][> File]

Address / [Count][Mode] [> File]

If the *Mode* variable is omitted, the previous mode specified is reused. The initial mode is **X**. The following modes are supported:

- b** Prints a byte in octal.
- c** Prints a byte as a character.
- D** Prints a long word in decimal.
- d** Prints a short word in decimal.
- Df** Prints a double-precision decimal float number.
- DDf** Prints a quad-precision decimal float number.
- f** Prints a single-precision floating-point number.

- g** Prints a double-precision floating-point number.
- Hf** Prints a single-precision decimal float number.
- h** Prints a byte in hexadecimal.
- i** Prints the machine instruction.
- lld** Prints an 8-byte signed decimal number.
- llo** Prints an 8-byte unsigned octal number.
- llu** Prints an 8-byte unsigned decimal number.
- llx** Prints an 8-byte unsigned hexadecimal number.
- O** Prints a long word in octal.
- o** Prints a short word in octal.
- q** Prints an extended-precision floating-point number.
- s** Prints a string of characters terminated by a null byte.
- X** Prints a long word in hexadecimal.
- x** Prints a short word in hexadecimal.

In the following example, expressions in parentheses can be used as an address:

```
(dbx) print &x
0x3fffe460
(dbx) &x/X
3fffe460: 31323300
(dbx) &x,&x+12/x
3fffe460: 3132 3300 7879 7a5a 5958 5756 003d 0032
(dbx) ($pc)/2i
100002cc (sub) 7c0802a6          mflr      r0
100002d0 (sub + 0x4)   bfc1fff8          stm       r30,-8(r1)
```

Running a program at the machine level

The commands for debugging your program at the machine-level are similar to those at the symbolic level. The **stopi** subcommand stops the machine when the address is reached, the condition is true, or the variable is changed. The **tracei** subcommands are similar to the symbolic trace commands. The **stepi** subcommand executes either one or the specified *Number* of machine instructions.

If you performed another **stepi** subcommand at this point, you would stop at address 0x10000618, identified as the entry point of procedure `printf`. If you do not intend to stop at this address, you could use the **return** subcommand to continue execution at the next instruction in `sub` at address 0x100002e0. At this point, the **nexti** subcommand will automatically continue execution to 0x10000428.

If your program has multiple threads, the symbolic thread name of the running thread is displayed when the program stops. For example:

```
stopped in sub at 0x100002d4 ($t4)
10000424 (sub+0x4) 480001f5 bl 0x10000618 (printf)
```

Debugging fdpr reordered executables

You can debug programs that have been reordered with **fdpr** (feedback directed program restructuring, part of Performance Toolbox for AIX) at the instruction level. If optimization options **-R0** or **-R2** are used, additional information is provided enabling **dbx** to map most reordered instruction addresses to the corresponding addresses in the original executable as follows:

```
0xRRRRRRRR = fdpr[0xYYYYYYYY]
```

In this example, `0xRRRRRRRR` is the reordered address and `0xYYYYYYYY` is the original address. In addition, **dbx** uses the traceback entries in the original instruction area to find associated procedure names for the `stopped in` message, the **func** subcommand, and the traceback.

```
(dbx) stepi
stopped in proc_d at 0x1000061c = fdpr[0x10000278]
0x1000061c (???) 9421ffc0      stwu   r1,-64(r1)
(dbx)
```

In the preceding example, **dbx** indicates the program is stopped in the `proc_d` subroutine at address `0x1000061c` in the reordered text section originally located at address `0x10000278`. For more information about **fdpr**, see the **fdpr** command.

Displaying assembly instructions

The **listi** subcommand for the **dbx** command displays a specified set of instructions from the source file. In the default mode, the **dbx** program lists the instructions for the architecture on which it is running. You can override the default mode with the **\$instructionset** and **\$mnemonics** variables of the **set** subcommand for the **dbx** command.

For more information on displaying instructions or disassembling instructions, see the **listi** subcommand for the **dbx** command. For more information on overriding the default mode, see the **\$instructionset** and **\$mnemonics** variables of the **set** subcommand for the **dbx** command.

Customizing the dbx debugging environment

You can customize the debugging environment by creating subcommand aliases and by specifying options in the **.dbxinit** file. You can read **dbx** subcommands from a file using the **-c** flag.

The following sections contain more information about customization options:

Defining a new dbx prompt

The **dbx** prompt is normally the name used to start the **dbx** program. If you specified `/usr/ucb/dbx a.out` on the command line, then the prompt is `/usr/ucb/dbx`.

You can change the prompt with the **prompt** subcommand, or by specifying a different prompt in the **prompt** line of the **.dbxinit** file. Changing the prompt in the **.dbxinit** file causes your prompt to be used instead of the default each time you initialize the **dbx** program.

For example, to initialize the **dbx** program with the debug prompt `debug->`, enter the following line in your **.dbxinit** file:

```
prompt "debug->"
```

Creating dbx subcommand aliases

You can build your own commands from the **dbx** primitive subcommand set. The following commands allow you to build a user alias from the arguments specified. All commands in the replacement string for the alias must be **dbx** primitive subcommands. You can then use your aliases in place of the **dbx** primitives.

The **alias** subcommand with no arguments displays the current aliases in effect; with one argument the command displays the replacement string associated with that alias.

alias [*AliasName* [*CommandName*]]

alias *AliasName* "*CommandString*"

alias *AliasName* (*Parameter1*, *Parameter2*, . . .) "*CommandString*"

The first two forms of the **alias** subcommand are used to substitute the replacement string for the **alias** each time it is used. The third form of aliasing is a limited macro facility. Each parameter specified in the **alias** subcommand is substituted in the replacement string.

The following aliases and associated subcommand names are defaults:

attr

attribute

bfth

stop (in given thread at specified function)

blth

stop (in given thread at specified source line)

c

cont

cv

condition

d

delete

e

edit

h

help

j

status

l

list

m

map

mu

mutex

n

next

p

print

q

quit

r

run

s

step

st

stop

t

where

th
thread

x
registers

You can remove an alias with the [unalias](#) command.

Using the `.dbxinit` file

Each time you begin a debugging session, the **dbx** program searches for special initialization files named **.dbxinit**, which contain lists of **dbx** subcommands to execute. These subcommands are executed before the **dbx** program begins to read subcommands from standard input. When the **dbx** command is started, it checks for a **.dbxinit** file in the user's current directory and in the user's **\$HOME** directory. If a **.dbxinit** file exists, its subcommands run at the beginning of the debug session. If a **.dbxinit** file exists in both the home and current directories, then both are read in that order. Because the current directory **.dbxinit** file is read last, its subcommands can supercede those in the home directory.

Normally, the **.dbxinit** file contains [alias](#) subcommands, but it can contain any valid **dbx** subcommands. For example:

```
$ cat .dbxinit
alias si "stop in"
prompt "dbg-->"
$ dbx a.out
dbx version 3.1
Type 'help' for help.
reading symbolic information . . .
dbg--> alias
si stop in
t where . . .
dbg-->
```

reading **dbx** subcommands from a file

The **-c** invocation option and **.dbxinit** file provide mechanisms for executing **dbx** subcommands before reading from standard input. When the **-c** option is specified, the **dbx** program does not search for a **.dbxinit** file. Use the [source](#) subcommand to read **dbx** subcommands from a file once the debugging session has begun.

After executing the list of commands in the **cmdfile** file, the **dbx** program displays a prompt and waits for input.

You can also use the **-c** option to specify a list of subcommands to be executed when initially starting the **dbx** program.

Debugging spinlocks

You can use the **dbx** program to debug spinlocks. To do so, set the **AIXTHREAD_SPINLOCKS** environment variable to **ON**.

Developing for the **dbx** plug-in framework

dbx provides a plug-in framework for developers who want to add new **dbx** subcommands and event handlers.

Any **dbx** user can create a plug-in that enhances **dbx** with application or library specific commands to aid in debugging.

Note:

1. Because the default **dbx** command is a 64-bit process, all plug-ins need to be compiled to 64-bit to be used with the **dbx** command. In order to load 32-bit plug-ins, use the 32-bit version of the **dbx** command, which is the **dbx32** command.
2. Care should be taken not to confuse `dbx callback routines` and `plug-in interface routines`.

3. `dbx callback routines` are the set of services offered by `dbx` to the plug-in. The plug-in is given access to these routines through a set of function pointers.
4. `plug-in interface routines` are the set of methods `dbx` requires to be implemented by the plug-in.

File format

Each plug-in must be a shared object file.

Naming

To correctly redirect subcommand input, **dbx** requires each plug-in to have a *unique name*.

The file name of the plug-in communicates this unique name to **dbx**. Upon initialization, **dbx** searches a set of predefined and user-specified directories for files whose base name matches the regular expression:

```
^libdbx_+\.so$
```

The following table shows examples of file names that are valid and not valid for **dbx plug-ins**. The corresponding *unique name* is shown for all valid examples:

File Name	Valid	Unique Name
libdbx_sample.so	Yes	sample
libdbx_xyz.so	Yes	xyz
libdbx_my_app.so	Yes	my_app
libdbx.so	No	
libdbx_.so	No	
libdbx_sample.so.plugin	No	
plugin_libdbx_sample.so	No	

Location

dbx allows the user to specify a list of directories to search using the `DBX_PLUGIN_PATH` environment variable. Each directory in the list should be separated by a colon. In the following example, the colon separates two directories.

```
$ export dbx_PLUGIN_PATH=$HOME/dbx_plugins:/mnt/share/dbx_plugins
```

Upon initialization, **dbx** searches for plug-ins. **dbx** also searches the directory of the executable file (if known). This directory is searched after the user-defined directories are searched.

Note: When you use `dbx` to attach to a process, the full path to the executable file cannot be determined.

Loading

A plug-in can be loaded in one of the following ways:

- A plug-in can be automatically loaded and initialized by placing it in a directory that is searched by **dbx**. This occurs at **dbx** initialization time.
- A plug-in can be manually loaded and initialized by specifying its location to the `pluginload dbx` subcommand. This can occur at any time during the **dbx** session.

After a successful automatic or manual plug-in load, a message similar to the following is displayed:

```
(dbx) pluginload /home/user/dbx_plugins/libdbx_sample.so
plug-in "/home/user/dbx_plugins/libdbx_sample.so" loaded
```


Any plug-in whose *unique name* is identical to that of a currently active plug-in is discarded and a warning message similar to the following is displayed.

```
(dbx) pluginload /mnt/share/dbx_plugins/libdbx_sample.so
could not load plug-in
"/mnt/share/dbx_plugins/libdbx_sample.so":
plug-in "/home/user/dbx_plugins/libdbx_sample.so" already loaded.
```

Unloading

Any plug-in, regardless of how it was loaded, can be manually unloaded by specifying its name to the `pluginunload dbx` subcommand. After a plug-in is successfully unloaded, a message similar to the following is displayed.

```
(dbx) pluginunload sample
plug-in "/home/user/dbx_plugins/libdbx_sample.so" unloaded.
```

Version Control

If changes are made to the plug-in framework that would otherwise break the compatibility of the existing plug-ins with earlier versions, a new version identifier is created. This process is true for any significant changes or additions done to the Plug-in Interface or Plug-in dbx callback routine.

To minimize the need for frequent plug-in version changes, some Plug-in dbx callback routines require an additional parameter that represents the size of the buffer. This practice is used for buffer parameters that are based on system structures whose size is not controlled by **dbx**. This allows the size of the system structures to change without requiring updates to the plug-in version.

Currently, the only version identifier is `DBX_PLUGIN_VERSION_1`.

Header File

dbx Plug-in developers can find function prototypes, data structure definitions and macro definitions in the following header file:

```
/usr/include/sys/dbx_plugin.h
```

Plug-in interface

Refer to the `dbx_plugin.h` header file for prototypes and definitions for the Plug-in Interface routines.

Each plug-in *must* implement and export all of the following routines:

- `int dbx_plugin_version(void)`
- `int dbx_plugin_session_init(dbx_plugin_session_t session, constdbx_plugin_service_t *servicep)`
- `void dbx_plugin_session_command(dbx_plugin_session_t session, int argc, char *const argv[])`
- `void dbx_plugin_session_event(dbx_plugin_session_t session, int event, dbx_plugin_event_info_t *event_info)`

int dbx_plugin_version(void)

This routine should return the **dbx Plug-in** version identifier corresponding to the version to which the plug-in conforms. Currently, the only version identifier is `DBX_PLUGIN_VERSION_1`.

int dbx_plugin_session_init(dbx_plugin_session_t session, constdbx_plugin_service_t *servicep)

This routine should perform any initialization needed for the plug-in to function properly before returning control back to **dbx**. This includes setting up any aliases for plug-in subcommands, if desired.

This routine should create a plug-in session that associates the given session identifier with the application program or core file. To identify the process or core file, the session identifier is used by

dbx in Plug-in Interface calls and by the plug-in for plugin **dbx callback routine** requests. This routine also accepts the callback routine structure.

This routine should return zero for successful initialization. If initialization is not successful, **dbx** unloads and discards the plug-in.

void dbx_plugin_session_command(dbx_plugin_session_t session, int argc, char *const argv[])

This routine should accept input from the **dbx** user in the form of arguments provided to the plugin subcommand. The syntax of the plugin subcommand is as follows:

```
plugin Name [arg0 arg1 arg2 ... argn]
```

This allows the **dbx** user to provide any input to any single plug-in. The plug-in has full control over what it accepts as input.

The plugin subcommand passes the command specified by the *arg** parameters to the plug-in specified by the *Name* parameter. (For example, the plug-in name could be `libdbx_Name`. so) Using this routine, **dbx** passes *arg0* through *argn* to the plug-in. *argv[0]* corresponds to *arg0*, *argv[1]* to *arg1*, and so on.

In most circumstances, *arg0* would represent the name of a subcommand defined by the plug-in and *arg1* through *argn* would represent additional flags or arguments. However, this is not a requirement.

Developers are encouraged to implement a `help` subcommand which displays usage information for the plug-in.

void dbx_plugin_session_event(dbx_plugin_session_t session, int event, dbx_plugin_event_info_t *event_info)

In response to application program events, this routine should perform any internal processing required by the plug-in. The routine is invoked once by **dbx** upon the occurrence of each event. The following table describes the event types for which a plug-in is notified:

ID (event)	Associated data (event_info)	Cause
DBX_PLUGIN_EVENT_RESTART	None	The dbx user executed the <code>run</code> subcommand.
DBX_PLUGIN_EVENT_EXIT	Exit code	The application program ended through the exit routine.
DBX_PLUGIN_EVENT_TERM	Terminating signal number	The application program terminated because of an unhandled signal.
DBX_PLUGIN_EVENT_LOAD	dbx_plugin_modinfo_t structure of loaded module	A module was loaded into the application program.
DBX_PLUGIN_EVENT_UNLOAD	dbx_plugin_modinfo_t structure of unloaded module	A module was unloaded from the application program.
DBX_PLUGIN_EVENT_BP	None	The application program has stopped because of a user or internal dbx breakpoint or data watchpoint.

ID (event)	Associated data (event_infp)	Cause
DBX_PLUGIN_EVENT_SIGNAL	Signal number	The application program stopped because of a signal delivery.
DBX_PLUGIN_EVENT_SWTHRD	Handle of current pthread	The dbx user executed the thread <code>current<handle></code> subcommand resulting in a change in the current pthread.

The DBX_PLUGIN_EVENT_BP and DBX_PLUGIN_EVENT_SIGNAL events imply that the application program was started but has stopped. These events are meant to signify that any cached data that the plug-in possesses might no longer be valid. Upon notification of these events, it is more efficient for plug-ins to simply invalidate any cached data rather than refreshing the data. A complete refresh of cached data should only occur when the data is needed. This is especially relevant because some signals might be ignored by dbx and some breakpoints might be internal breakpoints. If the user has no opportunity to run subcommands before the application program starts again, repeatedly refreshing data wastes resources.

```
void dbx_plugin_session_destroy(dbx_plugin_session_t session)
```

This routine should perform any final cleanup and memory management tasks required by the plug-in.

dbx callback routines

The following are the **dbx callback routines** provided for each plug-in through the `dbx_plugin_session_init` routine.

The dbx session callback routine allows you to get characteristics of the dbx session. **dbx** fills in the *flagsp* parameter.

```
typedef int (*dbx_plugin_session_service_t)(dbx_plugin_session_t session,
                                           dbx_plugin_session_flags_t *flagsp).
```

The dbx session callback routine parameters are:

Parameter

Description

session

Session identifier.

flagsp

Session characteristics in any combination of:

- DBX_PLUGIN_SESSION_64BIT

If set, the session represents a 64-bit application program. Otherwise, the session represents a 32-bit application program.

- DBX_PLUGIN_SESSION_CORE

If set, the session represents a core file. Otherwise, the session represents a live process.

The dbx session callback routine return codes are:

- DBX_PLUGIN_SUCCESS
- DBX_PLUGIN_BAD_SESSION session is not valid
- DBX_PLUGIN_BAD_POINTER *flagsp* is NULL

process

The dbx process callback routine allows you to get information on the process being debugged. **dbx** populates the *infop* parameter.

```
typedef int (*dbx_plugin_process_service_t)(dbx_plugin_session_t session,
                                           dbx_plugin_procinfo_t *infop,
                                           size_t procinfo_size)
```

The dbx process callback routine parameters are:

Parameter

Description

session

Session identifier

infop

Allocated dbx_plugin_procinfo_t structure

procinfo_size

Size of dbx_plugin_procinfo_t structure

The dbx process callback routine return codes are:

- DBX_PLUGIN_SUCCESS
- DBX_PLUGIN_BAD_SESSION *session* is not valid
- DBX_PLUGIN_BAD_POINTER *infop* is NULL
- DBX_PLUGIN_BAD_ARG *procinfo_size* is not valid
- DBX_PLUGIN_UNAVAILABLE process not active or info not in core

fds

The dbx fds callback routine allows you to get information on file descriptors for the process. You can either:

- Call iteratively to get information separately on each file descriptor. Or,
- Call once to get the total number of file descriptors and call once again to get information on all file descriptors simultaneously.

If the plug-in passes a non-NULL *infop* buffer, **dbx** populates the buffer with the number of entries requested in **countp*, starting with the file descriptor referenced by **indexp*.

If the plug-in passes a **countp* that is greater than the number of remaining entries, **dbx** retrieves all remaining entries. dbx updates *countp* to reflect the actual number of entries retrieved and *indexp* to reflect the next module index. If the last file descriptor was retrieved, *indexp* is set to -1. If the plug-in passes a NULL *infop* buffer, *indexp* and *countp* are still updated — just as if *infop* were non-NULL.

```
typedef int (*dbx_plugin_fds_service_t)(dbx_plugin_session_t session,
                                       dbx_plugin_fdinfo_t *infop,
                                       size_t fdinfo_size,
                                       unsigned int *indexp,
                                       unsigned int *countp)
```

The dbx fds callback routine parameters are:

Parameter

Description

session

Session identifier

infop

Allocated array of dbx_plugin_fdinfo_t structures or NULL

fdinfo_size

Size of a single dbx_plugin_fdinfo_t structure

indexp

Starting/next file descriptor (where zero corresponds to the first file descriptor)

countp

number of file descriptors

The dbx fds callback routine return codes are:

- DBX_PLUGIN_SUCCESS
- DBX_PLUGIN_BAD_SESSION *session* is not valid
- DBX_PLUGIN_BAD_POINTER *indexp* is NULL or *countp* is NULL
- DBX_PLUGIN_BAD_ARG *fdinfo_size* is not valid or **countp* == 0
- DBX_PLUGIN_UNAVAILABLE process not active or info not in core

modules

The dbx modules callback routine allows you to get information on loaded modules for the process. You can either:

- Call iteratively to get information separately on each module. Or,
- Call once to get the total number of modules and call once again to get information on all modules simultaneously.

If the plug-in passes a non-NULL infop buffer, **dbx** populates the buffer with the number of entries requested in **countp*, starting with the module referenced by **indexp*.

If the plug-in passes a **countp* that is greater than the number of remaining entries, **dbx** retrieves all remaining entries. **dbx** updates *countp* to reflect the actual number of entries retrieved and updates *indexp* to reflect the next module index. If the last module was retrieved, *indexp* is set to -1. If the plug-in passes a NULL infop buffer, *indexp* and *countp* are still updated — just as if *infop* were not NULL.

Note: This routine allocates memory to hold the file name and member character strings. The caller must free this memory when it is no longer needed.

```
typedef int (*dbx_plugin_modules_service_t)(dbx_plugin_session_t session,
                                           dbx_plugin_modinfo_t *infop,
                                           size_t modinfo_size,
                                           unsigned int *indexp,
                                           unsigned int *countp)
```

The dbx modules callback routine parameters are:

Parameter

Description

session

Session identifier

infop

Allocated array of dbx_plugin_modinfo_t structures or NULL

modinfo_size

Size of a single dbx_plugin_modinfo_t structure

indexp

Starting/next module (where zero corresponds to the first module)

countp

Number of modules

The dbx modules callback routine return codes are:

- DBX_PLUGIN_SUCCESS
- DBX_PLUGIN_BAD_SESSION *session* is not valid
- DBX_PLUGIN_BAD_POINTER *indexp* is NULL or *countp* is NULL
- DBX_PLUGIN_BAD_ARG *modinfo_size* is not valid or **countp* == 0

regions

The `dbx regions` callback routine allows you to get information on memory regions for the process.

Retrieved regions can include:

- Main thread stack region (DBX_PLUGIN_REGION_STACK)
- User data region (DBX_PLUGIN_REGION_DATA)
- Process private data region (DBX_PLUGIN_REGION_SDATA)
- Memory mapped region (DBX_PLUGIN_REGION_MMAP)
- Shared memory region (DBX_PLUGIN_REGION_SHM)

You can either:

- Call iteratively to get information separately on one region. Or,
- Call once to get the total number of regions and call once again to get information on all regions simultaneously.

If the plug-in passes a non-NULL *infop* buffer, **dbx** populates the buffer with the number of entries requested in **countp*, starting with the region referenced by **indexp*.

If the plug-in passes a **countp* that is greater than the number of remaining entries, **dbx** retrieves all remaining entries. **dbx** updates *countp* to reflect the actual number of entries retrieved and *indexp* to reflect the next region index.

If the last region was retrieved, *indexp* is set to -1. If the plug-in passes a NULL *infop* buffer, *indexp* and *countp* are still updated — just as if *infop* were non-NULL.

Note: Currently, this routine is only implemented for sessions representing core files. Sufficient information is not available to **dbx** for sessions representing live processes. Calls for such sessions return DBX_PLUGIN_UNAVAILABLE.

```
typedef int (*dbx_plugin_regions_service_t)(dbx_plugin_session_t session,
                                           dbx_plugin_reginfo_t *infop,
                                           size_t reginfo_size,
                                           unsigned int *indexp,
                                           unsigned int *countp)
```

The `dbx regions` callback routine parameters are:

Parameter

Description

session

Session identifier

infop

Allocated array of `dbx_plugin_region_t` structures or NULL

reginfo_size

Size of a single `dbx_plugin_reginfo_t` structure

indexp

Starting/next region (where zero corresponds to the first region)

countp

Number of regions

The `dbx regions` callback routine return codes are:

- DBX_PLUGIN_SUCCESS
- DBX_PLUGIN_BAD_SESSION *session* is not valid
- DBX_PLUGIN_BAD_POINTER *indexp* is NULL or *countp* is NULL
- DBX_PLUGIN_BAD_ARG *reginfo_size* is not valid or **countp* == 0
- DBX_PLUGIN_UNAVAILABLE *session* represents a live processes and regions not accessible

threads

The `dbx threads` callback routine allows you to get information on the kernel threads in the process.

You can either:

- Call iteratively to get information separately on one thread. Or,
- Call once to get the total number of threads and call once again to get information on all threads simultaneously.

If the plug-in passes a non-NULL *infor* buffer, **dbx** populates the buffer with the number of entries requested in **countp*, starting with the thread referenced by **indexp*.

If the plug-in passes a **countp* that is greater than or equal to the number of remaining entries, **dbx** retrieves all remaining entries and updates *countp* to reflect the actual number of entries retrieved.

If the last entry was retrieved, and *countp* is less than its passed value, *indexp* is set to -1. Otherwise, *indexp* is updated to reflect the thread id for the next request.

Note: If the value of *countp* passed is equal to the number of available entries, *countp* remains the same, but *indexp* is not set to -1.

If the plug-in passes a NULL *infor* buffer, *indexp* and *countp* are updated — just as if *infor* were non-NULL.

```
typedef int (*dbx_plugin_threads_service_t)(dbx_plugin_session_t session,
                                           dbx_plugin_thrinfo_t *infor,
                                           size_t thrinfo_size,
                                           tid64_t *indexp,
                                           unsigned int *countp)
```

The `dbx threads` callback routine parameters are:

Parameter

Description

session

Session identifier

infor

Allocated array of `dbx_plugin_thrinfo_t` structures or NULL

thrinfo_size

Size of a single `dbx_plugin_thrinfo_t` structure

indexp

Starting/next thread id (where, on input, zero corresponds to the first thread)

countp

Number of threads

The `dbx threads` callback routine return codes are:

- `DBX_PLUGIN_SUCCESS`
- `DBX_PLUGIN_BAD_SESSION` *session* is not valid
- `DBX_PLUGIN_BAD_POINTER` *indexp* is NULL or *countp* is NULL
- `DBX_PLUGIN_BAD_ID` **indexp* is not a valid id
- `DBX_PLUGIN_BAD_ARG` *thrinfo_size* is not valid or **countp* == 0
- `DBX_PLUGIN_UNAVAILABLE` process not active or entries not in core

pthreads

The `dbx pthreads` callback routine allows you to get information on pthreads in the process, including any kernel thread associations.

You can either:

- Call iteratively to get information separately on one pthread. Or,

- Call once to get the total number of pthreads and call once again to get information on all pthreads simultaneously.

If the plug-in passes a non-NULL *infop* buffer, **dbx** populates the buffer with the number of entries requested in **countp*, starting with the pthread referenced by **indexp*.

If the plug-in passes a **countp* that is greater than the number of remaining entries, **dbx** retrieves all remaining entries. **dbx** updates *countp* to reflect the actual number of entries retrieved and *indexp* to reflect the pthread handle for the next request.

If the last entry was retrieved, *indexp* is set to -1. If the plug-in passes a NULL *infop* buffer, *indexp* and *countp* are still updated — just as if *infop* were non-NULL.

If the first pthread is requested and *countp* is updated to zero, the process is not pthreaded.

```
typedef int (*dbx_plugin_pthreads_service_t)(dbx_plugin_session_t session,
                                             dbx_plugin_pthinfo_t *infop,
                                             size_t pthrdinfo_size,
                                             pthdb_thread_t *indexp,
                                             unsigned int *countp)
```

The dbx pthreads callback routine parameters are:

Parameter

Description

session

Session identifier

infop

Allocated array of dbx_plugin_pthinfo_t structures or NULL

pthrdinfo_size

Size of a single dbx_plugin_pthrdinfo_t structure

indexp

Starting/next pthread handle (where, on input, zero corresponds to the first pthread and DBX_PLUGIN_PTHREAD_CURRENT corresponds to the current pthread in **dbx**)

countp

Number of pthreads

The dbx pthreads callback routine return codes are:

- DBX_PLUGIN_SUCCESS
- DBX_PLUGIN_BAD_SESSION *session* is not valid
- DBX_PLUGIN_BAD_POINTER *indexp* is NULL or *countp* is NULL
- DBX_PLUGIN_BAD_ARG *pthrdinfo_size* is not valid or **countp* == 0

get_thread_context

The dbx *get_thread_context* callback routine allows you to read a kernel thread's general purpose, special purpose, and floating point registers. **dbx** populates the *contextp* parameter.

```
typedef int (*dbx_plugin_reg_service_t)(dbx_plugin_session_t session,
                                         uint64_t reg_flags,
                                         uint64_t id,
                                         dbx_plugin_context_t *contextp,
                                         size_t context_size)
```

The dbx *get_thread_context* callback routine parameters are:

Parameter

Description

session

Session identifier

reg_flags

Logical OR of at least one of DBX_PLUGIN_REG_GPRS, DBX_PLUGIN_REG_SPRS, DBX_PLUGIN_REG_FPRS, DBX_PLUGIN_REG_EXT

id

Kernel thread tid (tid64_t)

contextp

Allocated dbx_plugin_context_t structure

context_size

Size of dbx_plugin_context_t structure. If the DBX_PLUGIN_REG_EXT register flag is used, then the size of the dbx_plugin_extctx_t structure should be used. The dbx_plugin_extctx_t structure is an extended version of dbx_plugin_context_t structure.

The dbx get_thread_context callback routine return codes are:

- DBX_PLUGIN_SUCCESS.
- DBX_PLUGIN_BAD_SESSION *session* is not valid.
- DBX_PLUGIN_BAD_ID *id* is not valid.
- DBX_PLUGIN_BAD_ARG *reg_flags* is not valid or *context_size* is not valid.
- DBX_PLUGIN_BAD_POINTER *contextp* is NULL
- DBX_PLUGIN_UNAVAILABLE process is not active or thread is in kernel mode and registers are not accessible.

set_thread_context

The dbx set_thread_context callback routine allows you to write to a kernel thread's general purpose, special purpose and floating point registers.

```
typedef int (*dbx_plugin_reg_service_t) (dbx_plugin_session_t session,
                                       uint64_t reg_flags,
                                       uint64_t id,
                                       dbx_plugin_context_t *contextp,
                                       size_t context_size)
```

The dbx set_thread_context callback routine parameters are:

Parameter

Description

session

Session identifier

reg_flags

Logical OR of at least one of DBX_PLUGIN_REG_GPRS, DBX_PLUGIN_REG_SPRS, DBX_PLUGIN_REG_FPRS, DBX_PLUGIN_REG_EXT

id

Kernel thread tid (tid64_t)

contextp

Allocated dbx_plugin_context_t structure

context_size

Size of dbx_plugin_context_t structure. If the DBX_PLUGIN_REG_EXT register flag is used, then the size of the dbx_plugin_extctx_t structure should be used. The dbx_plugin_extctx_t structure is an extended version of dbx_plugin_context_t structure.

The dbx set_thread_context callback routine return codes are:

- DBX_PLUGIN_SUCCESS
- DBX_PLUGIN_BAD_SESSION *session* is not valid
- DBX_PLUGIN_BAD_ID *id* is not valid
- DBX_PLUGIN_BAD_ARG *reg_flags* is not valid or *context_size* is not valid

- DBX_PLUGIN_BAD_POINTER *contextp* is NULL
- DBX_PLUGIN_UNAVAILABLE process is not active or thread is in kernel mode and registers are not accessible

get_pthread_context

The dbx `get_pthread_context` callback routine allows you to read a pthread's general purpose, special purpose and floating point registers. dbx populates the *contextp* parameter.

```
typedef int (*dbx_plugin_reg_service_t)(dbx_plugin_session_t session,
                                       uint64_t reg_flags,
                                       uint64_t id,
                                       dbx_plugin_context_t *contextp,
                                       size_t context_size)
```

The dbx `get_pthread_context` callback routine parameters are:

Parameter

Description

session

Session identifier

reg_flags

logical OR of at least one of DBX_PLUGIN_REG_GPRS, DBX_PLUGIN_REG_SPRS, DBX_PLUGIN_REG_FPRS, DBX_PLUGIN_REG_EXT

id

pthread handle (pthdb_pthread_t)

contextp

Allocated dbx_plugin_context_t structure

context_size

size of dbx_plugin_context_t structure. If the DBX_PLUGIN_REG_EXT register flag is used, then the size of the dbx_plugin_extctx_t structure should be used. The dbx_plugin_extctx_t structure is an extended version of dbx_plugin_context_t structure.

The dbx `get_pthread_context` callback routine return codes are:

- DBX_PLUGIN_SUCCESS
- DBX_PLUGIN_BAD_SESSION *session* is not valid
- DBX_PLUGIN_BAD_ID *id* is not valid.
- DBX_PLUGIN_BAD_ARG *reg_flags* is not valid or *context_size* is not valid
- DBX_PLUGIN_BAD_POINTER *contextp* is NULL
- DBX_PLUGIN_UNAVAILABLE process is not active or thread is in kernel mode and registers are not accessible

set_pthread_context

The dbx `set_pthread_context` callback routine allows you to write to a pthread's general purpose, special purpose and floating point registers.

```
typedef int (*dbx_plugin_reg_service_t)(dbx_plugin_session_t session,
                                       uint64_t reg_flags,
                                       uint64_t id,
                                       dbx_plugin_context_t *contextp,
                                       size_t context_size)
```

The dbx `set_pthread_context` callback routine parameters are:

Parameter

Description

session

Session identifier

reg_flags

Logical OR of at least one of DBX_PLUGIN_REG_GPRS, DBX_PLUGIN_REG_SPRS, DBX_PLUGIN_REG_FPRS, DBX_PLUGIN_REG_EXT

id

Pthread handle (pthdb_pthread_t)

contextp

Allocated dbx_plugin_context_t structure

context_size

Size of dbx_plugin_context_t structure. If the DBX_PLUGIN_REG_EXT register flag is used, then the size of the dbx_plugin_extctx_t structure should be used. The dbx_plugin_extctx_t structure is an extended version of dbx_plugin_context_t structure.

The dbx set_pthread_context callback routine return codes are:

- DBX_PLUGIN_SUCCESS
- DBX_PLUGIN_BAD_SESSION *session* is not valid
- DBX_PLUGIN_BAD_ID *id* is not valid
- DBX_PLUGIN_BAD_ARG *reg_flags* is not valid or *context_size* is not valid
- DBX_PLUGIN_BAD_POINTER *contextp* is NULL
- DBX_PLUGIN_UNAVAILABLE process is not active or kernel thread associated with pthread is in kernel mode and registers are not accessible

read_memory

The dbx read_memory callback routine allows you to read from the process's address space. **dbx** populates the buffer parameter.

```
typedef int (*dbx_plugin_mem_service_t)(dbx_plugin_session_t session,
                                       uint64_t addr,
                                       void *buffer,
                                       size_t len)
```

The dbx read_memory callback routine parameters are:

Parameter

Description

session

Session identifier

addr

Address to read from

buffer

Allocated buffer to hold memory contents

len

Number of bytes to read

The dbx read_memory callback routine return codes are:

- DBX_PLUGIN_SUCCESS
- DBX_PLUGIN_BAD_SESSION *session* is not valid
- DBX_PLUGIN_BAD_POINTER *buffer* is NULL
- DBX_PLUGIN_UNAVAILABLE unable to read from *addr*

write_memory

The dbx write_memory callback routine allows you to write to the process's address space.

```
typedef int (*dbx_plugin_mem_service_t)(dbx_plugin_session_t session,
                                       uint64_t addr,
```

```
void *buffer,  
size_t len)
```

The dbx `write_memory` callback routine parameters are:

Parameter

Description

session

Session identifier

addr

Address to write to

buffer

Allocated and initialized buffer

len

Number of bytes to write

The dbx `write_memory` callback routine return codes are:

- `DBX_PLUGIN_SUCCESS`
- `DBX_PLUGIN_BAD_SESSION` *session* is not valid
- `DBX_PLUGIN_BAD_POINTER` *buffer* is NULL
- `DBX_PLUGIN_UNAVAILABLE` unable to write to *addr*

locate_symbol

The dbx `locate_symbol` callback routine allows you to convert symbol names to addresses.

The plug-in must initialize the *name* and *mod* fields of each entry in the symbols parameter array. The *name* field specifies the name of the symbol to be located. The *mod* field specifies the module index of the module in which the lookup should occur. A *mod* field initialized to -1 denotes that all modules should be searched.

dbx populates the *addr* field. Any unknown symbols have an address of zero. If the symbol is located and all modules searched, **dbx** updates the *mod* field with the actual module index of the symbol.

```
typedef int (*dbx_plugin_sym_service_t)(dbx_plugin_session_t session,  
                                       dbx_plugin_sym_t *symbols,  
                                       size_t syminfo_size,  
                                       unsigned int count)
```

The dbx `locate_symbol` callback routine parameters are:

Parameter

Description

session

Session identifier

symbols

Allocated array of `dbx_plugin_sym_t` structures with the *name* and *mod* fields initialized

syminfo_size

Size of `dbx_plugin_sym_t` structure

count

Number of symbols to locate

The dbx `locate_symbol` callback routine return codes are:

- `DBX_PLUGIN_SUCCESS`
- `DBX_PLUGIN_BAD_SESSION` *session* is not valid
- `DBX_PLUGIN_BAD_ARG` *syminfo_size* is not valid
- `DBX_PLUGIN_BAD_POINTER` *symbols* is NULL

what_function

The dbx `what_function` callback routine allows you to convert text addresses to symbols.

The plug-in must initialize the `addr` field of each entry in the symbols parameter array. The `addr` field specifies an instruction address within the function to be identified.

dbx populates the `name` field. Any unknown text address has a name of NULL. **dbx** populates the `mod` field with the actual module index of the text address.

```
typedef int (*dbx_plugin_sym_service_t)(dbx_plugin_session_t session,
                                       dbx_plugin_sym_t *symbols,
                                       size_t syminfo_size,
                                       unsigned int count)
```

The dbx `what_function` callback routine parameters are:

Parameter

Description

session

Session identifier

symbols

Allocated array of `dbx_plugin_sym_t` structures with the `addr` field(s) initialized with text address(es)

syminfo_size

Size of `dbx_plugin_sym_t` structure

count

Number of addresses to convert

The dbx `what_function` callback routine return codes are:

- `DBX_PLUGIN_SUCCESS`
- `DBX_PLUGIN_BAD_SESSION` `session` is not valid
- `DBX_PLUGIN_BAD_ARG` `syminfo_size` is not valid
- `DBX_PLUGIN_BAD_POINTER` `symbols` is NULL

print

The dbx `print` callback routine allows you to display informational output or error output.

```
typedef int (*dbx_plugin_print_service_t)(dbx_plugin_session_t session,
                                          int print_mode,
                                          char *message)
```

The dbx `print` callback routine parameters are:

Parameter

Description

session

session identifier

print_mode

Either `DBX_PLUGIN_PRINT_MODE_OUT` or `DBX_PLUGIN_PRINT_MODE_ERR`

message

Character string for **dbx** to display

The dbx `print` callback routine return codes are:

- `DBX_PLUGIN_SUCCESS`
- `DBX_PLUGIN_BAD_SESSION` `session` is not valid
- `DBX_PLUGIN_BAD_ARG` `print_mode` is not valid
- `DBX_PLUGIN_BAD_POINTER` `message` is NULL

alias

The dbx alias callback routine allows you to create an alias for a plug-in subcommand.

The syntax of the `plugin dbx` subcommand requires the dbx user to type a prefix of `plugin Name` for each plug-in subcommand invocation. To provide a way to shorten such invocations, **dbx** allows plug-ins to create new aliases.

The *alias* and *expansion* parameters should provide a description of the new alias. The syntax is the same as the syntax defined for the alias **dbx** subcommand.

The following are example invocations of the dbx alias callback routine:

```
alias("intprt", "plugin xyz interpret");
```

```
alias("intprt2(addr, count, format)", "addr / count format; plugin xyz interpret addr");
```

Note: If you try to create an alias that has the same name as an existing alias, the request is denied and a warning message is displayed. Plug-in developers are encouraged to carry out alias creation in a way that allows users to correct alias conflicts. One way to accomplish this is by reading alias definitions from a configuration file that is packaged with the plug-in.

```
typedef int (*dbx_plugin_alias_service_t)(dbx_plugin_session_t session,  
                                         const char *alias,  
                                         const char *expansion)
```

The dbx alias callback routine parameters are:

Parameter

Description

session

Session identifier

alias

Character string representing the alias name and optional parameter

expansion

Character string representing the alias expansion

The dbx alias callback routine return codes are:

- `DBX_PLUGIN_SUCCESS`
- `DBX_PLUGIN_BAD_SESSION` *session* is not valid
- `DBX_PLUGIN_BAD_ARG` *alias* is not valid
- `DBX_PLUGIN_BAD_POINTER` *alias* is NULL or *expansion* is NULL
- `DBX_PLUGIN_UNAVAILABLE` an alias with an identical name already exists

Example

1. The following example defines a `help` subcommand and a `hello` subcommand:

```
example.c:  
  
#include <sys/dbx_plugin.h>  
  
dbx_plugin_session_t sid;  
dbx_plugin_services_t dbx;  
  
static void usage(void);  
static void hello_cmd(void);  
  
int  
dbx_plugin_version(void) {  
    return DBX_PLUGIN_VERSION_1;  
}  
  
int dbx_plugin_session_init(dbx_plugin_session_t session,  
                           const dbx_plugin_services_t *servicep) {
```

```

/* record session identifier */
sid= session;

/* record dbx service */
memcpy(&dbx, servicep, sizeof(dbx_plugin_services_t));

(*(dbx.alias))(sid, "hello", "plugin example hello");
(*(dbx.alias))(sid, "help", "plugin example help");

return 0;
}

void
dbx_plugin_session_command(dbx_plugin_session_t session,
                           int argc,
                           char *const argv[]) {
    if (argc == 0 || (argc == 1 && strcmp(argv[0], "help") == 0)) {
        usage();
        return;
    }
    if (argc == 1 && strcmp(argv[0], "hello") == 0) {
        hello_cmd();
        return;
    }
    (*(dbx.print))(sid, DBX_PLUGIN_PRINT_MODE_ERR,
                  "unrecognized command\n");
}

void
dbx_plugin_session_event(dbx_plugin_session_t session,
                          int event,
                          dbx_plugin_event_info_t *event_infop) {
    /* ignore event notifications */
}

void
dbx_plugin_session_destroy(dbx_plugin_session_t session){
    /* no clean up to perform */
}

static
void
usage(void) {
    (*(dbx.print))(sid, DBX_PLUGIN_PRINT_MODE_OUT,
                  "Subcommands for Plug-in \"example\":\n\n" \
                  "  help - displays this output\n" \
                  "  hello - displays a greeting\n" \
                  "\n");
}

static
void
hello_cmd(void) {
    (*(dbx.print))(sid, DBX_PLUGIN_PRINT_MODE_OUT,
                  "Hello dbx World!\n");
}

```

example.exp:

```

dbx_plugin_version
dbx_plugin_session_init
dbx_plugin_session_command
dbx_plugin_session_event
dbx_plugin_session_destroy

```

2. To compile the example plug-in, type:

```

cc -q64 -o libdbx_example.so example.c -bM:Sre -bE:example.exp -bnoentry

```

List of dbx subcommands

The commands and sub commands for the **dbx** debug program are documented in the **dbx** command .

The **dbx** debug program provides subcommands for performing the following task categories:

Setting and deleting breakpoints

Subcommand	Description
<u>clear</u>	Removes all stops at a given source line.
<u>cleari</u>	Removes all breakpoints at an address.
<u>delete</u>	Removes the traces and stops corresponding to the specified numbers.
<u>status</u>	Displays the currently active trace and stop subcommands.
<u>stop</u>	Stops execution of the application program.

Running your program

Subcommand	Description
<u>cont</u>	Continues running the program from the current breakpoint until the program finishes or another breakpoint is encountered.
<u>detach</u>	Exits the debug program, but continues running the application.
<u>down</u>	Moves a function down the stack.
<u>goto</u>	Causes the specified source line to be the next line run.
<u>gotoi</u>	Changes program counter addresses.
<u>next</u>	Runs the application program up to the next source line.
<u>nexti</u>	Runs the application program up to the next source instruction.
<u>rerun</u>	Begins running an application.
<u>return</u>	Continues running the application program until a return to the specified procedure is reached.
<u>run</u>	Begins running an application.
<u>skip</u>	Continues execution from the current stopping point.
<u>step</u>	Runs one source line.
<u>stepi</u>	Runs one source instruction.
<u>up</u>	Move a function up the stack.

Tracing program execution

Subcommand	Description
<u>trace</u>	Prints tracing information.
<u>tracei</u>	Turns on tracing.
<u>where</u>	Displays a list of all active procedures and functions.

Ending program execution

Subcommand	Description
<u>quit</u>	Quits the dbx debug program.

Displaying the source file

Subcommand	Description
<u>edit</u>	Invokes an editor on the specified file.
<u>file</u>	Changes the current source file to the specified file.
<u>func</u>	Changes the current function to the specified function or procedure.
<u>list</u>	Displays lines of the current source file.
<u>listi</u>	Lists instructions from the application.
<u>move</u>	Changes the next line to be displayed.

Subcommand	Description
<u>/ (Search)</u>	Searches forward in the current source file for a pattern.
<u>? (Search)</u>	Searches backward in the current source file for a pattern.
<u>use</u>	Sets the list of directories to be searched when looking for a file.

Printing and modifying variables, expressions, and types

Subcommand	Description
<u>assign</u>	Assigns a value to a variable.
<u>case</u>	Changes the way in which dbx interprets symbols.
<u>dump</u>	Displays the names and values of variables in the specified procedure.
<u>print</u>	Prints the value of an expression or runs a procedure and prints the return code.
<u>set</u>	Assigns a value to a nonprogram variable.
<u>unset</u>	Deletes a nonprogram variable.
<u>whatis</u>	Displays the declaration of application program components.
<u>whereis</u>	Displays the full qualifications of all the symbols whose names match the specified identifier.
<u>which</u>	Displays the full qualification of the specified identifier.

Thread debugging

Subcommand	Description
<u>attribute</u>	Displays information about all or selected attributes objects.
<u>condition</u>	Displays information about all or selected condition variables.
<u>mutex</u>	Displays information about all or selected mutexes.
<u>thread</u>	Displays and controls threads.
<u>tstophwp</u>	Sets a thread-level hardware watchpoint stop.
<u>ttracehwp</u>	Sets a thread-level hardware watchpoint trace.
<u>tstop</u>	Sets a source-level breakpoint stop for a thread.
<u>tstopi</u>	Sets an instruction-level breakpoint stop for a thread.
<u>ttrace</u>	Sets a source-level trace for a thread.
<u>ttracei</u>	Sets an instruction-level trace for a thread.
<u>tnext</u>	Runs a thread up to the next source line.
<u>tnexti</u>	Runs a thread up to the next machine instruction.
<u>tstep</u>	Runs a thread one source line.
<u>tstepi</u>	Runs a thread one machine instruction.
<u>tskip</u>	Skips breakpoints for a thread.

Multiprocess debugging

Subcommand	Description
<u>multiproc</u>	Enables or disables multiprocess debugging.

Procedure calling

Subcommand	Description
<u>call</u>	Runs the object code associated with the named procedure or function.
<u>print</u>	Prints the value of an expression or runs a procedure and prints the return code.

Signal handling

Subcommand	Description
<u>catch</u>	Starts trapping a signal before that signal is sent to the application program.
<u>ignore</u>	Stops trapping a signal before that signal is sent to the application program.

Machine-level debugging

Subcommand	Description
<u>display memory</u>	Displays the contents of memory.
<u>gotoi</u>	Changes program counter addresses.
<u>map</u>	Displays address maps and loader information for the application program.
<u>nexti</u>	Runs the application program up to the next machine instruction.
<u>registers</u>	Displays the values of all general-purpose registers, system-control registers, floating-point registers, and the current instruction register.
<u>stepi</u>	Runs one source instruction.
<u>stopi</u>	Sets a stop at a specified location.
<u>tracei</u>	Turns on tracing.

Debugging environment control

Subcommand	Description
<u>alias</u>	Displays and assigns aliases for dbx subcommands.
<u>help</u>	Displays help information for dbx subcommands or topics.
<u>prompt</u>	Changes the dbx prompt to the specified string.
<u>screen</u>	Opens an Xwindow for dbx command output.
<u>sh</u>	Passes a command to the shell for execution.
<u>source</u>	Reads dbx commands from a file.
<u>unalias</u>	Removes an alias.

Error-logging overview

The error-logging process begins when an operating system module detects an error.

The error-detecting segment of code then sends error information to either the **errsave** and **errlast** kernel service or to the **errlog** subroutine. This error information is then written to the **/dev/error** special file. This process then adds a time stamp to the collected data. The **errdemon** daemon constantly checks the **/dev/error** file for new entries, and when new data is written, the daemon conducts a series of operations.

Before an entry is written to the error log, the **errdemon** daemon compares the label sent by the kernel or application code to the contents of the Error Record Template Repository. If the label matches an item in the repository, the daemon collects additional data from other parts of the system.

To create an entry in the error log, the **errdemon** daemon retrieves the appropriate template from the repository, the resource name of the unit that detected the error, and detail data. Also, if the error signifies a hardware-related problem and hardware vital product data (VPD) exists, the daemon retrieves the VPD from the Object Data Manager. When you access the error log, either through SMIT or with the **errpt** command, the error log is formatted according to the error template in the error template repository and presented in either a summary or detailed report. Entries can also be retrieved using the services provided in **liberrlog**, **errlog_open**, **errlog_close**, **errlog_find_first**, **errlog_find_next**, **errlog_find_sequence**, **errlog_set_direction**, and **errlog_write**. **errlog_write** provides a limited update capability.

Most entries in the error log are attributable to hardware and software problems, but informational messages can also be logged.

The **diag** command uses the error log to diagnose hardware problems. To correctly diagnose new system problems, the system deletes hardware-related entries older than 90 days from the error log. The system deletes software-related entries 30 days after they are logged.

You should be familiar with the following terms:

Term	Description
error ID	A 32-bit CRC hexadecimal code used to identify a particular failure. Each error record template has a unique error ID.
error label	The mnemonic name for an error ID.
error log	The file that stores instances of errors and failures encountered by the system.
error log entry	A record in the system error log that describes a hardware failure, a software failure, or an operator message. An error log entry contains captured failure data.
error record template	A description of information displayed when the error log is formatted for a report, including information on the type and class of the error, probable causes, and recommended actions. Collectively, the templates comprise the Error Record Template Repository.

Related information

[Error Logging Special Files](#)

[errsave Command](#)

[errlog Command](#)

[crontab Command](#)

[errclear Command](#)

[errdead Command](#)

[errdemon Command](#)

[errinstall Command](#)

[errlogger Command](#)

[errmsg Command](#)

[errpt Command](#)

[errstop Command](#)

[errupdate Command](#)

[odmadd Command](#)

[errstop Command](#)

[odmget Command](#)

[snap Command](#)

Error-logging facility

The error-logging facility records hardware and software failures in the error log for information purposes or for fault detection and corrective action.

Refer to the following to use the error-logging facility:

In AIX Version 4 some of the error log commands are delivered in an optionally installable package called **bos.sysmgt.serv_aid**. The base system (**bos.rte**) includes the following services for logging errors to the error log file:

- errlog subroutines
- **errsave** and **errlast** kernel service
- error device driver (**/dev/error**)
- **error** daemon
- **errstop** command

The commands required for licensed program installation (**errinstall** and **errupdate**) are also included in **bos.rte**. For information on transferring your system's error log file to a system that has the Software Service Aids package installed, see [Transferring Your Error Log to Another System](#).

Managing error logging

Error logging is automatically started by the **rc.boot** script during system initialization and is automatically stopped by the **shutdown** script during system shutdown.

The error log analysis performed by the **diag** command analyzes hardware error entries. The default length of time that hardware error entries remain in the error log is 90 days. If you remove hardware error entries less than 90 days old, you can limit the effectiveness of this error log analysis.

Transferring your error log to another system

The **errclear**, **errdead**, **errlogger**, **errmsg**, and **errpt** commands are part of the optionally installable Software Service Aids package (**bos.sysmgmt.serv_aid**). You need the Software Service Aids package to generate reports from the error log or to delete entries from the error log. You can install the Software Service Aids package on your system, or you can transfer your system's error log file to a system that has the Software Service Aids package installed.

Determine the path to your system's error log file by running the following command:

```
/usr/lib/errdemon -l
```

You can transfer the file to another system in a number of ways. You can:

- Copy the file to a remotely mounted file system using the **cp** command
- Copy the file across the network connection using the **rcp**, **ftp**, or **tftp** commands
- Copy the file to removable media using the **tar** or **backup** command and restore the file onto another system.

You can format reports for an error log copied to your system from another system by using the **-i** flag of the **errpt** command. The **-i** flag allows you to specify the path name of an error log file other than the default. Likewise, you can delete entries from an error log file copied to your system from another system by using the **-i** flag of the **errclear** command.

Configuring error logging

You can customize the name and location of the error log file and the size of the internal error buffer to suit your needs. You can also control the logging of duplicate errors.

Listing the current settings

To list the current settings, run **/usr/lib/errdemon -l**. The values for the error log file name, error log file size, and buffer size that are currently stored in the error-log configuration database display on your screen.

Customizing the log file location

To change the file name used for error logging, run the **/usr/lib/errdemon -i FileName** command. The specified file name is saved in the error log configuration database, and the error daemon is immediately restarted.

Customizing the log file size

To change the maximum size of the error log file, type:

```
/usr/lib/errdemon -s LogSize
```

The specified size limit for the log file is saved in the error-log configuration database, and the error daemon is immediately restarted. If the size limit for the log file is smaller than the size of the log file currently in use, the current log file is renamed by appending **.old** to the file name, and a new log file is created with the specified size limit. The amount of space specified is reserved for the error log file and is not available for use by other files. Therefore, be careful not to make the log excessively large. But, if you

make the log too small, important information may be overwritten prematurely. When the log file size limit is reached, the file *wraps*, that is, the oldest entries are overwritten by new entries.

Customizing the buffer size

To change the size of the error log device driver's internal buffer, type:

```
/usr/lib/errdemon -B BufferSize
```

The specified buffer size is saved in the error-log configuration database, and if it is larger than the buffer size currently in use, the in-memory buffer is immediately increased. If it is smaller than the buffer size currently in use, the new size is put into effect the next time that the error daemon is started after the system is rebooted. The buffer cannot be made smaller than the hard-coded default of 8 KB. The size you specify is rounded up to the next integral multiple of the memory page size (4 KBs). The memory used for the error log device driver's in-memory buffer is not available for use by other processes (the buffer is pinned).

Be careful not to impact your system's performance by making the buffer excessively large. But, if you make the buffer too small, the buffer may become full if error entries are arriving faster than they are being read from the buffer and put into the log file. When the buffer is full, new entries are discarded until space becomes available in the buffer. When this situation occurs, an error log entry is created to inform you of the problem, and you can correct the problem by enlarging the buffer.

Customizing duplicate error handling

By default, starting with AIX 5.1, the error daemon eliminates duplicate errors by looking at each error that is logged. An error is a duplicate if it is identical to the previous error, and if it occurs within the approximate time interval specified with **/usr/lib/errdemon -t time-interval**. The default time value is 10000, or 10 seconds. The value is in milliseconds.

The **-m maxdups** flag controls how many duplicates can accumulate before a duplicate entry is logged. The default value is 1000. If an error, followed by 1000 occurrences of the same error, is logged, a duplicate error is logged at that point rather than waiting for the time interval to expire or for a unique error to occur.

For example, if a device handler starts logging many identical errors rapidly, most will not appear in the log. Rather, the first occurrence will be logged. Subsequent occurrences will not be logged immediately, but are only counted. When the time interval expires, the **maxdups** value is reached, or when another error is logged, an alternate form of the error is logged, giving the times of the first and last duplicate and the number of duplicates.

Note: The time interval refers to the time since the last error, not the time since the first occurrence of this error, that is, it is reset each time an error is logged. Also, to be a duplicate, an error must exactly match the previous error. If, for example, anything about the detail data is different from the previous error, then that error is considered unique and logged as a separate error.

Removing error log entries

Entries are removed from the error log when the root user runs the **errclear** command, when the **errclear** command is automatically invoked by a daily **cron** job, or when the error log file wraps as a result of reaching its maximum size. When the error log file reaches the maximum size specified in the error-log configuration database, the oldest entries are overwritten by the newest entries.

Automatic removal

A **crontab** file provided with the system deletes hardware errors older than 90 days and other errors older than 30 days. To display the **crontab** entries for your system, type:

```
crontab -l Command
```

To change these entries, type:

```
crontab -e Command
```

errclear command

The **errclear** command can be used to selectively remove entries from the error log. The selection criteria you may specify include the error ID number, sequence number, error label, resource name, resource class, error class, and error type. You must also specify the age of entries to be removed. The entries that match the selection criteria you specified, and are older than the number of days you specified, will be removed.

Enabling and disabling logging for an event

You can disable logging or reporting of a particular event by modifying the **Log** or the **Report** field of the error template for the event. You can use the **errupdate** command to change the current settings for an event.

Showing events for which logging is disabled

To list all events for which logging is currently disabled, type:

```
errpt -t -F Log=0
```

Events for which logging is disabled are not saved in the error log file.

Showing events for which reporting is disabled

To list all events for which reporting is currently disabled, type:

```
errpt -t -F Report=0
```

Events for which reporting is disabled are saved in the error log file when they occur, but they are not displayed by the **errpt** command.

Changing the current setting for an event

To change the current settings for an event, you can use the **errupdate** command. The necessary input to the **errupdate** command can be in a file or from standard input.

The following example uses standard input. To disable the reporting of the **ERRLOG_OFF** event (error ID 192AC071), type the following to run the **errupdate** command:

```
errupdate <Enter>
=192AC071: <Enter>
Report=False <Enter>
<Ctrl-D>
<Ctrl-D>
```

Logging maintenance activities

The **errlogger** command allows the system administrator to record messages in the error log. Whenever you perform a maintenance activity, such as clearing entries from the error log, replacing hardware, or applying a software fix, it is a good idea to record this activity in the system error log.

The **ras_logger** command provides a way to log any error from the command line. It can be used to test newly created error templates and provides a way to log an error from a shell script.

Redirecting syslog messages to error log

Some applications use syslog for logging errors and other events. To list error log messages and syslog messages in a single report, redirect the syslog messages to the error log. You can do this by specifying *errlog* as the destination in the syslog configuration file (*/etc/syslog.conf*). See the **syslogd** daemon for more information.

Directing error log messages to syslog

You can log error log events in the **syslog** file by using the **logger** command with the concurrent error notification capabilities of error log. For example, to log system messages (syslog), add an **errnotify** object with the following contents:

```
errnotify:
  en_name = "syslog1"
  en_persistenceflg = 1
  en_method = "logger Msg from Error Log: `errpt -l $1 | grep -v 'ERROR_ID TIMESTAMP'`"
```

For example, create a file called **/tmp/syslog.add** with these contents. Then run the **odmadd /tmp/syslog.add** command (you must be logged in as root user).

For more information about concurrent error notification, see [Error Notification](#).

Error notification

The Error Notification object class specifies the conditions and actions to be taken when errors are recorded in the system error log. The user specifies these conditions and actions in an Error Notification object.

Each time an error is logged, the **error notification** daemon determines if the error log entry matches the selection criteria of any of the Error Notification objects. If matches exist, the daemon runs the programmed action, also called a *notify method*, for each matched object.

The Error Notification object class is located in the **/etc/objrepos/errnotify** file. Error Notification objects are added to the object class by using Object Data Manager (ODM) commands. Only processes running with the root user authority can add objects to the Error Notification object class. Error Notification objects contain the following descriptors:

en_alertflg

Identifies whether the error can be alerted. This descriptor is provided for use by alert agents associated with network management applications using the SNA Alert Architecture. The valid alert descriptor values are:

TRUE

can be alerted

FALSE

cannot be alerted

en_class

Identifies the class of the error log entries to match. The valid **en_class** descriptor values are:

H

Hardware Error class

S

Software Error class

O

Messages from the [errlogger](#) command

U

Undetermined

en_crcid

Specifies the error identifier associated with a particular error. An error identifier can be any numeric value that is valid as a [Predefined Attribute](#) (PdAt) object class attribute value. The **errpt** command displays error identifiers as hexadecimal. For example, to select an entry that the **errpt** command displays with IDENTIFIER: 67581038, specify **en_crcid = 0x67581038**.

en_dup

If set, identifies whether duplicate errors as defined by the kernel should be matched. The valid **en_dup** descriptor values are:

TRUE

Error is a duplicate.

FALSE

Error is not a duplicate.

en_err64

If set, identifies whether errors from a 64-bit or 32-bit environment should be matched. The valid **en_err64** descriptors value are:

TRUE

Error is from a 64-bit environment.

FALSE

Error is from a 32-bit environment.

en_label

Specifies the label associated with a particular error identifier as defined in the output of the **errpt -t** command.

en_method

Specifies a user-programmable action, such as a shell script or command string, to be run when an error matching the selection criteria of this Error Notification object is logged. The error notification daemon uses the **sh -c** command to execute the notify method.

The following key words are automatically expanded by the **error notification** daemon as arguments to the notify method.

\$1

Sequence number from the error log entry

\$2

Error ID from the error log entry

\$3

Class from the error log entry

\$4

Type from the error log entry

\$5

Alert flags value from the error log entry

\$6

Resource name from the error log entry

\$7

Resource type from the error log entry

\$8

Resource class from the error log entry

\$9

Error label from the error log entry

en_name

Uniquely identifies the object. This unique name is used when removing the object.

en_persistenceflg

Designates whether the Error Notification object should be automatically removed when the system is restarted. For example, to avoid erroneous signaling, Error Notification objects containing methods that send a signal to another process should not persist across system restarts. The receiving process and its process ID do not persist across system restarts.

The creator of the Error Notification object is responsible for removing the Error Notification object at the appropriate time. In the event that the process terminates and fails to remove the Error Notification object, the **en_persistenceflg** descriptor ensures that obsolete Error Notification objects are removed when the system is restarted.

The valid **en_persistenceflg** descriptor values are:

0

non-persistent (removed at boot time)

1

persistent (persists through boot)

en_pid

Specifies a process ID (PID) for use in identifying the Error Notification object. Objects that have a PID specified should have the **en_persistenceflg** descriptor set to 0.

en_rclass

Identifies the class of the failing resource. For the hardware error class, the resource class is the device class. The resource error class is not applicable for the software error class.

en_resource

Identifies the name of the failing resource. For the hardware error class, a resource name is the device name.

en_rtype

Identifies the type of the failing resource. For the hardware error class, a resource type is the device type by which a resource is known in the devices object class.

en_symptom

Enables notification of an error accompanied by a symptom string when set to **TRUE**.

en_type

Identifies the severity of error log entries to match. The valid **en_type** descriptor values are:

INFO

Informational

PEND

Impending loss of availability

PERM

Permanent

PERF

Unacceptable performance degradation

TEMP

Temporary

UNKN

Unknown

Examples

1. To create a notify method that mails a formatted error entry to root each time a disk error of type PERM is logged, create a file called `/tmp/en_sample.add` containing the following Error Notification object:

```
errnotify:
  en_name = "sample"
  en_persistenceflg = 0
  en_class = "H"
  en_type = "PERM"
  en_rclass = "disk"
  en_method = "errpt -a -l $1 | mail -s 'Disk Error' root"
```

To add the object to the Error Notification object class, type:

```
odmadd /tmp/en_sample.add
```

The **odmadd** command adds the Error Notification object contained in `/tmp/en_sample.add` to the **errnotify** file.

2. To verify that the Error Notification object was added to the object class, type:

```
odmget -q"en_name='sample'" errnotify
```

The **odmget** command locates the Error Notification object within the **errnotify** file that has an **en_name** value of "sample" and displays the object. The following output is returned:

```
errnotify:
  en_pid = 0
  en_name = "sample"
  en_persistenceflg = 0
  en_label = ""
  en_crcid = 0
  en_class = "H"
  en_type = "PERM"
  en_alertflg = ""
```

```
en_resource = ""
en_rtype = ""
en_rclass = "disk"
en_method = "errpt -a -l $1 | mail -s 'Disk Error' root"
```

3. To delete the sample Error Notification object from the Error Notification object class, type:

```
odmdelete -q"en_name='sample'" -o errnotify
```

The **odmdelete** command locates the Error Notification object within the **errnotify** file that has an **en_name** value of "sample" and removes it from the Error Notification object class.

4. To send an email to root when a duplicate error occurs, create a file called **/tmp/en_sample.add** containing the following error notification stanza:

```
errnotify:
  en_name = "errdupxmp"
  en_persistenceflg = 1
  en_dup = "TRUE"
  en_method = "/usr/lib/dupmethod $1"
```

Create the **/usr/lib/dupmethod** script as follows:

```
#!/bin/sh
# email root when a duplicate error is logged.
# We currently don't clear the duplicate from the log.
#
# Input:
# $1 contains the error log sequence number.
#
# Use errpt to generate the body of this email.
/usr/bin/errpt -al$1 | /usr/bin/mail -s "Duplicate Error Logged" root >/dev/null

# Now delete that error (currently not done)
#/usr/bin/errclear -l$1 0
exit $?
```

Error logging tasks

This section describes the error-logging tasks and information.

Error-logging tasks and information to assist you in using the error logging facility include:

- [Reading an Error Report](#)
- [Examples of Detailed Error Reports](#)
- [Example of a Summary Error Report](#)
- [Generating an Error Report](#)
- [Stopping an Error Log](#)
- [Cleaning an Error Log](#)
- [Copying an Error Log to Diskette or Tape](#)
- [Using the liberrlog Services](#)

Reading an error report

To obtain a report of all errors logged in the 24 hours prior to the failure, type:

```
errpt -a -s mmdhmmmy | pg
```

where *mmdhmmmy* represents the month, day, hour, minute, and year 24 hours prior to the failure.

An error-log report contains the following information:

Note: Not all errors generate information for each of the following categories.

LABEL

Predefined name for the event.

ID

Numerical identifier for the event.

Date/Time

Date and time of the event.

Sequence Number

Unique number for the event.

Machine ID

Identification number of your system processor unit.

Node ID

Mnemonic name of your system.

Class

General source of the error. The possible error classes are:

H

Hardware. (When you receive a hardware error, refer to your system operator guide for information about performing diagnostics on the problem device or other piece of equipment. The diagnostics program tests the device and analyzes the error log entries related to it to determine the state of the device.)

S

Software.

O

Informational messages.

U

Undetermined (for example, a network).

Type

Severity of the error that has occurred. The following types of errors are possible:

PEND

The loss of availability of a device or component is imminent.

PERF

The performance of the device or component has degraded to below an acceptable level.

PERM

Condition that could not be recovered from. Error types with this value are usually the most severe errors and are more likely to mean that you have a defective hardware device or software module. Error types other than PERM usually do not indicate a defect, but they are recorded so that they can be analyzed by the diagnostics programs.

TEMP

Condition that was recovered from after a number of unsuccessful attempts. This error type is also used to record informational entries, such as data transfer statistics for DASD devices.

UNKN

It is not possible to determine the severity of the error.

INFO

The error log entry is informational and was not the result of an error.

Resource name

Name of the resource that has detected the error. For software errors, this is the name of a software component or an executable program. For hardware errors, this is the name of a device or system component. It does not indicate that the component is faulty or needs replacement. Instead, it is used to determine the appropriate diagnostic modules to be used to analyze the error.

Resource class

General class of the resource that detected the failure (for example, a device class of disk).

Resource type

Type of the resource that detected the failure (for example, a device type of 355mb).

Location code

Path to the device. There may be up to four fields, which refer to drawer, slot, connector, and port, respectively.

VPD

Vital product data. The contents of this field, if any, vary. Error log entries for devices typically return information concerning the device manufacturer, serial number, Engineering Change levels, and Read Only Storage levels.

Description

Summary of the error.

Probable cause

List of some of the possible sources of the error.

User causes

List of possible reasons for errors due to user mistakes. An improperly inserted disk and external devices (such as modems and printers) that are not turned on are examples of user-caused errors.

Recommended actions

Description of actions for correcting a user-caused error.

Install causes

List of possible reasons for errors due to incorrect installation or configuration procedures. Examples of this type of error include hardware and software mismatches, incorrect installation of cables or cable connections becoming loose, and improperly configured systems.

Recommended actions

Description of actions for correcting an installation-caused error.

Failure causes

List of possible defects in hardware or software.

Note: A failure causes section in a software error log usually indicates a software defect. Logs that list user or installation causes or both, but not failure causes, usually indicate that the problem is not a software defect.

If you suspect a software defect, or are unable to correct user or installation causes, report the problem to your software service department.

Recommended actions

Description of actions for correcting the failure. For hardware errors, `PERFORM PROBLEM DETERMINATION PROCEDURES` is one of the recommended actions listed. For hardware errors, this will lead to running the diagnostic programs.

Detailed data

- Failure data that is unique for each error log entry, such as device sense data.
- Information on the current working directory of the process, such as `FILE SYSTEM SERIAL NUMBER` and `INODE NUMBER` when the process dumps the core.

To display a shortened version of the detailed report produced by the `-a` flag, use the `-A` flag. The `-A` flag is not valid with the `-a`, `-g`, or `-t` flags. The items reported when you use `-A` to produce the shortened version of the report are:

- Label
- Date and time
- Type
- Resource name
- Description
- Detail data

The example output of this flag is in the following format:

```
LABEL:          STOK_RCVRY_EXIT
Date/Time:      Tue Dec 14 15:25:33
```

```
Type:                TEMP Resource Name:  tok0
Description PROBLEM RESOLVED
Detail Data FILE NAME line: 273 file: stok_wdt.c
SENSE DATA
0000 0000 0000 0000 0000 0000 0000 0000 DEVICE ADDRESS 0004 AC62 25F1
```

Reporting can be turned off for some errors. To show which errors have reporting turned off, type:

```
errpt -t -F report=0 | pg
```

If reporting is turned off for any errors, enable reporting of all errors using the **errupdate** command.

Logging may also have been turned off for some errors. To show which errors have logging turned off, type:

```
errpt -t -F log=0 | pg
```

If logging is turned off for any errors, enable logging for all errors using the **errupdate** command. Logging all errors is useful if it becomes necessary to re-create a system error.

Examples of detailed error reports

The following are sample error-report entries that are generated by issuing the **errpt -a** command.

An error-class value of **H** and an error-type value of **PERM** indicate that the system encountered a hardware problem (for example, with a SCSI adapter device driver) and could not recover from it. Diagnostic information might be associated with this type of error. If so, it displays at the end of the error listing, as illustrated in the following example of a problem encountered with a device driver:

```
LABEL:      SCSI_ERR1
ID:         0502F666

Date/Time:   Jun 19 22:29:51
Sequence Number: 95
Machine ID:  123456789012
Node ID:     host1
Class:       H
Type:        PERM
Resource Name: scsi0
Resource Class: adapter
Resource Type: hscsi
Location:    00-08
VPD:
  Device Driver Level.....00
  Diagnostic Level.....00
  Displayable Message.....SCSI
  EC Level.....C25928
  FRU Number.....30F8834
  Manufacturer.....IBM97F
  Part Number.....59F4566
  Serial Number.....00002849
  ROS Level and ID.....24
  Read/Write Register Ptr.....0120

Description
ADAPTER ERROR

Probable Causes
ADAPTER HARDWARE CABLE
CABLE TERMINATOR DEVICE

Failure Causes
ADAPTER
CABLE LOOSE OR DEFECTIVE

Recommended Actions
PERFORM PROBLEM DETERMINATION PROCEDURES
CHECK CABLE AND ITS CONNECTIONS

Detail Data
SENSE DATA
0000 0000 0000 0000 0000 0000 0000 0000 0000 0000 0000 0000

Diagnostic Log sequence number: 153
Resource Tested:      scsi0
Resource Description:  SCSI I/O Controller
```

```

Location:          00-08
SRN:              889-191
Description:      Error log analysis indicates hardware failure.
Probable FRUs:
  SCSI Bus       FRU: n/a           00-08
                Fan Assembly
  SCSI2         FRU: 30F8834      00-08
                SCSI I/O Controller

```

An error-class value of **H** and an error-type value of **PEND** indicate that a piece of hardware (the Token Ring) may become unavailable soon due to numerous errors detected by the system.

```

LABEL:   TOK_ESERR
ID:      AF1621E8

Date/Time:   Jun 20 11:28:11
Sequence Number: 17262
Machine Id:  123456789012
Node Id:     host1
Class:      H
Type:       PEND
Resource Name: TokenRing
Resource Class: tok0
Resource Type: Adapter
Location:   TokenRing

Description
EXCESSIVE TOKEN-RING ERRORS

Probable Causes
TOKEN-RING FAULT DOMAIN

Failure Causes
TOKEN-RING FAULT DOMAIN

Recommended Actions
REVIEW LINK CONFIGURATION DETAIL DATA
CONTACT TOKEN-RING ADMINISTRATOR RESPONSIBLE FOR THIS LAN

Detail Data
SENSE DATA
0ACA 0032 A440 0001 0000 0000 0000 0000 0000 0000 0000 0000 0000
0000 2080 0000 0000 0010 0000 0000 0000 0000 0000 0000 0000 0000
0000 0000 78CC 0000 0000 0005 C88F 0304 F4E0 0000 1000 5A4F 5685
1000 5A4F 5685 3030 3030 0000 0000 0000 0000 0000 0000 0000 0000
0000 0000 0000 0000 0000 0000 0000 0000 0000 0000 0000 0000 0000
0000 0000 0000 0000 0000 0000 0000 0000 0000 0000 0000 0000 0000
0000 0000 0000 0000 0000 0000

```

An error-class value of **S** and an error-type value of **PERM** indicate that the system encountered a problem with software and could not recover from it.

```

LABEL:   DSI_PROC
ID:      20FAED7F

Date/Time:   Jun 28 23:40:14
Sequence Number: 20136
Machine Id:  123456789012
Node Id:     123456789012
Class:      S
Type:       PERM
Resource Name: SYSVMM

Description
Data Storage Interrupt, Processor

Probable Causes
SOFTWARE PROGRAM

Failure Causes
SOFTWARE PROGRAM

Recommended Actions
IF PROBLEM PERSISTS THEN DO THE FOLLOWING
CONTACT APPROPRIATE SERVICE REPRESENTATIVE

Detail Data
Data Storage Interrupt Status Register
4000 0000

```

```
Data Storage Interrupt Address Register
0000 9112
Segment Register, SEGREG
D000 1018
EXVAL
0000 0005
```

An error-class value of **S** and an error-type value of **TEMP** indicate that the system encountered a problem with software. After several attempts, the system was able to recover from the problem.

```
LABEL:          SCSI_ERR6
ID:            52DB7218

Date/Time:     Jun 28 23:21:11
Sequence Number: 20114
Machine Id:    123456789012
Node Id:      host1
Class:        S
Type:         INFO
Resource Name: scsi0

Description
SOFTWARE PROGRAM ERROR

Probable Causes
SOFTWARE PROGRAM

Failure Causes
SOFTWARE PROGRAM

Recommended Actions
IF PROBLEM PERSISTS THEN DO THE FOLLOWING
CONTACT APPROPRIATE SERVICE REPRESENTATIVE

Detail Data
SENSE DATA
0000 0000 0000 0000 0000 0011 0000 0008 000E 0900 0000 0000 FFFF
FFFE 4000 1C1F 01A9 09C4 0000 000F 0000 0000 0000 0000 FFFF FFFF
0325 0018 0040 1500 0000 0000 0000 0000 0000 0000 0000 0800
0000 0100 0000 0000 0000 0000 0000 0000 0000 0000 0000 0000
0000 0000
```

An error class value of **O** indicates that an informational message has been logged.

```
LABEL:          OPMSG
ID:            AA8AB241

Date/Time:     Jul 16 03:02:02
Sequence Number: 26042
Machine Id:    123456789012
Node Id:      host1
Class:        O
Type:         INFO
Resource Name: OPERATOR

Description
OPERATOR NOTIFICATION

User Causes
errlogger COMMAND

Recommended Actions
REVIEW DETAILED DATA

Detail Data
MESSAGE FROM errlogger COMMAND
hdisk1 : Error log analysis indicates a hardware failure.
```

Example of a summary error report

The following is an example of a summary error report generated using the **errpt** command. One line of information is returned for each error entry.

```
ERROR_
IDENTIFIER  TIMESTAMP  T CL RESOURCE_NAME ERROR_DESCRIPTION
192AC071    0101000070 I O  errdemon      Error logging turned off
0E017ED1    0405131090 P H  mem2          Memory failure
```

```
9DBCDFEE 0101000070 I 0 errdemon Error logging turned on
038F2580 0405131090 U H scdisk0 UNDETERMINED ERROR
AA8AB241 0405130990 I 0 OPERATOR OPERATOR NOTIFICATION
```

Generating an error report

To create an error report of software or hardware problems do the following:

1. Determine if error logging is on or off by determining if the error log contains entries:

```
errpt -a
```

The **errpt** command generates an error report from entries in the system error log.

If the error log does not contain entries, error logging has been turned off. Activate the facility by typing:

```
/usr/lib/errdemon
```

Note: You must have root user access to run this command.

The **errdemon** daemon starts error logging and writes error log entries in the system error log. If the daemon is not running, errors are not logged.

2. Generate an error log report using the **errpt** command. For example, to see all the errors for the `hdisk1` disk drive, type:

```
errpt -N hdisk1
```

3. Generate an error log report using SMIT. For example, use the **smit errpt** command:

```
smit errpt
```

- a. Select **1** to send the error report to standard output, or select **2** to send the report to the printer.
- b. Select **yes** to display or print error log entries as they occur. Otherwise, select **no**.
- c. Specify the appropriate device name in the **Select resource names** option (such as `hdisk1`).
- d. Select **Do**.

Stopping an error log

This procedure describes how to stop the error-logging facility.

To turn off error logging, use the **errstop** command. You must have root user authority to use this command.

Ordinarily, you would not want to turn off the error-logging facility. Instead, you should clean the error log of old or unnecessary entries.

Turn off the error-logging facility when you are installing or experimenting with new software or hardware. This way the error logging daemon does not use CPU time to log problems you know you are causing.

Cleaning an error log

Error-log cleaning is normally done for you as part of the daily **cron** command. If it is not done automatically, clean the error log yourself every couple of days after you have examined the contents to make sure there are no significant errors.

You can also clean up specific errors. For example, if you get a new disk and you do not want the old disk's errors in the log, you can clean just the old disk's errors.

Delete all entries in your error log by doing either of the following:

- Use the **errclear -d** command. For example, to delete all software errors, type:

```
errclear -d S 0
```


The **errclear** command deletes entries from the error log that are older than a specified number of days. The 0 in the previous example indicates that you want to delete entries for all days.

- Use the **smit errclear** command:

```
smit errclear
```

Copying an error log to diskette or tape

Copy an error log by doing one of the following:

- To copy the error log to diskette, use the **ls** and **backup** commands. Insert a formatted diskette into the diskette drive and type:

```
ls /var/adm/ras/errlog | backup -ivp
```

- To copy the error log to tape, insert a tape in the drive and type:

```
ls /var/adm/ras/errlog | backup -ivpf/dev/rmt0
```

- To gather system configuration information in a **tar** file and copy it to diskette, use the **snap** command. Insert a formatted diskette into the diskette drive and type:

```
snap -a -o /dev/rfd0
```

Note: To use the **snap** command, you need root user authority.

The **snap** command in this example uses the **-a** flag to gather all information about your system configuration. The **-o** flag copies the compressed **tar** file to the device you name. The `/dev/rfd0` names your disk drive.

To gather all configuration information in a **tar** file and copy it to tape, type:

```
snap -a -o /dev/rmt0
```

The `/dev/rmt0` names your tape drive.

Using the liberrlog services

The **liberrlog** services allow you to read entries from an error log, and provide a limited update capability. They are especially useful from an error notification method written in the C programming language, rather than a shell script. Accessing the error log using the **liberrlog** functions is much more efficient than using the **errpt** command.

Related information

[error_open](#)

[errorlog_close](#)

[errlog_find](#), [errlog_error_sequence](#), [errlog_find_next](#)

[errlog_set_direction](#)

[errlog_write](#)

Error logging and alerts

This section describes the process of error logging and receiving alert

If the **Alert** field of an error record template is set to **True**, programs that process alerts use the following fields in the error log to build an alert:

- **Class**
- **Type**
- **Description**
- **Probable Cause**
- **User Cause**

- **Install Cause**
- **Failure Cause**
- **Recommended Action**
- **Detail Data**

These template fields must be set up according to the SNA Generic Alert Architecture described in *SNA Formats*, order number GA27-3136. You can view the book at <http://publib.boulder.ibm.com/cgi-bin/bookmgr/BOOKS/D50A5007>. Alerts that are not set up according to the architecture cannot be processed properly by a receiving program, such as NetView®.

Messages added to the error-logging message sets must not conflict with the SNA Generic Alert Architecture. When the **errmsg** command is used to add messages, the command selects message numbers that do not conflict with the architecture.

If the **Alert** field of an error record template is set to **False**, you can use any of the messages in the error-logging message catalog.

Error logging controls

To control the error-logging facility, you can use error-logging commands, subroutines and kernel services, as well as files.

Error-logging commands

errclear

Deletes entries from the error log. This command can erase the entire error log. Removes entries with specified error ID numbers, classes, or types.

errdead

Extracts errors contained in the **/dev/error** buffer captured in the system dump. The system dump will contain error records if the **errdemon** daemon was not active prior to the dump.

errdemon

Reads error records from the **/dev/error** file and writes error log entries to the system error log. The **errdemon** also performs error notification as specified in the error notification objects in the Object Data Manager (ODM). This daemon is started automatically during system initialization.

errinstall

Can be used to add or replace messages in the error message catalog. Provided for use by software installation procedures. The system creates a backup file named *File.undo*. The **undo** file allows you to cancel the changes you made by issuing the **errinstall** command.

errlogger

Writes an operator message entry to the error log.

errmsg

Implements error logging in in-house applications. The **errmsg** command lists, adds, or deletes messages stored in the error message catalog. Using this command, text can be added to the Error Description, Probable Cause, User Cause, Install Cause, Failure Cause, Recommended Action, and Detailed Data message sets.

errpt

Generates an error report from entries in the system error log. The report can be formatted as a single line of data for each entry, or the report can be a detailed listing of data associated with each entry in the error log. Entries of varying classes and types can be omitted from or included in the report.

errstop

Stops the **errdemon** daemon, which is initiated during system initialization. Running the **errstop** command also disables some diagnostic and recovery functions of the system.

errupdate

Adds or deletes templates in the Error Record Template Repository. Modifies the Alert, Log, and Report attributes of an error template. Provided for use by software installation procedures.

Error logging subroutines and kernel services

errlog

Writes an error to the error log device driver

errsave and errlast

Allows the kernel and kernel extensions to write to the error log

errlog_open

Opens an error log

errlog_close

Closes an error log

errlog_find_first

Finds the first occurrence of an error log entry

errlog_find_next

Finds the next occurrence of an error log entry

errlog_find_sequence

Finds the error log entry with the specified sequence number

errlog_set_direction

Sets the direction for the error log find functions

errlog_write

Updates an error log entry

errresume

Resumes error logging after an **errlast** command was issued.

Error logging files

/dev/error

Provides standard device driver interfaces required by the error log component

/dev/errorctl

Provides nonstandard device driver interfaces for controlling the error logging system

/usr/include/sys/err_rec.h

Contains structures defined as arguments to the **errsave** kernel service and the **errlog** subroutine

/usr/include/sys/errlog.h

Defines the interface to the **liberrlog** subroutines

/var/adm/ras/errlog

Stores instances of errors and failures encountered by the system

/var/adm/ras/errtmplt

Contains the Error Record Template Repository

File systems and logical volumes

A file is a one-dimensional array of bytes that can contain ASCII or binary information.

AIX files can contain data, shell scripts, and programs. File names are also used to represent abstract objects such as sockets or device drivers.

Files are represented internally by index nodes (*i-nodes*). Within the journaled file system (JFS), an *i-node* is a structure that contains all access, timestamp, ownership, and data location information for each file. An *i-node* is 128-bytes in JFS and 512-bytes in the enhanced journaled file system (JFS2). Pointers within the *i-node* structure designate the real disk address of the data blocks associated with the file. An *i-node* is identified by an offset number (*i-number*) and has no file name information. The connection of *i-numbers* and file names is called a *link*.

File names exist only in directories. Directories are a unique type of file that give hierarchical structure to the file system. Directories contain directory entries. Each directory entry contains a file name and an *i-number*.

JFS and JFS2 are supported by this operating system. The file system links the file and directory data to the structure used by storage and retrieval mechanisms.

Related information

[ls](#)

[mkfs](#)

[pr](#)

[fullstat.h](#)

[stat](#)

[statfs](#)

File types

A file is a one-dimensional array of bytes with at least one hard link (file name). Files can contain ASCII or binary information.

Files contain data, shell scripts, or programs. File names are also used to represent abstract objects, such as sockets, pipes, and device drivers.

The kernel does not distinguish record boundaries in regular files, so programs can establish their own boundary markers.

Files are represented in the [journalized file system](#) (JFS and JFS2) by disk index nodes (i-node). Information about the file (such as ownership, access modes, access time, data addresses, and modification time) is stored in the i-node.

The journalized file system supports the following file types:

File types supported by journalized file system

Type of file	Macro name used in mode.h	Description
Regular	S_ISREG	A sequence of bytes with one or more names. Regular files can contain ASCII or binary data. These files can be randomly accessed (read from or written to) from any byte in the file.
Directory	S_ISDIR	Contains directory entries (file name and i-number pairs). Directory formats are determined by the file system. Processes read directories as they do ordinary files, but the kernel reserves the right to write to a directory. Special sets of subroutines control directory entries.
Block Special	S_ISBLK	Associates a structured device driver with a file name.
Character Special	S_ISCHR	Associates an unstructured device driver with a file name.
Pipes	S_ISFIFO	Designates an interprocess communication (IPC) channel. The mkfifo subroutine creates named pipes. The pipe subroutine creates unnamed pipes.

Type of file	Macro name used in mode.h	Description
Symbolic Links	S_ISLNK	A file that contains either an absolute or relative path name to another file name.
Sockets	S_ISSOCK	An IPC mechanism that allows applications to exchange data. The socket subroutine creates sockets, and the bind subroutine allows sockets to be named.

The maximum size of a regular file in a JFS file system enabled for large files is slightly less than 64 gigabytes (68589453312). In other file systems that are enabled for large files and in other JFS file system types, all files not listed as regular in the previous table have a maximum file size of 2 gigabytes minus 1 (2147483647). The maximum size of a file in JFS2 is limited by the size of the file system itself.

The architectural limit on the size of a JFS2 file system is 2^{52} bytes, or 4 petabytes. The maximum file size supported by the 64-bit kernel is 2^{44} - 4096 bytes, or just less than 16 terabytes.

The maximum length of a file name is 255 characters, and the maximum length of a path name is 1023 bytes.

Working with files

The operating system provides many subroutines that manipulate files. For brief descriptions of the most common file-control subroutines, see the following:

Creating files

The following subroutines are used when creating files:

creat

Creates a new, empty, regular file

link

Creates an additional name (directory entry) for an existing file

mkdir

Creates a directory

mkfifo

Creates a named pipe

mknod

Creates a file that defines a device

open

Creates a new, empty file if the **O_CREAT** flag is set

pipe

Creates an IPC

socket

Creates a socket

Manipulating files (programming)

The following subroutines can be used to manipulate files:

access

Determines the accessibility of a file.

chmod

Changes the access modes of a file.

chown

Changes ownership of a file.

close

Closes open file descriptors (including sockets).

fcntl

Creates space in a file.

fcntl, dup, or dup2

Control open file descriptors.

fsync

Writes changes in a file to permanent storage.

ioctl

Controls functions associated with open file descriptors, including special files, sockets, and generic device support, such as the termio general terminal interface.

lockf or flock

Control open file descriptors.

lseek or llseek

Move the I/O pointer position in an open file.

open

Returns a file descriptor used by other subroutines to refer to the opened file. The **open** operation takes a regular file name and a permission mode that indicates whether the file is to be read from, written to, or both.

read

Gets data from an open file if the appropriate permissions (**O_RDONLY** or **O_RDWR**) were set by the **open** subroutine.

rename

Changes the name of a file.

rmdir

Removes directories from the file system.

stat

Reports the status of a file, including the owner and access modes.

truncate

Changes the length of a file.

write

Puts data into an open file if the appropriate permissions (**O_WRONLY** or **O_RDWR**) were set by the **open** subroutine.

For more information on types and characteristics of file systems, see [File systems](#) in *Operating system and device management*.

Working with JFS directories

Directories provide a hierarchical structure to the file system, link files, and i-node subdirectory names. There is no limit on the depth of nested directories. Disk space is allocated for directories in 4096-byte blocks, but the operating system allocates directory space in 512-byte records.

Processes can read directories as regular files. However, the kernel can write directories. For this reason, directories are created and maintained by a set of subroutines unique to them.

JFS directory structures

Directories contain a sequence of directory entries. Each directory entry contains three fixed-length fields (the index number associated with the file's i-node, the length of the file name, and the number of bytes for the entry) and one variable-length field for the file name. The file name field is null-terminated and padded to 4 bytes. File names can be up to 255 bytes long.

Directory entries are of variable length to allow file names the greatest flexibility. However, all directory space is allocated at all times.

No directory entry can span 512-byte sections of a directory. When a directory requires more than 512 bytes, another 512-byte record is appended to the original record. If all of the 512-byte records in the allocated data block are filled, an additional data block (4096 bytes) is allotted.

When a file is removed, the space that the file occupied in the directory structure is added to the preceding directory entry. The information about the removed directory remains until a new entry fits into the space vacated.

Every directory contains the entries `.` (dot) and `..` (dot, dot). The `.` (dot) directory entry points to the i-node for the directory itself. The `..` (dot, dot) directory entry points to the i-node for the parent directory. The **mkfs** program initializes a file system so that the `.` (dot) and `..` (dot, dot) entries in the new root directory point to the root i-node of the file system.

Directories have the following access modes:

Mode	Description
read	Allows a process to read directory entries
write	Allows a process to create new directory entries or remove old ones by using the creat , mknod , link , and unlink subroutines
execute	Allows a process to use the directory as a current working directory or to search below the directory in the file tree

Working with JFS directories (programming)

The following is a list of subroutines available for working with directories:

closedir

Closes a directory stream and frees the structure associated with the *DirectoryPointer* parameter

mkdir

Creates directories

opendir

Opens the directory designated by the *DirectoryName* parameter and associates a directory stream with it

readdir

Returns a pointer to the next directory entry

rewinddir

Resets the position of the specified directory stream to the beginning of the directory

rmdir

Removes directories

seekdir

Sets the position of the next **readdir** subroutine operation on the directory stream

telldir

Returns the current location associated with the specified directory stream

Changing the current directory of a process

When the system is booted, the first process uses the root directory of the root file system as its current directory. New processes created with the **fork** subroutine inherit the current directory used by the parent process. The **chdir** subroutine changes the current directory of a process.

The **chdir** subroutine parses the path name to ensure that the target file is a directory and that the process owner has permissions to the directory. After the **chdir** subroutine is run, the process uses the new current directory to search all path names that do not begin with a / (slash).

Changing the root directory of a process

You can cause the directory named by a process *Path* parameter to become the effective root directory by using the **chroot** subroutine. Child processes of the calling process consider the directory indicated by the **chroot** subroutine as the logical root directory of the file system.

Processes use the global file system root directory for all path names starting with a / (slash). All path name searches beginning with a / (slash) begin at this new root directory.

Subroutines that control JFS directories

Due to the unique nature of directory files, directories are controlled by a special set of subroutines. The following subroutines are designed to control directories:

chdir

Changes the current working directory

chroot

Changes the effective root directory

getcwd or getwd

Gets path to current directory

mkdir

Creates a directory

opendir, readdir, telldir, seekdir, rewinddir, or closedir

Perform various actions on directories

rename

Renames a directory

rmdir

Removes a directory

Working with JFS2 directories

Directories provide a hierarchical structure to the file system, link files, and i-node subdirectory names. There is no limit on the depth of nested directories.

Disk space is allocated for directories in file system blocks.

Processes can read directories as regular files. However, the kernel can write directories. For this reason, directories are created and maintained by a set of subroutines unique to them.

JFS2 directory structures

A directory contains entries that describe the objects contained in the directory. A directory entry has a fixed length and contains the following:

- The i-node number
- The name (up to 22 bytes long)
- A name length field
- A field to continue the entry if the name is longer than 22 bytes

The directory entries are stored in a B+ tree sorted by name. The self (.) and parent (..) information is contained in the i-node instead of in a directory entry.

Directories have the following access modes:

Mode	Description
read	Allows a process to read directory entries
write	Allows a process to create new directory entries or remove old ones, by using the creat , mknod , link , and unlink subroutines

Mode	Description
execute	Allows a process to use the directory as a current working directory or to search below the directory in the file tree

Working with JFS2 directories (programming)

The following is a list of subroutines available for working with directories:

closedir

Closes a directory stream and frees the structure associated with the *DirectoryPointer* parameter

mkdir

Creates directories

opendir

Returns a structure pointer that is used by the **readdir** subroutine to obtain the next directory entry, by **rewinddir** to reset the read position to the beginning, and by **closedir** to close the directory.

readdir

Returns a pointer to the next directory entry

rewinddir

Resets the position of the specified directory stream to the beginning of the directory

rmdir

Removes directories

seekdir

Returns to a position previously obtained with the **telldir** subroutine

telldir

Returns the current location associated with the specified directory stream

Do not use the **open**, **read**, **lseek**, and **close** subroutines to access directories.

Changing current directory of a process

When the system is booted, the first process uses the root directory of the root file system as its current directory. New processes created with the **fork** subroutine inherit the current directory used by the parent process. The **chdir** subroutine changes the current directory of a process.

The **chdir** subroutine parses the path name to ensure that the target file is a directory and that the process owner has permissions to the directory. After the **chdir** subroutine is run, the process uses the new current directory to search all path names that do not begin with a / (slash).

Changing the root directory of a process

Processes can change their understanding of the root directory through the **chroot** subroutine. Child processes of the calling process consider the directory indicated by the **chroot** subroutine as the logical root directory of the file system.

Processes use the global file system root directory for all path names starting with a / (slash). All path name searches beginning with a / (slash) begin at this new root directory.

Subroutines that control JFS2 directories

Due to the unique nature of directory files, directories are controlled by a special set of subroutines. The following subroutines are designed to control directories:

chdir

Changes the current working directory

chroot

Changes the effective root directory

opendir, readdir, telldir, seekdir, rewinddir, or closedir

Perform various actions on directories

getcwd or getwd

Gets path to current directory

mkdir

Creates a directory

rename

Renames a directory

rmdir

Removes a directory

Working with JFS i-nodes

Files in the journaled file system (JFS) are represented internally as index nodes (i-nodes). JFS i-nodes exist in a static form on disk and contain access information for the file, as well as pointers to the real disk addresses of the file's data blocks.

The number of disk i-nodes available to a file system is dependent on the size of the file system, the allocation group size (8 MB by default), and the number of bytes per i-node ratio (4096 by default). These parameters are given to the **mkfs** command at file system creation. When enough files have been created to use all the available i-nodes, no more files can be created, even if the file system has free space.

To determine the number of available i-nodes, use the **df -v** command. Disk i-nodes are defined in the **/usr/include/jfs/ino.h** file.

Disk i-node structure for JFS

Each disk i-node in JFS is a 128-byte structure. The offset of a particular i-node within the i-node list of the file system produces the unique number (i-number) by which the operating system identifies the i-node. A bit map, known as the *i-node map*, tracks the availability of free disk i-nodes for the file system.

Disk i-nodes include the following information:

Field	Contents
i_mode	Type of file and access permission mode bits
i_size	Size of file in bytes
i_uid	Access permissions for the user ID
i_gid	Access permissions for the group ID
i_nblocks	Number of blocks allocated to the file
i_mtime	Last time the file was modified
i_atime	Last time the file was accessed
i_ctime	Last time the i-node was modified
i_nlink	Number of hard links to the file
i_rdaddr[8]	Real disk addresses of the data
i_rindirect	Real disk address of the indirect block, if any

You cannot change file data without changing the i-node, but it is possible to change the i-node without changing the contents of the file. For example, when permission is changed, the information within the i-node (**i_mode**) is modified, but the data in the file remains the same.

The **i_rdaddr** field within the disk i-node contains 8 disk addresses. These addresses point to the first 8 data blocks assigned to the file. The **i_rindirect** field address points to an indirect block. Indirect blocks are either single indirect or double indirect. Thus, there are three possible geometries of block allocation for a file: direct, indirect, or double indirect.

Disk i-nodes do not contain file or path name information. Directory entries are used to link file names to i-nodes. Any i-node can be linked to many file names by creating additional directory entries with the **link** or **symlink** subroutine. To determine the i-node number assigned to a file, use the **ls -i** command.

The i-nodes that represent files that define devices contain slightly different information from i-nodes for regular files. Files associated with devices are called *special files*. There are no data block addresses in special device files, but the major and minor device numbers are included in the **i_rdev** field.

A disk i-node is released when the link count (**i_nlink**) to the i-node equals 0. Links represent the file names associated with the i-node. When the link count to the disk i-node is 0, all the data blocks associated with the i-node are released to the bit map of free data blocks for the file system. The i-node is then placed on the free i-node map.

JFS In-core i-node structure

When a file is opened, the operating system creates an in-core i-node. The *in-core i-node* contains a copy of all the fields defined in the disk i-node, plus additional fields for tracking and managing access to the in-core i-node. When a file is opened, the information in the disk i-node is copied into an in-core i-node for easier access. In-core i-nodes are defined in the **/usr/include/jfs/inode.h** file. Some of the additional information tracked by the in-core i-node is as follows:

- Status of the in-core i-node, including flags that indicate:
 - An i-node lock
 - A process waiting for the i-node to unlock
 - Changes to the file's i-node information
 - Changes to the file's data
- Logical device number of the file system that contains the file
- i-number used to identify the i-node
- Reference count. When the reference count field equals 0, the in-core i-node is released.

When an in-core i-node is released (for example, with the **close** subroutine), the in-core i-node reference count is reduced by 1. If this reduction results in the reference count to the in-core i-node becoming 0, the i-node is released from the in-core i-node table, and the contents of the in-core i-node are written to the disk copy of the i-node (if the two versions differ).

Working with JFS2 i-nodes

Files in JFS2 are represented internally as index nodes (i-nodes).

JFS2 disk i-nodes exist in a static form on the disk and contain access information for the files, as well as pointers to the real disk addresses of the file's data blocks. The i-nodes are allocated dynamically by JFS2. Disk-inodes are defined in the **/usr/include/j2/j2_dinode.h** file.

When a file is opened, an in-core i-node is created by the operating system. The in-core i-node contains a copy of all the fields defined in the disk i-node, plus additional fields for tracking the in-core i-node. In-core i-nodes are defined in the **/usr/include/j2/j2_inode.h** file.

Disk i-node structure for JFS2

Each disk i-node in JFS2 is a 512-byte structure. The index of a particular i-node allocation map of the file system produces the unique number (i-number) by which the operating system identifies the i-node. The i-node allocation map tracks the location of the i-nodes on the disk, as well as their availability.

Disk i-nodes include the following information:

Field	Contents
di_mode	Type of file and access permission mode bits
di_size	Size of file in bytes
di_uid	Access permissions for the user ID

Field	Contents
di_gid	Access permissions for the group ID
di_nblocks	Number of blocks allocated to the file
di_mtime	Last time the file was modified
di_atime	Last time the file was accessed
di_ctime	Last time the i-node was modified
di_nlink	Number of hard links to the file
di_broot	Root of B+ tree describing the disk addresses of the data

You cannot change the file data without changing the i-node, but it is possible to change the i-node without changing the contents of the file. For example, when permission is changed, the information within the i-node (**di_mode**) is modified, but the data in the file remains the same.

The **di_broot** describes the root of the B+ tree. It describes the data for the i-node. **di_broot** has a field indicating how many of its entries in the i-node are being used and another field describing whether they are leaf nodes or internal nodes for the B+ tree.

Disk i-nodes do not contain file or path name information. Directory entries are used to link file names to i-nodes. Any i-node can be linked to many file names by creating additional directory entries with the **link** or **symlink** subroutine. To determine the i-node number assigned to a file, use the **ls -i** command.

The i-nodes that represent files that define devices contain slightly different information from i-nodes for regular files. Files associated with devices are called *special files*. There are no data block addresses in special device files, but the major and minor device numbers are included in the **di_rdev** field.

A disk i-node is released when the link count (**di_nlink**) to the i-node equals 0. Links represent the file names associated with the i-node. When the link count to the disk i-node is 0, all the data blocks associated with the i-node are released to the bitmap of free data blocks for the file system. The i-node is then placed on the free i-node map.

JFS2 in-core i-node structure

When a file is opened, the information in the disk i-node is copied into an in-core i-node for easier access. The in-core i-node structure contains additional fields that manage access to the disk i-node's valuable data. The fields of the in-core i-node are defined in the **j2_inode.h** file. Some of the additional information tracked by the in-core i-node is as follows:

- Status of the in-core i-node, including flags that indicate:
 - An i-node lock
 - A process waiting for the i-node to unlock
 - Changes to the file's i-node information
 - Changes to the file's data
- Logical device number of the file system that contains the file
- i-number used to identify the i-node
- Reference count. When the reference count field equals 0, the in-core i-node is released.

When an in-core i-node is released (for example, with the **close** subroutine), the in-core i-node reference count is reduced by 1. If this reduction results in the reference count to the in-core i-node becoming 0, the i-node is released from the in-core i-node table, and the contents of the in-core i-node are written to the disk copy of the i-node (if the two versions differ).

Allocating JFS file space

File space allocation is the method by which data is apportioned physical storage space in the operating system.

The kernel allocates disk space to a file or directory in the form of logical blocks. A *logical block* for JFS refers to the division of a file or directory's contents into 4096-byte units. Logical blocks are not tangible entities; however, the data in a logical block consumes physical storage space on the disk. Each file or directory consists of 0 or more logical blocks. Fragments, instead of logical blocks, are the basic units for allocated disk space in JFS.

Full and partial logical blocks

A file or directory may contain full or partial logical blocks. A full logical block contains 4096 bytes of data. Partial logical blocks occur when the last logical block of a file or directory contains less than 4096 bytes of data.

For example, a file of 8192 bytes is two logical blocks. The first 4096 bytes reside in the first logical block and the following 4096 bytes reside in the second logical block. Likewise, a file of 4608 bytes consists of two logical blocks. However, the last logical block is a partial logical block, containing the last 512 bytes of the file's data. Only the last logical block of a file can be a partial logical block.

Allocation in fragmented file systems

The default fragment size is 4096 bytes. You can specify smaller fragment sizes with the **mkfs** command during a file system's creation. Allowable fragment sizes are: 512, 1024, 2048, and 4096 bytes. You can use only one fragment size in a file system.

To maintain efficiency in file system operations, the JFS allocates 4096 bytes of fragment space to files and directories that are 32 KB or larger in size. A fragment that covers 4096 bytes of disk space is allocated to a full logical block. When data is added to a file or directory, the kernel allocates disk fragments to store the logical blocks. Thus, if the file system's fragment size is 512 bytes, a full logical block is the allocation of eight fragments.

The kernel allocates disk space so that only the last bytes of data receive a partial block allocation. As the partial block grows beyond the limits of its current allocation, additional fragments are allocated. If the partial block increases to 4096 bytes, the data stored in its fragments reallocated into 4096 file-system block allocations. A partial logical block that contains less than 4096 bytes of data is allocated the number of fragments that best matches its storage requirements.

Block reallocation also occurs if data is added to logical blocks that represent file holes. A *file hole* is an empty logical block located prior to the last logical block that stores data. (File holes do not occur within directories.) These empty logical blocks are not allocated fragments. However, as data is added to file holes, allocation occurs. Each logical block that was not previously allocated disk space is allocated 4096 bytes of fragment space.

Additional block allocation is not required if existing data in the middle of a file or directory is overwritten. The logical block containing the existing data has already been allocated fragments.

JFS tries to maintain contiguous allocation of a file or directory's logical blocks on the disk. Maintaining contiguous allocation lessens seek time because the data for a file or directory can be accessed sequentially and found on the same area of the disk. However, disk fragments for one logical block are not always contiguous to the disk fragments for another logical block. The disk space required for contiguous allocation may not be available if it has already been written to by another file or directory. An allocation for a single logical block, however, always contains contiguous fragments.

The file system uses a bitmap called the *block allocation map* to record the status of every block in the file system. When the file system needs to allocate a new fragment, it refers to the fragment allocation map to identify which fragments are available. A fragment can only be allocated to a single file or directory at a time.

Allocation in compressed JFS file systems

In a file system that supports data compression, directories are allocated disk space. Data compression also applies to regular files and symbolic links whose size is larger than that of their i-nodes.

The allocation of disk space for compressed file systems is the same as that of fragments in fragmented file systems. A logical block is allocated 4096 bytes when it is modified. This allocation guarantees that there will be a place to store the logical block if the data does not compress. The system requires that a write or store operation report an out-of-disk-space condition into a memory-mapped file at a logical block's initial modification. After modification is complete, the logical block is compressed before it is written to a disk. The compressed logical block is then allocated only the number of fragments required for its storage.

In a fragmented file system, only the last logical block of a file (not larger than 32 KB) can be allocated less than 4096 bytes. The logical block becomes a partial logical block. In a compressed file system, every logical block can be allocated less than a full block.

A logical block is no longer considered modified after it is written to a disk. Each time a logical block is modified, a full disk block is allocated again, according to the system requirements. Reallocation of the initial full block occurs when the logical block of compressed data is successfully written to a disk.

Allocation in JFS file systems enabled for large files

In a file system enabled for large files, the JFS allocates two sizes of fragments for regular files. A "large" fragment (32 X 4096) is allocated for logical blocks after the 4 MB boundary, and a 4096 bytes fragment is allocated for logical blocks before the 4 MB boundary. All nonregular files allocate 4096 bytes fragments. This geometry allows a maximum file size of slightly less than 64 gigabytes (68589453312).

A *large* fragment is made up of 32 contiguous 4096 bytes fragments. Because of this requirement, it is recommended that file systems enabled for large files have predominantly large files in them. Storing many small files (files less than 4 MB) can cause free-space fragmentation problems. This can cause large allocations to fail with an ENOSPC error condition because the file system does not contain 32 contiguous disk addresses.

Disk address format

JFS fragment support requires fragment-level addressability. As a result, disk addresses have a special format for mapping where the fragments of a logical block reside on the disk. Fragmented and compressed file systems use the same method for representing disk addresses. Disk addresses are contained in the **i_rdaddr** field of the i-nodes or in the indirect blocks. All fragments referenced in a single address must be contiguous on the disk.

The disk address format consists of the **nfrags** and **addr** fields. These fields describe the area of disk covered by the address:

addr

Indicates which fragment on the disk is the starting fragment

nfrags

Indicates the total number of contiguous fragments not used by the address

For example, if the fragment size for the file system is 512 bytes and the logical block is divided into eight fragments, the **nfrags** value is 3, indicating that five fragments are included in the address.

The following examples illustrate possible values for the **addr** and **nfrags** fields for different disk addresses. These values assume a fragment size of 512 bytes, indicating that the logical block is divided into eight fragments.

Address for a single fragment:

```
addr: 143
```

```
nfrags: 7
```

This address indicates that the starting location of the data is fragment 143 on the disk. The **nfrags** value indicates that the total number of fragments included in the address is one. The **nfrags** value changes in a file system that has a fragment size other than 512 bytes. To correctly read the **nfrags** value, the system, or any user examining the address, must know the fragment size of the file system.

Address for five fragments:

```
addr: 1117
```

```
nfrags: 3
```

In this case, the address starts at fragment number 1117 on the disk and continues for five fragments (including the starting fragment). Three fragments are remaining, as illustrated by the **nfrags** value.

The disk addresses are 32 bits in size. The bits are numbered from 0 to 31. The 0 bit is always reserved. Bits 1 through 3 contain the **nfrags** field. Bits 4 through 31 contain the **addr** field.

JFS indirect blocks

The JFS uses the indirect blocks to address the disk space allocated to larger files. Indirect blocks allow the greatest flexibility for file sizes and the fastest retrieval time. The indirect block is assigned using the **i_rindirect** field of the disk i-node. This field allows for the following geometries or methods for addressing the disk space:

- Direct
- Single indirect
- Double indirect

Each of these methods uses the same disk address format as compressed and fragmented file systems. Because files larger than 32 KB are allocated fragments of 4096 bytes, the **nfrags** field for addresses using the single indirect or double indirect method has a value of 0.

Direct method

When the direct method of disk addressing is used, each of the eight addresses listed in the **i_rdaddr** field of the disk i-node points directly to a single allocation of disk fragments. The maximum size of a file using direct geometry is 32,768 bytes (32KB), or 8 x 4096 bytes. When the file requires more than 32 KB, an indirect block is used to address the file's disk space.

Single indirect method

The **i_rindirect** field contains an address that points to either a single indirect block or a double indirect block. When the single indirect disk-addressing method is used, the **i_rindirect** field contains the address of an indirect block containing 1024 addresses. These addresses point to the disk fragments for each allocation. Using the single indirect block geometry, the file can be up to 4,194,304 bytes (4 MB), or 1024 x 4096 bytes.

Double indirect method

The double indirect disk-addressing method uses the **i_rindirect** field to point to a double indirect block. The double indirect block contains 512 addresses that point to indirect blocks, which contain pointers to the fragment allocations. The largest file size that can be used with the double indirect geometry in a file system not enabled for large files is 2,147,483,648 bytes (2 GB), or 512(1024 x 4096) bytes.

Note: The maximum file size that the **read** and **write** system calls would allow is 2 GB minus 1 ($2^{31}-1$). When memory map interface is used, 2 GB can be addressed.

File systems enabled for large files allow a maximum file size of slightly less than 64 gigabytes (68589453312). The first single indirect block contains 4096 byte fragments, and all subsequent single indirect blocks contain (32 X 4096) byte fragments. The following produces the maximum file size for file systems enabling large files:

```
(1 * (1024 * 4096)) + (511 * (1024 * 131072))
```

The fragment allocation assigned to a directory is divided into records of 512 bytes each and grows in accordance with the allocation of these records.

Quotas

Disk quotas restrict the amount of file system space that any single user or group can monopolize.

The **quotactl** subroutine sets limits on both the number of files and the number of disk blocks allocated to each user or group on a file system. Quotas enforce the following types of limits:

hard

Maximum limit allowed. When a process hits its hard limit, requests for more space fail.

soft

Practical limit. If a process hits the soft limit, a warning is printed to the user's terminal. The warning is often displayed at login. If the user fails to correct the problem after several login sessions, the soft limit can become a hard limit.

System warnings are designed to encourage users to heed the soft limit. However, the quota system allows processes access to the higher hard limit when more resources are temporarily required.

Allocating JFS2 file space

File space allocation is the method by which data is apportioned physical storage space in the operating system.

The kernel allocates disk space to a file or directory in the form of *logical blocks*. A logical block refers to the division of a file or directory contents into 512, 1024, 2048, or 4096 byte units. When a JFS2 file system is created, the logical block size is specified to be one of 512, 1024, 2048, or 4096 bytes. Logical blocks are not tangible entities; however, the data in a logical block consumes physical storage space on the disk. Each file or directory consists of 0 or more logical blocks.

Full and partial logical blocks

A file or directory may contain full or partial logical blocks. A full logical block contains 512, 1024, 2048, or 4096 bytes of data, depending on the file system block size specified when the JFS2 file system was created. Partial logical blocks occur when the last logical block of a file or directory contains less than a file-system block size of data.

For example, a JFS2 file system with a logical block size of 4096 with a file of 8192 bytes is two logical blocks. The first 4096 bytes reside in the first logical block and the following 4096 bytes reside in the second logical block. Likewise, a file of 4608 bytes consists of two logical blocks. However, the last logical block is a partial logical block containing the last 512 bytes of the file's data.

JFS2 File space allocation

The default block size is 4096 bytes. You can specify smaller block sizes with the **mkfs** command during a file system's creation. Allowable block sizes are 512, 1024, 2048, and 4096 bytes. You can use only one block size in a file system.

The kernel allocates disk space so only the last file system block of data receives a partial block allocation. As the partial block grows beyond the limits of its current allocation, additional blocks are allocated.

Block reallocation also occurs if data is added to logical blocks that represent file holes. A *file hole* is an empty logical block located prior to the last logical block that stores data. (File holes do not occur within directories.) These empty logical blocks are not allocated blocks. However, as data is added to file holes, allocation occurs. Each logical block that was not previously allocated disk space is allocated a file system block of space.

Additional block allocation is not required if existing data in the middle of a file or a directory is overwritten. The logical block containing the existing data has already been allocated file system blocks.

JFS2 tries to maintain contiguous allocation of a file or directory's logical blocks on the disk. Maintaining contiguous allocation lessens seek time because the data for a file or directory can be accessed sequentially and found on the same area of the disk. The disk space required for contiguous allocation might not be available if another file or directory has already written to it.

The file system uses a bitmap called the *block allocation map* to record the status of every block in the file system. When the file system needs to allocate a new block, it refers to the block allocation map to identify which blocks are available. A block can only be allocated to a single file or directory at a time.

Extents

An extent is a contiguous variable-length sequence of file system blocks allocated to a JFS2 object as a unit. Large extents may span multiple allocation groups.

An i-node represents every JFS2 object. I-nodes contain the expected object-specific information such as time stamps or file type (regular or directory, and so on). They also contain a B+ tree to record the allocation of extents.

The length and address values are necessary to define an extent. The length is measured in units of the file system block size. A 24-bit value represents the length of an extent, so an extent can range in size from 1 to $2^{24} - 1$ file system blocks. Therefore, the size of the maximum extent depends on the file system block size. The address is the address of the first block of the extent. The address is also in units of file system blocks; it is the block offset from the beginning of the file system.

An extent-based file system combined with user-specified file system block size allows JFS2 to not have separate support for internal fragmentation. You can configure the file system with a small file system block size, such as 512 bytes, to minimize internal fragmentation for file systems with large numbers of small-sized files.

In general, the allocation policy for JFS2 tries to maximize contiguous allocation by allowing a minimum number of extents, with each extent as large and contiguous as possible. This allows for larger I/O transfer, resulting in improved performance. However, in some cases, this is not always possible.

B+ trees

The B+ tree data structure is used for file layout. The most common operations that JFS2 performs are reading and writing extents. B+ trees are used to help with performance of these operations.

An extent allocation descriptor (**xad_t** structure) describes the extent and adds two more fields that are needed for representing files: an offset field, describing the logical byte address the extent represents, and a flags field. The **xad_t** structure is defined in the `/usr/include/j2/j2_xtree.h` file.

An **xad** structure describes two abstract ranges:

- The physical range of disk blocks. This starts at file system block number `addressXAD(xadp)` address and extends for `lengthXAD(xadp)` file system blocks.
- The logical range of bytes within a file. This starts at byte number `offsetXAD(xadp)*(file system block size)` and extends for `lengthXAD(xadp)*(file system block size.)`

The physical range and the logical range are both the same number of bytes in length. Note that offset is stored in units of file system block size (for example, a value of 3) in offset means 3 file system blocks, not 3 bytes. Extents within a file are always aligned on file system block size boundaries.

JFS2 limitation

JFS2 requires contiguous free space of at least a page, or 4 KB, in length when extending files. If you do not have contiguous free space of at least 4 KB, then the file system does not allow the extension of the file, even if there is enough available storage space in smaller blocks.

JFS file system layout

A file system is a set of files, directories, and other structures.

File systems maintain information and identify where a file or directory's data is located on the disk. In addition to files and directories, JFS file systems contain a boot block, a superblock, bitmaps, and one or more allocation groups. Each file system occupies one logical volume.

JFS boot block

The boot block occupies the first 4096 bytes of the file system, starting at byte offset 0 on the disk. The boot block is available to start the operating system.

JFS superblock

The superblock is 4096 bytes in size and starts at byte offset 4096 on the disk. The superblock maintains information about the entire file system and includes the following fields:

- Size of the file system
- Number of data blocks in the file system
- A flag indicating the state of the file system
- Allocation group sizes

JFS allocation bitmaps

The file system contains the following allocation bitmaps:

- The fragment allocation map records the allocation state of each fragment.
- The disk i-node allocation map records the status of each i-node.

JFS fragments

Many file systems have disk blocks or data blocks. These blocks divide the disk into units of equal size to store the data in a file or directory's logical blocks. The disk block may be further divided into fixed-size allocation units called *fragments*. Some systems do not allow fragment allocations to span the boundaries of the disk block. In other words, a logical block cannot be allocated fragments from different disk blocks.

The journaled file system (JFS), however, provides a view of the file system as a contiguous series of fragments. JFS fragments are the basic allocation unit and the disk is addressed at the fragment level. Thus, fragment allocations can span the boundaries of what might otherwise be a disk block. The default JFS fragment size is 4096 bytes, although you can specify smaller sizes. In addition to containing data for files and directories, fragments also contain disk addresses and data for indirect blocks.

JFS allocation groups

The set of fragments making up the file system are divided into one or more fixed-sized units of contiguous fragments. Each unit is an *allocation group*. The first of these groups begins the file system and contains a reserved area occupying the first 32 x 4096 bytes of the group. The first 4096 bytes of this area hold the boot block, and the second 4096 bytes hold the file system superblock.

Each allocation group contains a static number of contiguous disk i-nodes that occupy some of the group's fragments. These fragments are reserved for the i-nodes at file-system creation and extension time. For the first allocation group, the disk i-nodes occupy the fragments immediately following the reserved block area. For subsequent groups, the disk i-nodes are found at the start of each group. Disk i-nodes are 128 bytes in size and are identified by a unique disk i-node number or i-number. The i-number maps a disk i-node to its location on the disk or to an i-node within its allocation group.

A file system's fragment allocation group size and the disk i-node allocation group size are specified as the number of fragments and disk i-nodes that exist in each allocation group. The default allocation group size is 8 MB, but can be as large as 64 MB. These values are stored in the file system superblock, and they are set at file system creation.

Allocation groups allow the JFS resource allocation policies to use effective methods for to achieve optimum file system I/O performance. These allocation policies try to cluster disk blocks and disk i-nodes for related data to achieve good locality for the disk. Files are often read and written sequentially, and files within a directory are often accessed together. Also, these allocation policies try to distribute unrelated data throughout the file system in an attempt to minimize free-space fragmentation.

JFS disk i-nodes

A logical block contains a file or directory's data in units of 4096 bytes. Each logical block is allocated fragments for the storage of its data. Each file and directory has an i-node that contains access information such as file type, access permissions, owner's ID, and number of links to that file. These i-nodes also contain "addresses" for finding the location on the disk where the data for a logical block is stored.

Each i-node has an array of numbered sections. Each section contains an address for one of the file or directory's logical blocks. These addresses indicate the starting fragment and the total number of

fragments included in a single allocation. For example, a file with a size of 4096 bytes has a single address on the i-node's array. Its 4096 bytes of data are contained in a single logical block. A larger file with a size of 6144 bytes has two addresses. One address contains the first 4096 bytes and a second address contains the remaining 2048 bytes (a partial logical block). If a file has a large number of logical blocks, the i-node does not contain the disk addresses. Instead, the i-node points to an indirect block that contains the additional addresses.

JFS2 file system layout

A file system is a set of files, directories and other structures.

The file systems maintain information and identify where the data is located on the disk for a file or directory. In addition to files and directories a JFS2 file system contains a superblock, allocation maps and one or more allocation groups. Each file system occupies one logical volume.

JFS2 superblock

The superblock is 4096 bytes in size and starts at byte offset 32768 on the disk. The superblock maintains information about the entire file system and includes the following fields:

- Size of the file system
- Number of data blocks in the file system
- A flag indicating the state of the file system
- Allocation group sizes
- File system block size

JFS2 allocation maps

The file system contains the following allocation maps:

- The i-node allocation map records the location and allocation of all i-nodes in the file system.
- The block allocation map records the allocation state of each file system block.

JFS2 disk i-nodes

A logical block contains a file or directory's data in units of file system blocks. Each logical block is allocated file system blocks for the storage of its data. Each file and directory has an i-node that contains access information such as file type, access permissions, owner's ID, and number of links to that file. These i-nodes also contain a B+ tree for finding the location on the disk where the data for a logical block is stored.

JFS2 allocation groups

Allocation groups divide the space on a file system into chunks. Allocation groups are used only for a problem-solving technique in which the most appropriate solution, found by attempting alternative methods, is selected at successive stages of a program for using in the next step of the program. Allocation groups allow JFS2 resource-allocation policies to use well-known methods for achieving optimum I/O performance. First, the allocation policies try to cluster disk blocks and disk i-nodes for related data to achieve good locality for the disk. Files are often read and written sequentially and the files within a directory are often accessed together. Second, the allocation policies try to distribute unrelated data throughout the file system in order to accommodate disk locality.

Allocation groups within a file system are identified by a zero-based allocation group index, the allocation group number.

Allocation group sizes must be selected that yield allocation groups that are large enough to provide for contiguous resource allocation over time. Allocation groups are limited to a maximum number of 128 groups. The minimum allocation group size is 8192 file-system blocks.

Partial allocation groups

A file system whose size is not a multiple of the allocation group size will contain a partial allocation group; the last allocation group of the file system is not fully covered by disk blocks. This partial allocation

group will be treated as a complete allocation group, except that the nonexistent disk blocks will be marked as allocated in the block allocation map.

Writing programs that access large files

AIX supports files that are larger than 2 gigabytes (2 GB). This section assists programmers in understanding the implications of large files on their applications and to assist them in modifying their applications. Application programs can be modified, through programming interfaces, to be aware of large files. The file system programming interfaces generally are based on the **off_t** data type.

Implications for existing programs

The 32-bit application environment that all applications used prior to AIX 4.2 remains unchanged. However, existing application programs cannot handle large files.

For example, the **st_size** field in the **stat** structure, which is used to return file sizes, is a signed, 32-bit long. Therefore, that **stat** structure cannot be used to return file sizes that are larger than **LONG_MAX**. If an application attempts to use the **stat** subroutine with a file that is larger than **LONG_MAX**, the **stat** subroutine will fail, and **errno** will be set to **EOverflow**, indicating that the file size overflows the size field of the structure being used by the program.

This behavior is significant because existing programs that might not appear to have any impacts as a result of large files will experience failures in the presence of large files even though the file size is irrelevant.

The **errno EOverflow** can also be returned by an **lseek** pointer and by the **fcntl** subroutine if the values that need to be returned are larger than the data type or structure that the program is using. For **lseek**, if the resulting offset is larger than **LONG_MAX**, **lseek** will fail and **errno** will be set to **EOverflow**. For the **fcntl** subroutine, if the caller uses **F_GETLK** and the blocking lock's starting offset or length is larger than **LONG_MAX**, the **fcntl** call will fail, and **errno** will be set to **EOverflow**.

Open protection

Many existing application programs could have unexpected behavior, including data corruption, if allowed to operate on large files. AIX uses an open-protection scheme to protect applications from this class of failure.

In addition to open protection, a number of other subroutines offer protection by providing an execution environment, which is identical to the environment under which these programs were developed. If an application uses the **write** family of subroutines and the **write** request crosses the 2 GB boundary, the **write** subroutines will transfer data only up to 2 GB minus 1. If the application attempts to write at or beyond the 2GB -1 boundary, the **write** subroutines will fail and set **errno** to **EFBIG**. The behavior of the **mmap**, **ftruncate**, and **fclear** subroutines are similar.

The **read** family of subroutines also participates in the open-protection scheme. If an application attempts to read a file across the 2 GB threshold, only the data up to 2 GB minus 1 will be read. Reads at or beyond the 2GB -1 boundary will fail, and **errno** will be set to **EOverflow**.

Open protection is implemented by a flag associated with an open file description. The current state of the flag can be queried with the **fcntl** subroutine using the **F_GETFL** command. The flag can be modified with the **fcntl** subroutine using the **F_SETFL** command.

Because open file descriptions are inherited across the **exec** family of subroutines, application programs that pass file descriptors that are enabled for large-file access to other programs should consider whether the receiving program can safely access the large file.

Porting applications to the large file environment

AIX provides two methods for applications to be enabled for large-file access. Application programmers must decide which approach best suits their needs:

- Define **_LARGE_FILES**, which carefully redefines all of the relevant data types, structures, and subroutine names to their large-file enabled counterparts. Defining **_LARGE_FILES** has the advantage of maximizing application portability to other platforms because the application is still written to the

normal POSIX and XPG interfaces. It has the disadvantage of creating some ambiguity in the code because the size of the various data items cannot be determined from looking at the code.

- Recode the application to explicitly call the large-file enabled subroutines. Recoding the application has the disadvantages of requiring more effort and reducing application portability. It can be used when the redefinition effect of `_LARGE_FILES` would have a considerable negative impact on the program or when it is desirable to convert only a very small portion of the program.

In either case, the application program *must* be carefully audited to ensure correct behavior in the new environment.

Using `_LARGE_FILES`

In the default compilation environment, the `off_t` data type is defined as a signed, 32-bit long. If the application defines `_LARGE_FILES` before the inclusion of any header files, then the large-file programming environment is enabled and `off_t` is defined to be a signed, 64-bit long long. In addition, all of the subroutines that deal with file sizes or file offsets are redefined to be their large-file enabled counterparts. Similarly, all of the data structures with embedded file sizes or offsets are redefined.

The following table shows the redefinitions that occur in the `_LARGE_FILES` environment:

Entity	Redefined as	Header file
<code>off_t</code> Object	<code>long long</code>	<code><sys/types.h></code>
<code>fpos_t</code> Object	<code>long long</code>	<code><sys/types.h></code>
<code>struct stat</code> Structure	<code>struct stat64</code>	<code><sys/stat.h></code>
<code>stat</code> Subroutine	<code>stat64()</code>	<code><sys/stat.h></code>
<code>fstat</code> Subroutine	<code>fstat64()</code>	<code><sys/stat.h></code>
<code>lstat</code> Subroutine	<code>lstat64()</code>	<code><sys/stat.h></code>
<code>mmap</code> Subroutine	<code>mmap64()</code>	<code><sys/mman.h></code>
<code>lockf</code> Subroutine	<code>lockf64()</code>	<code><sys/lockf.h></code>
<code>struct flock</code> Structure	<code>struct flock64</code>	<code><sys/flock.h></code>
<code>open</code> Subroutine	<code>open64()</code>	<code><fcntl.h></code>
<code>creat</code> Subroutine	<code>creat64()</code>	<code><fcntl.h></code>
<code>F_GETLK</code> Command Parameter	<code>F_GETLK64</code>	<code><fcntl.h></code>
<code>F_SETLK</code> Command Parameter	<code>F_SETLK64</code>	<code><fcntl.h></code>
<code>F_SETLKW</code> Command Parameter	<code>F_SETLKW64</code>	<code><fcntl.h></code>
<code>ftw</code> Subroutine	<code>ftw64()</code>	<code><ftw.h></code>
<code>nftw</code> Subroutine	<code>nftw64()</code>	<code><ftw.h></code>
<code>fseeko</code> Subroutine	<code>fseeko64()</code>	<code><stdio.h></code>
<code>ftello</code> Subroutine	<code>ftello64()</code>	<code><stdio.h></code>
<code>fgetpos</code> ubroutine	<code>fgetpos64()</code>	<code><stdio.h></code>
<code>fsetpos</code> Subroutine	<code>fsetpos64()</code>	<code><stdio.h></code>
<code>fopen</code> Subroutine	<code>fopen64()</code>	<code><stdio.h></code>
<code>freopen</code> Subroutine	<code>freopen64()</code>	<code><stdio.h></code>

Entity	Redefined as	Header file
lseek Subroutine	lseek64()	<unistd.h>
ftruncate Subroutine	ftruncate64()	<unistd.h>
truncate Subroutine	truncate64()	<unistd.h>
fclear Subroutine	fclear64()	<unistd.h>
pwrite Subroutine	pwrite64()	<unistd.h>
pread Subroutine	pread64()	<unistd.h>
struct aiocb Structure	struct aiocb64	<sys/aio.h>
aio_read Subroutine	aio_read64()	<sys/aio.h>
aio_write Subroutine	aio_write64()	<sys/aio.h>
aio_cancel Subroutine	aio_cancel64()	<sys/aio.h>
aio_suspend Subroutine	aio_suspend64()	<sys/aio.h>
aio_return Subroutine	aio_return64()	<sys/aio.h>
aio_error Subroutine	aio_error64()	<sys/aio.h>
liocb Structure	liocb64	<sys/aio.h>
lio_listio Subroutine	lio_listio64()	<sys/aio.h>

Using 64-bit file system subroutines

Using the `_LARGE_FILES` environment may be impractical for some applications due to the far-reaching implications of changing the size of `off_t` to 64 bits. If the number of changes is small, it may be more practical to convert a relatively small part of the application to be large-file enabled. The 64-bit file system data types, structures, and subroutines are listed below:

```
<sys/types.h>
typedef long long off64_t;
typedef long long fpos64_t;

<fcntl.h>

extern int      open64(const char *, int, ...);
extern int      creat64(const char *, mode_t);

#define F_GETLK64
#define F_SETLK64
#define F_SETLKW64

<ftw.h>
extern int ftw64(const char *, int (*)(const char *, const struct stat64 *, int), int);
extern int nftw64(const char *, int (*)(const char *, const struct stat64 *, int, struct FTW *), int, int);

<stdio.h>

extern int      fgetpos64(FILE *, fpos64_t *);
extern FILE    *fopen64(const char *, const char *);
extern FILE    *freopen64(const char *, const char *, FILE *);
extern int      fseek64(FILE *, off64_t, int);
extern int      fsetpos64(FILE *, fpos64_t *);
extern off64_t  ftello64(FILE *);

<unistd.h>

extern off64_t  lseek64(int, off64_t, int);
extern int      ftruncate64(int, off64_t);
extern int      truncate64(const char *, off64_t);
extern off64_t  fclear64(int, off64_t);
extern ssize_t  pread64(int, void *, size_t, off64_t);
extern ssize_t  pwrite64(int, const void *, size_t, off64_t);
extern int      fsync_range64(int, int, off64_t, off64_t);

<sys/flock.h>

struct flock64;

<sys/lockf.h>
```

```

extern int lockf64 (int, int, off64_t);

<sys/mman.h>

extern void      *mmap64(void *, size_t, int, int, int, off64_t);

<sys/stat.h>

struct stat64;

extern int      stat64(const char *, struct stat64 *);
extern int      fstat64(int, struct stat64 *);
extern int      lstat64(const char *, struct stat64 *);

<sys/aio.h>

struct aiocb64
int      aio_read64(int, struct aiocb64 *):
int      aio_write64(int, struct aiocb64 *);
int      aio_listio64(int, struct aiocb64 *[],
int, struct sigevent *);
int      aio_cancel64(int, struct aiocb64 *);
int      aio_suspend64(int, struct aiocb64 *[]);

struct liocb64
int      lio_listio64(int, struct liocb64 *[], int, void *);

```

Common pitfalls in using the large file environment

Porting of application programs to the large-file environment can expose a number of different problems in the application. These problems are frequently the result of poor coding practices, which are harmless in a 32-bit **off_t** environment, but which can manifest themselves when compiled in a 64-bit **off_t** environment. Some of the more common problems and solutions are discussed in this section.

Note: In the following examples, **off_t** is assumed to be a 64-bit file offset.

Improper use of data types

A common source of problems with application programs is a failure to use the proper data types. If an application attempts to store file sizes or file offsets in an integer variable, the resulting value will be truncated and lose significance. To avoid this problem, use the **off_t** data type to store file sizes and offsets.

Incorrect:

```

int file_size;
struct stat s;

file_size = s.st_size;

```

Better:

```

off_t file_size;
struct stat s;
file_size = s.st_size;

```

When you are passing 64-bit integers to functions as arguments or when you are returning 64-bit integers from functions, both the caller and the called function *must* agree on the types of the arguments and the return value.

Passing a 32-bit integer to a function that expects a 64-bit integer causes the called function to misinterpret the caller's arguments, leading to unexpected behavior. This type of problem is especially severe if the program passes scalar values to a function that expects to receive a 64-bit integer.

You can avoid problems by using function prototypes carefully. In the code fragments below, **example()** is a function that takes a 64-bit file offset as a parameter. In the first example, the compiler generates the normal 32-bit integer function linkage, which would be incorrect because the receiving function expects 64-bit integer linkage. In the second example, the LL specifier is added, forcing the compiler to use the proper linkage. In the last example, the function prototype causes the compiler to promote the scalar value to a 64-bit integer. This is the preferred approach because the source code remains portable between 32-bit and 64-bit environments.

Incorrect:

```
fexample(0);
```

Better:

```
fexample(0LL);
```

Best:

```
\est:
```

Linking for programmers

A link is a connection between a file name and an i-node (hard link) or between file names (symbolic link).

Linking allows access to an i-node from multiple file names. Directory entries pair file names with i-nodes. File names are easy for users to identify, and i-nodes contain the real disk addresses of the file's data. A reference count of all links into an i-node is maintained in the **i_nlink** field of the i-node. Subroutines that create and destroy links use file names, not file descriptors. Therefore, it is not necessary to open files when creating a link.

Processes can access and change the contents of the i-node by any of the linked file names. AIX supports hard links and symbolic links.

Hard links

Subroutine**Description****link**

Subroutine that creates hard links. The presence of a hard link guarantees the existence of a file because a hard link increments the link count in the **i_nlink** field of the i-node.

unlink

Subroutine that releases links. When all hard links to an i-node are released, the file is no longer accessible.

Hard links must link file names and i-nodes within the same file system because the i-node number is relative to a single file system. Hard links always refer to a specific file because the directory entry created by the hard link pairs the new file name to an i-node. The user ID that created the original file owns the file and retains access mode authority over the file. Otherwise, all hard links are treated equally by the operating system.

Example: If the **/u/tom/bob** file is linked to the **/u/jack/foo** file, the link count in the **i_nlink** field of the **foo** file is 2. Both hard links are equal. If **/u/jack/foo** is removed, the file continues to exist by the name **/u/tom/bob** and can be accessed by users with access to the **tom** directory. However, the owner of the file is **jack** even though **/u/jack/foo** was removed. The space occupied by the file is charged to **jack's** quota account. To change file ownership, use the **chown** subroutine.

Symbolic links

Symbolic links are implemented as a file that contains a path name by using the **symlink** command. When a process encounters a symbolic link, the path contained in the symbolic link is prepended to the path that the process was searching. If the path name in the symbolic link is an absolute path name, the process searches from the root directory for the named file. If the path name in the symbolic link does not begin with a / (slash), the process interprets the rest of the path relative to the position of the symbolic link. The **unlink** subroutine also removes symbolic links.

Symbolic links can traverse file systems because they are treated as regular files by the operating system rather than as part of the file system structure. The presence of a symbolic link does not guarantee the existence of the target file because a symbolic link has no effect on the **i_nlink** field of the i-node.

Subroutine Description

readlink Subroutine that reads the contents of a symbolic link. Many subroutines (including the **open** and **stat** subroutines) follow symbolic paths.

lstat Subroutine created to report on the status of the file containing the symbolic link and does not follow the link. See the **symlink** subroutine for a list of subroutines that traverse symbolic links.

Symbolic links are also called *soft links* because they link to a file by path name. If the target file is renamed or removed, the symbolic link cannot resolve.

Example: The symbolic link to **/u/joe/foo** is a file that contains the literal data **/u/joe/foo**. When the owner of the **foo** file removes this file, subroutine calls made to the symbolic link cannot succeed. If the file owner then creates a new file named **foo** in the same directory, the symbolic link leads to the new file. Therefore, the link is considered soft because it is linked to interchangeable i-nodes.

In the **ls -l** command listing, an **l** in the first position indicates a linked file. In the final column of that listing, the links between files are represented as Path2 -> Path1 (or Newname -> Oldname).

Subroutine

Description

unlink

Subroutine that removes a directory entry. The *Path* parameter in the subroutine identifies the file to be disconnected. At the completion of the **unlink** call, the link count of the i-node is reduced by the value of 1.

remove

Subroutine that also removes a file name by calling either the **unlink** or **rmdir** subroutine.

Directory links

Subroutine

Description

mkdir

Subroutine that creates directory entries for new directories, which creates hard links to the i-node representing the directory

Symbolic links are recommended for creating additional links to a directory. Symbolic links do not interfere with the **.** and **..** directory entries and will maintain the empty, well-formed directory status. The following illustrates an example of the empty, well-formed directory **/u/joe/foo** and the **i_nlink** values.

`/u`

Values	empty values	empty values	Directory			
68			j	o	e	0

`/u/joe
mkdir ("foo", 0666)`

Values	empty values	empty values	Directory			
68			n	0	0	0
			n	n	0	0
235			f	o	o	0

/u/joe/fo

Values	empty values	empty values	Directory			
235			n	0	0	0
68			n	n	0	0

i_nlink Values

i = 68

n_link 3

For i = 68, the n_link value is 3 (/u; /u/joe; /u/joe/fo).

i = 235

n_link 2

For i = 235, the n_link value is 2 (/u/joe; /u/joe/fo).

Using file descriptors

A file descriptor is an unsigned integer used by a process to identify an open file.

The number of file descriptors available to a process is limited by the **/OPEN_MAX** control in the **sys/limits.h** file. The number of file descriptors is also controlled by the **ulimit -n** flag. The **open**, **pipe**, **creat**, and **fcntl** subroutines all generate file descriptors. File descriptors are generally unique to each process, but they can be shared by child processes created with a **fork** subroutine or copied by the **fcntl**, **dup**, and **dup2** subroutines.

File descriptors are indexes to the file descriptor table in the **u_block** area maintained by the kernel for each process. The most common ways for processes to obtain file descriptors are through **open** or **creat** operations or through inheritance from a parent process. When a **fork** operation occurs, the descriptor table is copied for the child process, which allows the child process equal access to the files used by the parent process.

File descriptor tables and system open file tables

The file descriptor and open file table structures track each process' access to a file and ensure data integrity.

Table

file descriptor table

Description

Translates an index number (file descriptor) in the table to an open file. File descriptor tables are created for each process and are located in the **u_block** area set aside for that process. Each of the entries in a file descriptor table has the following fields: the flags area and the file pointer. There are up to **OPEN_MAX** file descriptors. The structure of the file descriptor table is as follows:

```
struct ufd
{
    struct file *fp;
    int flags;
} *u_ufd
```

Table	Description
system open file table	<p>Contains entries for each open file. A file table entry tracks the current offset referenced by all read or write operations to the file and the open mode (O_RDONLY, O_WRONLY, or O_RDWR) of the file.</p> <p>The open file table structure contains the current I/O offset for the file. The system treats each read/write operation as an implied seek to the current offset. Thus if <i>x</i> bytes are read or written, the pointer advances <i>x</i> bytes. The lseek subroutine can be used to reassign the current offset to a specified location in files that are randomly accessible. Stream-type files (such as pipes and sockets) do not use the offset because the data in the file is not randomly accessible.</p>

Managing file descriptors

Because files can be shared by many users, it is necessary to allow related processes to share a common offset pointer and have a separate current offset pointer for independent processes that access the same file. The open file table entry maintains a reference count to track the number of file descriptors assigned to the file.

Multiple references to a single file can be caused by any of the following:

- A separate process opening the file
- Child processes retaining the file descriptors assigned to the parent process
- The **fcntl** or **dup** subroutine creating copies of the file descriptors

Sharing open files

Each open operation creates a system open file table entry. Separate table entries ensure each process has separate current I/O offsets. Independent offsets protect the integrity of the data.

When a file descriptor is duplicated, two processes then share the same offset and interleaving can occur, in which bytes are not read or written sequentially.

Duplicating file descriptors

File descriptors can be duplicated between processes in the following ways: the **dup** or **dup2** subroutine, the **fork** subroutine, and the **fcntl** (file descriptor control) subroutine.

dup and dup2 subroutines

The **dup** subroutine creates a copy of a file descriptor. The duplicate is created at an empty space in the user file descriptor table that contains the original descriptor. A **dup** process increments the reference count in the file table entry by 1 and returns the index number of the file-descriptor where the copy was placed.

The **dup2** subroutine scans for the requested descriptor assignment and closes the requested file descriptor if it is open. It allows the process to designate which descriptor entry the copy will occupy, if a specific descriptor-table entry is required.

fork subroutine

The **fork** subroutine creates a child process that inherits the file descriptors assigned to the parent process. The child process then execs a new process. Inherited descriptors that had the **close-on-exec** flag set by the **fcntl** subroutine close.

fcntl (file descriptor control) subroutine

The **fcntl** subroutine manipulates file structure and controls open file descriptors. It can be used to make the following changes to a descriptor:

- Duplicate a file descriptor (identical to the **dup** subroutine).
- Get or set the close-on-exec flag.

- Set nonblocking mode for the descriptor.
- Append future writes to the end of the file (**O_APPEND**).
- Enable the generation of a signal to the process when it is possible to do I/O.
- Set or get the process ID or the group process ID for **SIGIO** handling.
- Close all file descriptors.

Preset file descriptor values

When the shell runs a program, it opens three files with file descriptors 0, 1, and 2. The default assignments for these descriptors are as follows:

Descriptor	Explanation
0	Represents standard input.
1	Represents standard output.
2	Represents standard error.

These default file descriptors are connected to the terminal, so that if a program reads file descriptor 0 and writes file descriptors 1 and 2, the program collects input from the terminal and sends output to the terminal. As the program uses other files, file descriptors are assigned in ascending order.

If I/O is redirected using the < (less than) or > (greater than) symbols, the shell's default file descriptor assignments are changed. For example, the following changes the default assignments for file descriptors 0 and 1 from the terminal to the appropriate files:

```
prog < FileX > FileY
```

In this example, file descriptor 0 now refers to FileX and file descriptor 1 refers to FileY. File descriptor 2 has not been changed. The program does not need to know where its input comes from nor where it is sent, as long as file descriptor 0 represents the input file and 1 and 2 represent output files.

The following sample program illustrates the redirection of standard output:

```
#include <fcntl.h>
#include <stdio.h>

void redirect_stdout(char *);

main()
{
    printf("Hello world\n");           /*this printf goes to
                                     * standard output*/
    fflush(stdout);
    redirect_stdout("foo");           /*redirect standard output*/
    printf("Hello to you too, foo\n"); /*printf goes to file foo */
    fflush(stdout);
}

void
redirect_stdout(char *filename)
{
    int fd;
    if ((fd = open(filename,O_CREAT|O_WRONLY,0666)) < 0)
        /*open a new file */
    {
        perror(filename);
        exit(1);
    }
    close(1);                          /*close old */
    /*standard output*/
    if (dup(fd) !=1)                    /*dup new fd to
    /*standard input*/
    {
        fprintf(stderr,"Unexpected dup failure\n");
        exit(1);
    }
    close(fd);                          /*close original, new fd,*/
}
```

```
* no longer needed*/
```

```
}
```

Within the file descriptor table, file descriptor numbers are assigned the lowest descriptor number available at the time of a request for a descriptor. However, any value can be assigned within the file descriptor table by using the **dup** subroutine.

File descriptor resource limit

The number of file descriptors that can be allocated to a process is governed by a resource limit. The default value is set in the **/etc/security/limits** file and is typically set at **2000**. The limit can be changed by the **ulimit** command or the **setrlimit** subroutine. The maximum size is defined by the constant **OPEN_MAX**.

Creating and removing files

This section describes the internal procedures performed by the operating system when creating, opening, closing, or removing files.

Creating a file

Different subroutines create specific types of files, as follows:

Subroutine	Type of File Created
creat	Regular
open	Regular (when the O_CREAT flag is set)
mknod	Regular, first-in-first-out (FIFO), or special
mkfifo	Named pipe (FIFO)
pipe	Unnamed pipe
socket	Sockets
mkdir	Directories
symlink	Symbolic link

Creating a regular file (creat, open, or mknod subroutines)

You use the **creat** subroutine to create a file according to the values set in the *Pathname* and *Mode* parameters. If the file named in the *Pathname* parameter exists and the process has write permission to the file, the **creat** subroutine truncates the file. Truncation releases all data blocks and sets the file size to 0. You can also create new, regular files using the **open** subroutine with the **O_CREAT** flag.

Files created with the **creat**, **mkfifo**, or **mknod** subroutine take the access permissions set in the *Mode* parameter. Regular files created with the **open** subroutine take their access modes from the **O_CREAT** flag *Mode* parameter. The **umask** subroutine sets a file-mode creation mask (set of access modes) for new files created by processes and returns the previous value of the mask.

The permission bits on a newly created file are a result of the reverse of the **umask** bits ANDed with the file-creation mode bits set by the creating process. When a new file is created by a process, the operating system performs the following actions:

- Determines the permissions of the creating process
- Retrieves the appropriate **umask** value
- Reverses the **umask** value
- Uses the AND operation to combine the permissions of the creating process with the reverse of the **umask** value

Creating a special file (mknod or mkfifo subroutine)

You can use the **mknod** and **mkfifo** subroutines to create new special files. The **mknod** subroutine handles named pipes (FIFO), ordinary, and device files. It creates an i-node for a file identical to that created by the **creat** subroutine. When you use the **mknod** subroutine, the file-type field is set to indicate the type of file being created. If the file is a block or character-type device file, the names of the major and minor devices are written into the i-node.

The **mkfifo** subroutine is an interface for the **mknod** subroutine and is used to create named pipes.

Opening a file

The **open** subroutine is the first step required for a process to access an existing file. The **open** subroutine returns a file descriptor. Reading, writing, seeking, duplicating, setting I/O parameters, determining file status, and closing the file all use the file descriptor returned by the **open** call. The **open** subroutine creates entries for a file in the file descriptor table when assigning file descriptors.

The **open** subroutine does the following:

- Checks for appropriate permissions that allow the process access to the file.
- Assigns an entry in the file descriptor table for the open file. The **open** subroutine sets the initial read/write byte offset to 0, the beginning of the file.

The **ioctl** or **ioctlx** subroutines perform control operations on opened special device files.

Closing a file

When a process no longer needs access to the open file, the **close** subroutine removes the entry for the file from the table. If more than one file descriptor references the file table entry for the file, the reference count for the file is decreased by 1, and the close completes. If a file has only 1 reference to it, the file table entry is freed. Attempts by the process to use the disconnected file descriptor result in errors until another **open** subroutine reassigns a value for that file descriptor value. When a process exits, the kernel examines its active user file descriptors and internally closes each one. This action ensures that all files close before the process ends.

Removing a file

When a file is no longer needed, you can use the **unlink** subroutine to remove the specified file from the directory containing it. If there are multiple hard links to the same file, the **unlink** subroutine removes the specified link. If there is only one link, the **unlink** subroutine removes the file itself. For more information, see the **unlink** subroutine.

Working with file I/O

All input and output (I/O) operations use the current file offset information stored in the system file structure.

The current I/O offset designates a byte offset that is constantly tracked for every open file. The current I/O offset signals a read or write process where to begin operations in the file. The **open** subroutine resets it to 0. The pointer can be set or changed using the **lseek** subroutine.

Manipulating the current offset

Read and write operations can access a file sequentially because the current I/O offset of the file tracks the byte offset of each previous operation. The offset is stored in the system file table.

You can adjust the offset on files that can be randomly accessed, such as regular and special-type files, using the **lseek** subroutine.

Subroutine

Description

lseek

Allows a process to position the offset at a designated byte. The **lseek** subroutine positions the pointer at the byte designated by the *Offset* variable. The *Offset* value can be calculated from the following places in the file (designated by the value of the *Whence* variable):

absolute offset

Beginning byte of the file

relative offset

Position of the former pointer

end_relative offset

End of the file

The return value for the **lseek** subroutine is the current value of the pointer's position in the file. For example:

```
cur_off= lseek(fd, 0, SEEK_CUR);
```

The **lseek** subroutine is implemented in the file table. All subsequent read and write operations use the new position of the offset as their starting location.

Note: The offset cannot be changed on pipes or socket-type files.

Subroutine

Description

fclear

Subroutine that creates an empty space in a file. It sets to zero the number of bytes designated in the *NumberOfBytes* variable beginning at the current offset. The **fclear** subroutine cannot be used if the **O_DEFER** flag was set at the time the file was opened.

Reading a file

This section describes the read subroutines.

Reading a file

Subroutine Description

read

Subroutine that copies a specified number of bytes from an open file to a specified buffer. The copy begins at the point indicated by the current offset. The number of bytes and buffer are specified by the *NBytes* and *Buffer* parameters.

The **read** subroutine does the following:

1. Ensures that the *FileDescriptor* parameter is valid and that the process has **read** permissions. The subroutine then gets the file table entry specified by the *FileDescriptor* parameter.
2. Sets a flag in the file to indicate a read operation is in progress. This action locks other processes out of the file during the operation.
3. Converts the offset byte value and the value of the *NBytes* variables into a block address.
4. Transfers the contents of the identified block into a storage buffer.
5. Copies the contents of the storage buffer into the area designated by the *Buffer* variable.
6. Updates the current offset according to the number of bytes actually read. Resetting the offset ensures that the data is read in sequence by the next read process.
7. Deducts the number of bytes read from the total specified in the *NByte* variable.
8. Loops until the number of bytes to be read is satisfied.
9. Returns the total number of bytes read.

The cycle completes when the file to be read is empty, the number of bytes requested is met, or a reading error is encountered during the process.

To avoid an extra iteration in the read loop, start read requests at the beginning of data block boundaries and to be multiples of the data block size. If a process reads blocks sequentially, the operating system assumes all subsequent reads will also be sequential.

During the read operation, the i-node is locked. No other processes are allowed to modify the contents of the file while a read is in progress. However the file is unlocked immediately on completion of the read operation. If another process changes the file between two read operations, the resulting data is different, but the integrity of the data structure is maintained.

The following example illustrates how to use the read subroutine to count the number of null bytes in the **foo** file:

```
#include <fcntl.h>
#include <sys/param.h>

main()
{
    int fd;
    int nbytes;
    int nnulls;
    int i;
    char buf[PAGESIZE];      /*A convenient buffer size*/
    nnulls=0;
    if ((fd = open("foo",O_RDONLY)) < 0)
        exit();
    while ((nbytes = read(fd,buf,sizeof(buf))) > 0)
        for (i = 0; i < nbytes; i++)
            if (buf[i] == '\0')
                nnulls++;
    printf("%d nulls found\n", nnulls);
}
```

Writing a file

This section describes the write subroutines.

Subroutine

Description

write

Subroutine that adds the amount of data specified in the *NBytes* variable from the space designated by the *Buffer* variable to the file described by the *FileDescriptor* variable. It functions similar to the **read** subroutine. The byte offset for the write operation is found in the system file table's current offset.

If you write to a file that does not contain a block corresponding to the byte offset resulting from the write process, the **write** subroutine allocates a new block. This new block is added to the i-node information that defines the file. Adding the new block might allocate more than one block if the underlying file system needs to add blocks for addressing the file blocks.

During the write operation, the i-node is locked. No other processes are allowed to modify the contents of the file while a write is in progress. However, the file is unlocked immediately on completion of the write operation. If another process changes the file between two write operations, the resulting data is different, but the integrity of the data structure is maintained.

The **write** subroutine loops in a way similar to the **read** subroutine, logically writing one block to disk for each iteration. At each iteration, the process either writes an entire block or only a portion of one. If only a portion of a data block is required to accomplish an operation, the **write** subroutine reads the block from disk to avoid overwriting existing information. If an entire block is required, it does not read the block because the entire block is overwritten. The write operation proceeds block by block until the number of bytes designated in the *NBytes* parameter is written.

Delayed write

You can designate a delayed write process with the **O_DEFER** flag. The data is then transferred to disk as a temporary file. The delayed write feature caches the data in case another process reads or writes the data sooner. Delayed write saves extra disk operations. Many programs, such as mail and editors, create temporary files in the **/tmp** directory and quickly remove them.

When a file is opened with the deferred update (**O_DEFER**) flag, the data is not written to permanent storage until a process issues an **fsync** subroutine call or a process issues a synchronous **write** to the file (opened with **O_SYNC** flag). The **fsync** subroutine saves all changes in an open file to disk. See the **open** subroutine for a description of the **O_DEFER** and **O_SYNC** flags.

Truncating files

The **truncate** or **ftruncate** subroutines change the length of regular files. The truncating process must have write permission to the file. The *Length* variable value indicates the size of the file after the truncation operation is complete. All measures are relative to the first byte of the file, not the current offset. If the new length (designated in the *Length* variable) is less than the previous length, the data between the two is removed. If the new length is greater than the existing length, zeros are added to extend the file size. When truncation is complete, full blocks are returned to the file system, and the file size is updated.

Direct I/O vs. normal cached I/O

Normally, the JFS or JFS2 caches file pages in kernel memory. When the application does a file read request, if the file page is not in memory, the JFS or JFS2 reads the data from the disk into the file cache, then copies the data from the file cache to the user's buffer..

For application writes, the data is merely copied from the user's buffer into the cache. The actual writes to disk are done later.

This type of caching policy can be extremely effective when the cache hit rate is high. It also enables read-ahead and write-behind policies. Lastly, it makes file writes asynchronous, allowing the application to continue processing instead of waiting for I/O requests to complete.

Direct I/O is an alternative caching policy that causes the file data to be transferred directly between the disk and the user's buffer. Direct I/O for files is functionally equivalent to raw I/O for devices. Applications can use direct I/O on JFS or JFS2 files.

Benefits of direct I/O

The primary benefit of direct I/O is to reduce CPU utilization for file reads and writes by eliminating the copy from the cache to the user buffer.

This can also be a benefit for file data which has a very poor cache hit rate. If the cache hit rate is low, then most read requests have to go to the disk. Direct I/O can also benefit applications that must use synchronous writes because these writes have to go to disk. In both of these cases, CPU usage is reduced because the data copy is eliminated.

A second benefit of direct I/O is that it allows applications to avoid diluting the effectiveness of caching of other files. Anytime a file is read or written, that file competes for space in the cache. This situation may cause other file data to be pushed out of the cache. If the newly cached data has very poor reuse characteristics, the effectiveness of the cache can be reduced. Direct I/O gives applications the ability to identify files where the normal caching policies are ineffective, thus releasing more cache space for files where the policies are effective.

Performance costs of direct I/O

Although direct I/O can reduce CPU usage, using it typically results in the process taking longer to complete, especially for relatively small requests. This penalty is caused by the fundamental differences between normal cached I/O and direct I/O.

Direct I/O reads

Every direct I/O read causes a synchronous read from disk; unlike the normal cached I/O policy where read may be satisfied from the cache. This can result in very poor performance if the data was likely to be in memory under the normal caching policy.

Direct I/O also bypasses the normal JFS or JFS2 read-ahead algorithms. These algorithms can be extremely effective for sequential access to files by issuing larger and larger read requests and by overlapping reads of future blocks with application processing.

Applications can compensate for the loss of JFS or JFS2 read-ahead by issuing larger read requests. At a minimum, direct I/O readers should issue read requests of at least 128k to match the JFS or JFS2 read-ahead characteristics.

Applications can also simulate JFS or JFS2 read-ahead by issuing asynchronous direct I/O read-ahead either by use of multiple threads or by using the **aio_read** subroutine.

Direct I/O writes

Every direct I/O write causes a synchronous write to disk; unlike the normal cached I/O policy where the data is merely copied and then written to disk later. This fundamental difference can cause a significant performance penalty for applications that are converted to use direct I/O.

Conflicting file access modes

To avoid consistency issues between programs that use direct I/O and programs that use normal cached I/O, direct I/O is an exclusive use mode. If there are multiple opens of a file and some of them are direct and others are not, the file will stay in its normal cached access mode. Only when the file is open exclusively by direct I/O programs will the file be placed in direct I/O mode.

Similarly, if the file is mapped into virtual memory through the **shmat** or **mmap** system calls, the file will stay in normal cached mode.

The JFS or JFS2 will attempt to move the file into direct I/O mode anytime the last conflicting or non-direct access is eliminated (either by the **close**, **munmap**, or **shmdt** subroutines). Changing the file from normal mode to direct I/O mode can be rather expensive because it requires writing all modified pages to disk and removing all the file's pages from memory.

Enabling applications to use direct I/O

Applications enable direct I/O access to a file by passing the **O_DIRECT** flag to the **open** subroutine. This flag is defined in the **fcntl.h** file. Applications must be compiled with **_ALL_SOURCE** enabled to see the definition of **O_DIRECT**.

Offset/Length/Address alignment requirements of the target buffer

For direct I/O to work efficiently, the request should be suitably conditioned. Applications can query the offset, length, and address alignment requirements by using the **finfo** and **ffinfo** subroutines. When the **FI_DIOCAP** command is used, the **finfo** and **ffinfo** subroutines return information in the **diocapbuf** structure as described in the **sys/finfo.h** file. This structure contains the following fields:

dio_offset

Recommended offset alignment for direct I/O writes to files in this file system

dio_max

Recommended maximum write length for direct I/O writes to files in this system

dio_min

Recommended minimum write length for direct I/O writes to files in this file system

dio_align

Recommended buffer alignment for direct I/O writes to files in this file system

Failure to meet these requirements may cause file reads and writes to use the normal cached model and may cause direct I/O to be disabled for the file. Different file systems may have different requirements, as the following table illustrates.

Table 1. File system format

File system format	dio_offset	dio_max	dio_min	dio_align
JFS fixed, 4 K blk	4 K	2 MB	4 K	4 K
JFS fragmented	4 K	2 MB	4 K	4 K
JFS compressed	n/a	n/a	n/a	n/a
JFS big file	128 K	2 MB	128 K	4 K
JFS2	4 K	4 GB	4 K	4 K

Direct I/O limitations

Direct I/O is not supported for files in a compressed-file file system. Attempts to open these files with `O_DIRECT` will be ignored and the files will be accessed with the normal cached I/O methods.

Direct I/O and data I/O integrity completion

Although direct I/O writes are done synchronously, they do not provide synchronized I/O data integrity completion, as defined by POSIX. Applications that need this feature should use `O_DSYNC` in addition to `O_DIRECT`. `O_DSYNC` guarantees that all of the data and enough of the metadata (for example, indirect blocks) have written to the stable store to be able to retrieve the data after a system crash. `O_DIRECT` only writes the data; it does not write the metadata.

Working with pipes

Pipes are unnamed objects created to allow two processes to communicate.

One process reads and the other process writes to the pipe file. This unique type of file is also called a first-in-first-out (FIFO) file. The data blocks of the FIFO are manipulated in a circular queue, maintaining read and write pointers internally to preserve the FIFO order of data. The `PIPE_BUF` system variable, defined in the `limits.h` file, designates the maximum number of bytes guaranteed to be atomic when written to a pipe.

The shell uses unnamed pipes to implement command pipelining. Most unnamed pipes are created by the shell. The `|` (vertical) symbol represents a pipe between processes. In the following example, the output of the `ls` command is printed to the screen:

```
ls | pr
```

Pipes are treated as regular files as much is possible. Normally, the current offset information is stored in the system file table. However, because pipes are shared by processes, the read/write pointers must be specific to the file, not to the process. File table entries are created by the `open` subroutine and are unique to the open process, not to the file. Processes with access to pipes share the access through common system file table entries.

Using pipe subroutines

The `pipe` subroutine creates an interprocess channel and returns two file descriptors. File descriptor 0 is opened for reading. File descriptor 1 is opened for writing. The read operation accesses the data on a FIFO basis. These file descriptors are used with `read`, `write`, and `close` subroutines.

In the following example, a child process is created and sends its process ID back through a pipe:

```
#include <sys/types.h>
main()
{
    int p[2];
    char buf[80];
    pid_t pid;

    if (pipe(p))
    {
        perror("pipe failed");
        exit(1);
    }
}
```

```

}
if ((pid=fork()) == 0)
{
    close(p[0]);          /* in child process */
                        /*close unused read */
                        /*side of the pipe */
    sprintf(buf, "%d", getpid());
                        /*construct data */
                        /*to send */
    write(p[1], buf, strlen(buf)+1);
                        /*write it out, including
                        /*null byte */
    exit(0);
}
close(p[1]);          /*in parent process*/
read(p[0], buf, sizeof(buf)); /*close unused write side of pipe */
printf("Child process said: %s/n", buf); /*read the pipe*/
exit(0);              /*display the result */
}

```

If a process reads an empty pipe, the process waits until data arrives. If a process writes to a pipe that is too full (**PIPE_BUF**), the process waits until space is available. If the write side of the pipe is closed, a subsequent read operation to the pipe returns an end-of-file.

Other subroutines that control pipes are the **popen** and **pclose** subroutines:

popen

Creates the pipe (using the **pipe** subroutine) then forks to create a copy of the caller. The child process decides whether it is supposed to read or write, closes the other side of the pipe, then calls the shell (using the **execl** subroutine) to run the desired process.

The parent closes the end of the pipe it did not use. These closes are necessary to make end-of-file tests work correctly. For example, if a child process intended to read the pipe does not close the write end of the pipe, it will never see the end of file condition on the pipe, because there is one write process potentially active.

The conventional way to associate the pipe descriptor with the standard input of a process is:

```

close(p[1]);
close(0);
dup(p[0]);
close(p[0]);

```

The **close** subroutine disconnects file descriptor 0, the standard input. The **dup** subroutine returns a duplicate of an already open file descriptor. File descriptors are assigned in ascending order and the first available one is returned. The effect of the **dup** subroutine is to copy the file descriptor for the pipe (read side) to file descriptor 0, thus standard input becomes the read side of the pipe. Finally, the previous read side is closed. The process is similar for a child process to write from a parent.

pclose

Closes a pipe between the calling program and a shell command to be executed. Use the **pclose** subroutine to close any stream opened with the **popen** subroutine.

The **pclose** subroutine waits for the associated process to end, then closes and returns the exit status of the command. This subroutine is preferable to the **close** subroutine because **pclose** waits for child processes to finish before closing the pipe. Equally important, when a process creates several children, only a bounded number of unfinished child processes can exist, even if some of them have completed their tasks. Performing the wait allows child processes to complete their tasks.

Synchronous I/O

By default, writes to files in JFS or JFS2 file systems are asynchronous.

However, JFS and JFS2 file systems support the following types of synchronous I/O:

- Specified by the `O_DSYNC` open flag. When a file is opened using the `O_DSYNC` open mode, the `write ()` system call will not return until the file data and all file system meta-data required to retrieve the file data are both written to their permanent storage locations.
- Specified by the `O_SYNC` open flag. In addition to items specified by `O_DSYNC`, `O_SYNC` specifies that the `write ()` system call will not return until all file attributes relative to the I/O are written to their permanent storage locations, even if the attributes are not required to retrieve the file data.

Before the `O_DSYNC` open mode existed, AIX applied `O_DSYNC` semantics to `O_SYNC`. For binary compatibility reasons, this behavior still exists. If true `O_SYNC` behavior is required, then both `O_DSYNC` and `O_SYNC` open flags must be specified. Exporting the `XPG_SUS_ENV=ON` environment variable also enables true `O_SYNC` behavior.

- Specified by the `O_RSYNC` open flag, and it simply applies the behaviors associated with `O_SYNC` or `_DSYNC` to reads. For files in JFS and JFS2 file systems, only the combination of `O_RSYNC | O_SYNC` has meaning, indicating that the read system call will not return until the file's access time is written to its permanent storage location.

File status

File status information resides in the i-node.

The **stat** subroutines are used to return information on a file. The **stat** subroutines report file type, file owner, access mode, file size, number of links, i-node number, and file access times. These subroutines write information into a data structure designated by the *Buffer* variable. The process must have search permission for the directories in the path to the designated file.

Subroutine	Description
stat	Subroutine that returns the information about files named by the <i>Path</i> parameter. If the size of the file cannot be represented in the structure designated by the <i>Buffer</i> variable, stat will fail with the <code>errno</code> set to <code>E_OVERFLOW</code> .
lstat	Subroutine that provides information about a symbolic link, and the stat subroutine returns information about the file referenced by the link.
fstat	Returns information from an open file using the file descriptor.

The **statfs**, **fstatfs**, and **ustat** subroutines return status information about a file system.

Subroutine	Description
fstatfs	Returns the information about the file system that contains the file associated with the given file descriptor. The structure of the returned information is described in the <code>/usr/include/sys/statfs.h</code> file for the statfs and fstatfs subroutines and in the <code>ustat.h</code> file for the ustat subroutine.
statfs	Returns information about the file system that contains the file specified by the <i>Path</i> parameter.
ustat	Returns information about a mounted file system designated by the <i>Device</i> variable. This device identifier is for any given file and can be determined by examining the <code>st_dev</code> field of the stat structure defined in the <code>/usr/include/sys/stat.h</code> file. The ustat subroutine is superseded by the statfs and fstatfs subroutines.
utimes and utime	Also affect file status information by changing the file access and modification time in the i-node.

File accessibility

Every file is created with an access mode. Each access mode grants read, write, or execute permission to users, the user's group, and all other users.

The access bits on a newly created file are a result of the reverse of the **umask** bits ANDed with the file-creation mode bits set by the creating process. When a new file is created by a process, the operating system performs the following actions:

- Determines the permissions of the creating process
- Retrieves the appropriate **umask** value
- Reverses the **umask** value
- Uses the AND operation to combine the permissions of the creating process with the reverse of the **umask** value

For example, if an existing file has the 027 permissions bits set, the user is not allowed any permissions. Write permission is granted to the group. Read, write, and execute access is set for all others. The **umask** value of the 027 permissions modes would be 750 (the opposite of the original permissions). When 750 is ANDed with 666 (the file creation mode bits set by the system call that created the file), the actual permissions for the file are 640. Another representation of this example is:

027 =	\bar{r}	\bar{w}	\bar{x}	\bar{r}	\bar{w}	\bar{x}	Existing file access mode
750 =	\bar{r}	\bar{w}	\bar{x}	\bar{r}	\bar{w}	\bar{x}	Reverse (umask) of original permissions
666 =	R	W		R	W	X	File creation access mode
ANDED TO							
750 =	R	W	X	R		X	The umask value
640 =	R	W		R			Resulting file access mode

Subroutine	Description
access subroutine	Investigates and reports on the accessibility mode of the file named in the <i>Pathname</i> parameter. This subroutine uses the real user ID and the real group ID instead of the effective user and group ID. Using the real user and group IDs allows programs with the set-user-ID and set-group-ID access modes to limit access only to users with proper authorization.
chmod and fchmod subroutines	Changes file access permissions.
chown subroutine	Resets the ownership field of the i-node for the file and clears the previous owner. The new information is written to the i-node and the process finishes. The chmod subroutine works in similar fashion, but the permission mode flags are changed instead of the file ownership.
umask	Sets and gets the value of the file creation mask.

In the following example, the user does not have access to the file `secrets`. However, when the program `special` is run and the access mode for the program is set-uID root, the program can access the file. The program must use the **access** subroutine to prevent subversion of system security.

```
$ ls -l
total 0
-r-s--x--x   1 root  system   8290 Jun 09 17:07 special
-rw-----   1 root  system   1833 Jun 09 17:07 secrets
$ cat secrets
cat: cannot open secrets
```

The **access** subroutine must be used by any set-uID or set-gID program to forestall this type of intrusion. Changing file ownership and access modes are actions that affect the i-node, not the data in the file. To make these changes, the owner of the process must have root user authority or own the file.

Creating new file system types

If it is necessary to create a new type of file system, file system helpers and mount helpers must be created.

This section provides information about the implementation specifics and execution syntax of file system helpers and mount helpers.

File system helpers

To enable support of multiple file system types, most file system commands do not contain the code that communicates with individual file systems. Instead, the commands collect parameters, file system names, and other information not specific to one file system type and then pass this information to a back-end program. The back end understands specific information about the relevant file system type and does the detail work of communicating with the file system. Back-end programs used by file system commands are known as *file system helpers* and *mount helpers*.

To determine the appropriate file system helper, the front-end command looks for a helper under the **/sbin/helpers/vfstype/command** file, where *vfstype* matches the file system type found in the **/etc/vfs** file and *command* matches the name of the command being executed. The flags passed to the front-end command are passed to the file system helper.

One required file system helper that needs to be provided, called *fstype*, does not match a command name. This helper is used to identify if a specified logical volume contains a file system of the *vfstype* of the helper.

- The helper returns 0 if the logical volume does not contain a file system of its type. A return value of 0 indicates the logical volume does not contain a log.
- The helper returns 1 if the logical volume does contain a file system of its type and the file system does not need a separate device for a log. A return value of 1 indicates the logical volume does contain a log.
- The helper returns 2 if the logical volume does contain a file system of its type and the file system does need a separate device for a log. If the **-l** flag is specified, the *fstype* helper should check for a log of its file system type on the specified logical volume.

Obsolete file system helper mechanism

This section describes the obsolete file system helper mechanism that was used on previous versions of AIX. This mechanism is still available but should not be used anymore.

File system helper operations

The following table lists the possible operations requested of a helper in the **/usr/include/fshelp.h** file:

Helper operations	Value
#define FSHOP_NULL	0
#define FSHOP_CHECK	1
#define FSHOP_CHGSIZ	2
#define FSHOP_FINDATA	3
#define FSHOP_FREE	4
#define FSHOP_MAKE	5
#define FSHOP_REBUILD	6
#define FSHOP_STATFS	7
#define FSHOP_STAT	8
#define FSHOP_USAGE	9
#define FSHOP_NAMEI	10
#define FSHOP_DEBUG	11

However, the JFS file system supports only the following operations:

Operation Value Corresponding Command

```
#define FSHOP_CHECK      1    fsck
#define FSHOP_CHGSIZ    2    chfs
#define FSHOP_MAKE      5    mkfs
#define FSHOP_STATFS    7    df
#define FSHOP_NAMEI     10   ff
```

File system helper execution syntax

The execution syntax of the file system helper is as follows:

```
OpName OpKey FilsysFileDescriptor PipeFileDescriptor Modeflags
DebugLevel OpFlags
```

OpName

Specifies the *arg0* parameter when the program invokes the helper. The value of the **OpName** field appears in a list of processes (see the **ps** command).

OpKey

Corresponds to the available helper operations. Thus, if the **OpKey** value is 1, the **fsck** (file system check) operation is being requested.

FilsysFileDescriptor

Indicates the file descriptor on which the file system has been opened by the program.

PipeFileDescriptor

Indicates the file descriptor of the pipe (see the **pipe** subroutine) that is open between the original program and the helper program. This channel allows the helper to communicate with the front-end program.

Example: The helper sends an indication of its success or failure through the pipe, which can affect further front-end processing. Also, if the debug level is high enough, the helper can have additional information to send to the original program.

Modeflags

Provides an indication of how the helper is being invoked and can affect the behavior of the helper, especially in regard to error reporting. Mode flags are defined in the **/usr/include/fshelp.h** file:

```
Flags                                Indicator
#define FSHMOD_INTERACT_FLAG         "i"
#define FSHMOD_FORCE_FLAG            "f"
#define FSHMOD_NONBLOCK_FLAG         "n"
#define FSHMOD_PERROR_FLAG           "p"
#define FSHMOD_ERRDUMP_FLAG          "e"
#define FSHMOD_STANDALONE_FLAG       "s"
#define FSHMOD_IGNDEVTYPE_FLAG       "I"
```

Example: The **FSHMOD_INTERACT** flag indicates whether the command is being run interactively (determined by testing the **isatty** subroutine on the standard input). Not every operation uses all (or any) of these modes.

DebugLevel

Determines the amount of debugging information required: the higher the debugging level, the more detailed the information returned.

OpFlags

Includes the actual device (or devices) on which the operation is to be performed and any other options specified by the front end.

Sample helper invocations

The actual invocation of the file system helper for `fsck -fp /user` is as follows:

```
execl("/etc/helpers/v3fshelpers", "fshop_check", "1", "3", "5", "ifp",
      "0", "devices=/dev/lv02,fast,preen,mounted")
```

In this example:

- The executable **execd** is `/etc/helper/v3fshelpers`.
- The name appearing in a listing of processes (**ps** command) is `fshop_check`.
- The operation requested is **FSHOP_CHECK** as represented by the value "1".
- The file system is open on file descriptor "3".
- The pipe through which the helper can send information to the front-end program is open on file descriptor "5".
- The ModeFlags string is `"-ifp"`, indicating interactive mode plus force and perror modes.
- The DebugLevel is 0, so no extra debugging output is returned to the **fsck** command.
- The OpFlags string tells the back-end programs which device the operation is to be performed on (`/dev/lv02`), the options requested (`fast` and `preen`), and notes that the device is mounted. For the **fsck** command, no changes will be made because the **fsck** command does not work on mounted file systems.

Another example of how file system helpers are invoked uses the **mkfs** command. To create a JFS file system on the existing logical volume named `/dev/lv02`, type the following to create a mount point:

```
mkfs /junk
```

If you want to create a file system and know only the device you want to mount it to, type:

```
mkfs /dev/lv02
```

In either case, the following file system helper is invoked:

```
execl ("/etc/helpers/v3fshelpers", "fshop_make", "5", "3", "5", "-ip", \
      "0", "name=/junk,label=/junk,dev=/dev/lv02")
```

The operation requested is now **FSHOP_MAKE**. The mode is `interactive` and `perror`. The OpFlags string includes both the mount point and the device.

Mount helpers

The **mount** command is a front-end program that uses a helper to communicate with specific file systems. Helper programs for the **mount** and **umount** (or **unmount**) commands are called *mount helpers*.

Like other file system-specific commands, the **mount** command collects the parameters and options given at the command line and interprets that information within the context of the file system configuration information found in the `/etc/filesystems` file. Using the information in the `/etc/filesystems` file, the command invokes the appropriate mount helper for the type of file system involved. For example, if the user types the following, the **mount** command checks the `/etc/filesystems` file for the stanza that describes the `/test` file system.

```
mount /test
```

From the `/etc/filesystems` file, the **mount** command determines that the `/test` file system is a remote NFS mount from the node named `host1`. The **mount** command also notes any options associated with the mount.

An example `/etc/filesystems` file stanza is as follows:

```
/test:
    dev          = /export
    vfs          = nfs
```

```
nodename      = host1
options       = ro,fg,hard,intr
```

The file system type (nfs in our example) determines which mount helper to invoke. The command compares the file system type to the first fields in the [/etc/vfs](#) file. The field that matches will have the mount helper as its third field.

Mount helper execution syntax

The following is a sample of the execution syntax of the mount helper:

```
/etc/helpers/nfsmnthelp M 0 host1 /export /test ro,fg,hard,intr
```

Both the **mount** and **unmount** commands have six parameters. The first four parameters are the same for both commands:

operation

Indicates operation requested of the helper. Values are either M (mount operation), Q (query operation), or U (unmount operation). The query operation is obsolete.

debuglevel

Gives the numeric parameter for the **-D** flag. Neither the **mount** nor the **unmount** commands support the **-D** flag, so the value is 0.

nodename

Names the node for a remote mount or a null string for local mounts. The **mount** or **unmount** commands do not invoke a mount helper if the nodename parameter is null.

object

Names the local or remote device, directory, or file that is being mounted or unmounted. Not all file systems support all combinations. For example, most remote file systems do not support device mounts, while most local file systems do not support anything else.

The remaining parameters for the **mount** command are as follows:

mount point

Names the local directory or file where the object is to be mounted.

options

Lists any file system-specific options, separated by commas. Information for this parameter comes from the **options** field of the relevant stanza in the [/etc/filesystems](#) file or from the **-o Options** flag on the command line (**mount -o Options**). The **mount** command also recognizes the **-r** (read-only) flag and converts it to the string **ro** in this field.

The remaining parameters for the **unmount** command are as follows:

vfsNumber

Gives the unique number that identifies the mount being undone. The unique number is returned by the **vmount** call and can be retrieved by calling the **mntctl** or **stat** subroutine. For the mount helper, this parameter is used as the first parameter to the [uvmount](#) subroutine call that actually does the unmount.

flag

Gives the value of the second parameter to the **uvmount** subroutine. The value is 1 if the unmount operation is forced using the **-f** flag (**umount -f**). Otherwise, the value is 0. Not all file systems support forced unmounts.

Logical volume programming

The Logical Volume Manager (LVM) consists of the library of LVM subroutines and the logical volume device driver, described as follows:

- Library of LVM subroutines. These subroutines define volume groups and maintain the logical and physical volumes of volume groups.
- Logical volume device driver. The logical volume device driver is a pseudo-device driver that processes all logical I/O. It exists as a layer between the file system and the disk device drivers. The logical

volume device driver converts a logical address to a physical address, handles mirroring and bad-block relocation, and then sends the I/O request to the specific disk device driver. Interfaces to the logical volume device driver are provided by the **open**, **close**, **read**, **write**, and **ioctl** subroutines.

For a description of the **readx** and **writex** extension parameters and those **ioctl** operations specific to the logical volume device driver, see the *Kernel Extensions and Device Support Programming Concepts*.

For more information about logical volumes, see the *Operating system and device management*.

Library of logical volume subroutines

LVM subroutines define and maintain the logical and physical volumes of a volume group.

System management commands use these subroutines to perform system management for the logical and physical volumes of a system. The programming interface for the library of LVM subroutines is available to provide alternatives to or expand the function of the system management commands for logical volumes.

Note: The LVM subroutines use the **sysconfig** system call, which requires root user authority, to query and update kernel data structures describing a volume group. You must have root user authority to use the services of the LVM subroutine library.

The following services are available:

Services	Description
<u>lvm_querylv</u>	Queries a logical volume and returns all pertinent information.
<u>lvm_querypv</u>	Queries a physical volume and returns all pertinent information.
<u>lvm_queryvg</u>	Queries a volume group and returns pertinent information.
<u>lvm_queryvgs</u>	Queries the volume groups of the system and returns information for groups that are varied online.

J2_CFG_ASSIST ioctl operation

The **J2_CFG_ASSIST** ioctl operation returns the performance statistics of a JFS2 file system.

The **J2_CFG_ASSIST** ioctl operation returns a **cfg_assist** structure, as defined in the **/usr/include/sys/lvdd.h** file. The structure contains the following fields:

Field	Description
<i>throughput</i>	Average throughput of disks under the file system in MB/sec. For supported storage devices, throughput is obtained from the device; otherwise runtime throughput of the file system is returned.
<i>latency</i>	Average latency of all disks under the file system in milliseconds. For supported storage devices, throughput is obtained from the device; otherwise runtime latency of the file system is returned.
<i>flags</i>	Flags to be used. For a list of valid flags, see the /usr/include/sys/lvdd.h file.
<i>vg_max_transfer</i>	The maximum transfer size of the volume group (VG), in KB. The <i>vg_max_transfer</i> field value is the maximum amount of data that can be transferred in one I/O request to the disks of the volume group.
<i>write_atomicity</i>	Write atomicity in bytes. The <i>write_atomicity</i> field value is the largest number of bytes that are not broken up when they are written on aligned boundaries.

The **J2_CFG_ASSIST** ioctl operation returns the following parameters only for supported storage devices; otherwise it returns null values.

Parameter	Description
<i>atomicWriteAlignment</i>	Required alignment for write atomicity in bytes.
<i>ideal_sequential_read_size</i>	Ideal, sequential, read size of the disks under the file system in KB.
<i>ideal_sequential_write_size</i>	Ideal, sequential, write size of the disks under the file system in KB.
<i>ideal_random_read_size</i>	Ideal, random, read size of the disks under the file system in KB.
<i>ideal_random_write_size</i>	Ideal, random, write size of the disks under the file system in KB.
<i>stripsize</i>	The strip size of the disks under the file system in KB. This is the amount of data that is contiguous on a single spindle in the raid array.
<i>stripesize</i>	The value of the <i>Stripesize</i> parameter is in KB. (<i>Stripesize</i> = <i>stripsize</i> × number of spindles in a RAID array - parity.)
<i>parallelism</i>	The number of spindles that comprise the RAID device that can be concurrently read from, or written to, in parallel.

Return values

When you complete this operation, a value of 0 is returned. If the operation fails, a value of -1 is returned and the **errno** global variable is set to one of the following values:

Value	Description
EFAULT	Indicates that the copy of the parameter failed.
ENOMEM	Indicates that the allocation of the memory failed.
EAGAIN	Indicates that the runtime statistics are unavailable for any of the physical volumes under the file system. Try again after more I/O has been issued to the file system.

Floating-point exceptions

This topic provides information about floating-point exceptions and how your programs can detect and handle them.

The Institute of Electrical and Electronics Engineers (IEEE) defines a standard for floating-point exceptions called the IEEE Standard for Binary Floating-Point Arithmetic (IEEE 754). This standard defines five types of floating-point exception that must be signaled when detected:

- Invalid operation
- Division by zero
- Overflow
- Underflow
- Inexact calculation

When one of these exceptions occurs in a user process, it is signaled either by setting a flag or taking a trap. By default, the system sets a status flag in the Floating-Point Status and Control registers (FPSCR), indicating the exception has occurred. Once the status flags are set by an exception, they are cleared only when the process clears them explicitly or when the process ends. The operating system provides subroutines to query, set, or clear these flags.

The system can also cause the floating-point exception signal (**SIGFPE**) to be raised if a floating-point exception occurs. Because this is not the default behavior, the operating system provides subroutines to change the state of the process so the signal is enabled. When a floating-point exception raises the **SIGFPE** signal, the process terminates and produces a core file if no signal-handler subroutine is present in the process. Otherwise, the process calls the signal-handler subroutine.

Floating-point exception subroutines

Floating-point exception subroutines can be used to:

- Change the execution state of the process
- Enable the signaling of exceptions
- Disable exceptions or clear flags
- Determine which exceptions caused the signal
- Test the exception sticky flags

The following subroutines are provided to accomplish these tasks:

Subroutine	Task
<u>fp_any_xcp</u> or <u>fp_divbyzero</u>	Test the exception sticky flags
<u>fp_enable</u> or <u>fp_enable_all</u>	Enable the signaling of exceptions
<u>fp_inexact</u> , <u>fp_invalid_op</u> , <u>fp_iop_convert</u> , <u>fp_iop_infdinf</u> , <u>fp_iop_infmzr</u> , <u>fp_iop_infsinf</u> , <u>fp_iop_invcmp</u> , <u>fp_iop_snan</u> , <u>fp_iop_sqrt</u> , <u>fp_iop_vxsoft</u> , <u>fp_iop_zrdzr</u> , or <u>fp_overflow</u>	Test the exception sticky flags
<u>fp_sh_info</u>	Determines which exceptions caused the signal
<u>fp_sh_set_stat</u>	Disables exceptions or clear flags
<u>fp_trap</u>	Changes the execution state of the process
<u>fp_underflow</u>	Tests the exception sticky flags
<u>sigaction</u>	Installs signal handler

Floating-point trap handler operation

To generate a trap, a program must change the execution state of the process using the [**fp_trap**](#) subroutine and enable the exception to be trapped using the [**fp_enable**](#) or [**fp_enable_all**](#) subroutine.

Changing the execution state of the program may slow performance because floating-point trapping causes the process to execute in serial mode.

When a floating-point trap occurs, the **SIGFPE** signal is raised. By default, the **SIGFPE** signal causes the process to terminate and produce a core file. To change this behavior, the program must establish a signal handler for this signal. See the [sigaction](#), [sigvec](#), or [signal](#) subroutines for more information on signal handlers.

Exceptions: disabled and enabled comparison

Refer to the following lists for an illustration of the differences between the disabled and enabled processing states and the subroutines that are used.

Exceptions-disabled model

The following subroutines test exception flags in the disabled processing state:

- [**fp_any_xcp**](#)
- [**fp_clr_flag**](#)
- [**fp_divbyzero**](#)
- [**fp_inexact**](#)
- [**fp_invalid_op**](#)
- [**fp_iop_convert**](#)
- [**fp_iop_infdinf**](#)

- **fp_iop_infmzr**
- **fp_iop_infsi**
- **fp_iop_invcmp**
- **fp_iop_snan**
- **fp_iop_sqrt**
- **fp_iop_vxsoft**
- **fp_iop_zrdzr**
- **fp_overflow**
- **fp_underflow**

Exceptions-enabled model

The following subroutines function in the enabled processing state:

Subroutine	Processing state
fp_enable or fp_enable_all	Enable the signaling of exceptions
fp_sh_info	Determines which exceptions caused the signal
fp_sh_set_stat	Disables exceptions or clear flags
fp_trap	Changes the execution state of the process
sigaction	Installs signal handler

Imprecise trapping modes

Some systems have *imprecise trapping modes*. This means the hardware can detect a floating-point exception and deliver an interrupt, but processing may continue while the signal is delivered. As a result, the instruction address register (IAR) is at a different instruction when the interrupt is delivered.

Imprecise trapping modes cause less performance degradation than *precise trapping mode*. However, some recovery operations are not possible, because the operation that caused the exception cannot be determined or because subsequent instruction may have modified the argument that caused the exception.

To use imprecise exceptions, a signal handler must be able to determine if a trap was precise or imprecise.

Precise traps

In a precise trap, the instruction address register (IAR) points to the instruction that caused the trap. A program can modify the arguments to the instruction and restart it, or fix the result of the operation and continue with the next instruction. To continue, the IAR must be incremented to point to the next instruction.

Imprecise traps

In an imprecise trap, the IAR points to an instruction beyond the one that caused the exception. The instruction to which the IAR points has not been started. To continue execution, the signal handler does not increment the IAR.

To eliminate ambiguity, the `trap_mode` field is provided in the **fp_sh_info** structure. This field specifies the trapping mode in effect in the user process when the signal handler was entered. This information can also be determined by examining the Machine Status register (MSR) in the **mstsave** structure.

The **fp_sh_info** subroutine allows a floating-point signal handler to determine if the floating-point exception was precise or imprecise.

Note: Even when precise trapping mode is enabled some floating-point exceptions may be imprecise (such as operations implemented in software). Similarly, in imprecise trapping mode some exceptions may be precise.

When using imprecise exceptions, some parts of your code may require that all floating-point exceptions are reported before proceeding. The **fp_flush_imprecise** subroutine is provided to accomplish this. It is also recommended that the **atexit** subroutine be used to register the **fp_flush_imprecise** subroutine to run at program exit. Running at exit ensures that the program does not exit with unreported imprecise exceptions.

Hardware-specific subroutines

Some systems have hardware instructions to compute the square root of a floating-point number and to convert a floating-point number to an integer. Models not having these hardware instructions use software subroutines to do this. Either method can cause a trap if the invalid operation exception is enabled. The software subroutines report, through the **fp_sh_info** subroutine, that an imprecise exception occurred, because the IAR does not point to a single instruction that can be restarted to retry the operation.

Example of a floating-point trap handler

```

/*
 * This code demonstrates a working floating-point exception
 * trap handler. The handler simply identifies which
 * floating-point exceptions caused the trap and return.
 * The handler will return the default signal return
 * mechanism longjmp().
 */

#include <signal.h>
#include <setjmp.h>
#include <fpexp.h>
#include <fptrap.h>
#include <stdlib.h>
#include <stdio.h>

#define EXIT_BAD -1
#define EXIT_GOOD 0

/*
 * Handshaking variable with the signal handler. If zero,
 * then the signal handler returns via the default signal
 * return mechanism; if non-zero, then the signal handler
 * returns via longjmp.
 */
static int fpsigexit;
#define SIGRETURN_EXIT 0
#define LONGJUMP_EXIT 1

static jmp_buf jump_buffer; /* jump buffer */
#define JMP_DEFINED 0 /* setjmp rc on initial call */
#define JMP_FPE 2 /* setjmp rc on return from */
/* signal handler */

/*
 * The fp_list structure allows text descriptions
 * of each possible trap type to be tied to the mask
 * that identifies it.
 */

typedef struct
{
    fpflag_t mask;
    char *text;
} fp_list_t;

/* IEEE required trap types */

fp_list_t
trap_list[] =
{
    { FP_INVALID, "FP_INVALID"},
    { FP_OVERFLOW, "FP_OVERFLOW"},

```

```

    { FP_UNDERFLOW, "FP_UNDERFLOW"},
    { FP_DIV_BY_ZERO, "FP_DIV_BY_ZERO"},
    { FP_INEXACT, "FP_INEXACT"}
};

/* INEXACT detail list -- this is an system extension */

fp_list_t
detail_list[] =
{
    { FP_INV_SNAN, "FP_INV_SNAN" },
    { FP_INV_ISI, "FP_INV_ISI" },
    { FP_INV_IDI, "FP_INV_IDI" },
    { FP_INV_ZDZ, "FP_INV_ZDZ" },
    { FP_INV_IMZ, "FP_INV_IMZ" },
    { FP_INV_CMP, "FP_INV_CMP" },
    { FP_INV_SQRT, "FP_INV_SQRT" },
    { FP_INV_CVI, "FP_INV_CVI" },
    { FP_INV_VXSOFT, "FP_INV_VXSOFT" }
};

/*
 * the TEST_IT macro is used in main() to raise
 * an exception.
 */

#define TEST_IT(WHAT, RAISE_ARG) \
{ \
    puts(strcat("testing: ", WHAT)); \
    fp_clr_flag(FP_ALL_XCP); \
    fp_raise_xcp(RAISE_ARG); \
}

/*
 * NAME: my_div
 *
 * FUNCTION: Perform floating-point division.
 */

double
my_div(double x, double y)
{
    return x / y;
}

/*
 * NAME: sigfpe_handler
 *
 * FUNCTION: A trap handler that is entered when
 * a floating-point exception occurs. The
 * function determines what exceptions caused
 * the trap, prints this to stdout, and returns
 * to the process which caused the trap.
 *
 * NOTES: This trap handler can return either via the
 * default return mechanism or via longjmp().
 * The global variable fpsigexit determines which.
 *
 * When entered, all floating-point traps are
 * disabled.
 *
 * This sample uses printf(). This should be used
 * with caution since printf() of a floating-
 * point number can cause a trap, and then
 * another printf() of a floating-point number
 * in the signal handler will corrupt the static
 * buffer used for the conversion.
 *
 * OUTPUTS: The type of exception that caused
 * the trap.
 */

```



```

static void
sigfpe_handler(int sig,
               int code,
               struct sigcontext *SCP)
{
    struct mtsave * state = &SCP->sc_jmpbuf.jmp_context;
    fp_sh_info_t flt_context; /* structure for fp_sh_info()
                             /* call */
    int i; /* loop counter */
    extern int fpsigexit; /* global handshaking variable */
    extern jmp_buf jump_buffer /* */

    /*
     * Determine which floating-point exceptions caused
     * the trap. fp_sh_info() is used to build the floating signal
     * handler info structure, then the member
     * flt_context.trap can be examined. First the trap type is
     * compared for the IEEE required traps, and if the trap type
     * is an invalid operation, the detail bits are examined.
     */

    fp_sh_info(SCP, &flt_context, FP_SH_INFO_SIZE);

```

```

static void
sigfpe_handler(int sig,
               int code,
               struct sigcontext *SCP)
{
    struct mtsave * state = &SCP->sc_jmpbuf.jmp_context;
    fp_sh_info_t flt_context; /* structure for fp_sh_info()
                             /* call */
    int i; /* loop counter */
    extern int fpsigexit; /* global handshaking variable */
    extern jmp_buf jump_buffer; /* */

    /*
     * Determine which floating-point exceptions caused
     * the trap. fp_sh_info() is used to build the floating signal
     * handler info structure, then the member
     * flt_context.trap can be examined. First the trap type is
     * compared for the IEEE required traps, and if the trap type
     * is an invalid operation, the detail bits are examined.
     */

    fp_sh_info(SCP, &flt_context, FP_SH_INFO_SIZE);

```

```

for (i = 0; i < (sizeof(trap_list) / sizeof(fp_list_t)); i++)
{
    if (flt_context.trap & trap_list[i].mask)
        (void) printf("Trap caused by %s error\n", trap_list[i].text);
}

```

```

if (flt_context.trap & FP_INVALID)
{
    for (i = 0; i < (sizeof(detail_list) / sizeof(fp_list_t)); i++)
    {
        if (flt_context.trap & detail_list[i].mask)
            (void) printf("Type of invalid op is %s\n", detail_list[i].text);
    }
}

```

```

/* report which trap mode was in effect */

```

```

switch (flt_context.trap_mode)
{
    case FP_TRAP_OFF:
        puts("Trapping Mode is OFF");
        break;

    case FP_TRAP_SYNC:
        puts("Trapping Mode is SYNC");
        break;

    case FP_TRAP_IMP:
        puts("Trapping Mode is IMP");
        break;
}

```

```

case FP_TRAP_IMP_REC:
    puts("Trapping Mode is IMP_REC");
    break;

default:
    puts("ERROR: Invalid trap mode");
}

```

```

if (fpsigexit == LONGJUMP_EXIT)
{
    /*
    * Return via longjmp. In this instance there is no need to
    * clear any exceptions or disable traps to prevent
    * recurrence of the exception, because on return the
    * process will have the signal handler's floating-point
    * state.
    */
    longjmp(jump_buffer, JMP_FPE);
}
else
{
    /*
    * Return via default signal handler return mechanism.
    * In this case you must take some action to prevent
    * recurrence of the trap, either by clearing the
    * exception bit in the fpSCR or by disabling the trap.
    * In this case, clear the exception bit.
    * The fp_sh_set_stat routine is used to clear
    * the exception bit.
    */
}

```

```

fp_sh_set_stat(SCP, (flt_context.fpSCR & ((fpstat_t) ~flt_context.trap)));

```

```

/*
* Increment the iar of the process that caused the trap,
* to prevent re-execution of the instruction.
* The FP_IAR_STAT bit in flt_context.flags indicates if
* state->iar points to an instruction that has logically
* started. If this bit is true, state->iar points to
* an operation that has started and will cause another
* exception if it runs again. In this case you want to
* continue execution and ignore the results of that
* operation, so the iar is advanced to point to the
* next instruction. If the bit is false, the iar already
* points to the next instruction that must run.
*/

if ( flt_context.flags & FP_IAR_STAT )
{
    puts("Increment IAR");
    state->iar += 4;
}
}
return;
}

/*
* NAME: main
*
* FUNCTION: Demonstrate the sigfpe_handler trap handler.
*
*/

```

```

int
main(void)
{
    struct sigaction response;
    struct sigaction old_response;
    extern int fpsigexit;
    extern jmp_buf jump_buffer;
    int jump_rc;
    int trap_mode;
    double arg1, arg2, r;

    /*
    * Set up for floating-point trapping. Do the following:
    */
}

```

```

* 1. Clear any existing floating-point exception flags.
* 2. Set up a SIGFPE signal handler.
* 3. Place the process in synchronous execution mode.
* 4. Enable all floating-point traps.
*/

```

```

fp_clr_flag(FP_ALL_XCP);
(void) sigaction(SIGFPE, NULL, &old_response);
(void) sigemptyset(&response.sa_mask);
response.sa_flags = FALSE;
response.sa_handler = (void (*)(int)) sigfpe_handler;
(void) sigaction(SIGFPE, &response, NULL);
fp_enable_all();

```

```

/*
* Demonstrate trap handler return via default signal handler
* return. The TEST_IT macro will raise the floating-point
* exception type given in its second argument. Testing
* is done in this case with precise trapping, because
* it is supported on all platforms to date.
*/

```

```

trap_mode = fp_trap(FP_TRAP_SYNC);
if ((trap_mode == FP_TRAP_ERROR) ||
    (trap_mode == FP_TRAP_UNIMPL))
{
    printf("ERROR: rc from fp_trap is %d\n",
          trap_mode);
    exit(-1);
}

(void) printf("Default signal handler return: \n");
fp_sigexit = SIGRETURN_EXIT;

```

```

TEST_IT("div by zero", FP_DIV_BY_ZERO);
TEST_IT("overflow", FP_OVERFLOW);
TEST_IT("underflow", FP_UNDERFLOW);
TEST_IT("inexact", FP_INEXACT);

```

```

TEST_IT("signaling nan", FP_INV_SNaN);
TEST_IT("INF - INF", FP_INV_ISI);
TEST_IT("INF / INF", FP_INV_IDI);
TEST_IT("ZERO / ZERO", FP_INV_ZDZ);
TEST_IT("INF * ZERO", FP_INV_IMZ);
TEST_IT("invalid compare", FP_INV_CMP);
TEST_IT("invalid sqrt", FP_INV_SQRT);
TEST_IT("invalid conversion", FP_INV_CVI);
TEST_IT("software request", FP_INV_VXSOFT);

```

```

/*
* Next, use fp_trap() to determine what the
* the fastest trapmode available is on
* this platform.
*/

```

```

trap_mode = fp_trap(FP_TRAP_FASTMODE);
switch (trap_mode)
{
    case FP_TRAP_SYNC:
        puts("Fast mode for this platform is PRECISE");
        break;

    case FP_TRAP_OFF:
        puts("This platform doesn't support trapping");
        break;

    case FP_TRAP_IMP:
        puts("Fast mode for this platform is IMPRECISE");
        break;
}

```

```

case FP_TRAP_IMP_REC:
    puts("Fast mode for this platform is IMPRECISE RECOVERABLE");
    break;

default:
    printf("Unexpected return code from fp_trap(FP_TRAP_FASTMODE): %d\n",
        trap_mode);
    exit(-2);
}

/*
 * if this platform supports imprecise trapping, demonstrate this.
 */

trap_mode = fp_trap(FP_TRAP_IMP);
if (trap_mode != FP_TRAP_UNIMPL)
{
    puts("Demonstrate imprecise FP exceptions");
    arg1 = 1.2;
    arg2 = 0.0;
    r = my_div(arg1, arg2);
    fp_flush_imprecise();
}

/* demonstrate trap handler return via longjmp().
 */

(void) printf("longjmp return: \n");
fptestexit = LONGJUMP_EXIT;
jump_rc = setjmp(jump_buffer);

switch (jump_rc)
{
case JMP_DEFINED:
    (void) printf("setjmp has been set up; testing ...\n");
    TEST_IT("div by zero", FP_DIV_BY_ZERO);
    break;

case JMP_FPE:
    (void) printf("back from signal handler\n");
    /*
     * Note that at this point the process has the floating-
     * point status inherited from the trap handler. If the
     * trap handler did not enable trapping (as the example
     * did not) then this process at this point has no traps
     * enabled. We create a floating-point exception to
     * demonstrate that a trap does not occur, then re-enable
     * traps.
     */

    (void) printf("Creating overflow; should not trap\n");
    TEST_IT("Overflow", FP_OVERFLOW);
    fp_enable_all();
    break;

default:
    (void) printf("unexpected rc from setjmp: %d\n", jump_rc);
    exit(EXIT_BAD);
}
exit(EXIT_GOOD);
}

```

Related information

[fp_clr_flag](#), [fp_set_flag](#), [fp_read_flag](#), [fp_swag](#)
[fp_raise_xcp](#)
[sigaction](#), [sigvec](#), [snap Command](#)

Input and output handling

This topic provides an introduction to programming considerations for input and output handling and the input and output (I/O) handling subroutines.

The I/O library subroutines can send data to or from either devices or files. The system treats devices as if they were I/O files. For example, you must also open and close a device just as you do a file.

Some of the subroutines use standard input and standard output as their I/O channels. For most of the subroutines, however, you can specify a different file for the source or destination of the data transfer. For some subroutines, you can use a file pointer to a structure that contains the name of the file; for others, you can use a file descriptor (that is, the positive integer assigned to the file when it is opened).

The I/O subroutines stored in the C Library (**libc.a**) provide stream I/O. To access these stream I/O subroutines, you must include the **stdio.h** file by using the following statement:

```
#include <stdio.h>
```

Some of the I/O library subroutines are macros defined in a header file and some are object modules of functions. In many cases, the library contains a macro and a function that do the same type of operation. Consider the following when deciding whether to use the macro or the function:

- You cannot set a breakpoint for a macro using the **dbx** program.
- Macros are usually faster than their equivalent functions because the preprocessor replaces the macros with actual lines of code in the program.
- Macros result in larger object code after being compiled.
- Functions can have side effects to avoid.

The files, commands, and subroutines used in I/O handling provide the following interfaces:

Low-level

The low-level interface provides basic open and close functions for files and devices.

Stream

The stream interface provides read and write I/O for pipes and FIFOs.

Terminal

The terminal interface provides formatted output and buffering.

Asynchronous

The asynchronous interface provides concurrent I/O and processing.

Input Language

The input language interface uses the **lex** and **yacc** commands to generate a lexical analyzer and a parser program for interpreting I/O.

Low-level I/O interfaces

Low-level I/O interfaces are direct entry points into a kernel, providing functions such as opening files, reading to and writing from files, and closing files.

The **line** command provides the interface that allows one line from standard input to be read, and the following subroutines provide other low-level I/O functions:

open, openx, or creat

Prepare a file, or other path object, for reading and writing by means of an assigned file descriptor

read, readx, readv, or readvx

Read from an open file descriptor

write, writex, writev, or writevx

Write to an open file descriptor

close

Relinquish a file descriptor

The **open** and **creat** subroutines set up entries in three system tables. A file descriptor indexes the first table, which functions as a per process data area that can be accessed by read and write subroutines. Each entry in this table has a pointer to a corresponding entry in the second table.

The second table is a per-system database, or file table, that allows an open file to be shared among several processes. The entries in this table indicate if the file was open for reading, writing, or as a pipe, and when the file was closed. There is also an offset to indicate where the next read or write will take place and a final pointer to indicates entry to the third table, which contains a copy of the file's i-node.

The file table contains entries for every instance of an **open** or **create** subroutine on the file, but the i-node table contains only one entry for each file.

Note: While processing an **open** or **creat** subroutine for a special file, the system always calls the device's **open** subroutine to allow any special processing (such as rewinding a tape or turning on a data-terminal-ready modem lead). However, the system uses the **close** subroutine only when the last process closes the file (that is, when the i-node table entry is deallocated). This means that a device cannot maintain or depend on a count of its users unless an exclusive-use device (that prevents a device from being reopened before its closed) is implemented.

When a read or write operation occurs, the user's arguments and the file table entry are used to set up the following variables:

- User address of the I/O target area
- Byte count for the transfer
- Current location in the file

If the file referred to is a character-type special file, the appropriate read or write subroutine is called to transfer data, as well as update the count and current location. Otherwise, the current location is used to calculate a logical block number in the file.

If the file is an ordinary file, the logical block number must be mapped to a physical block number. A block-type special file need not be mapped. The resulting physical block number is used to read or write the appropriate device.

Block device drivers can provide the ability to transfer information directly between the user's core image and the device in block sizes as large as the caller requests without using buffers. The method involves setting up a character-type special file corresponding to the raw device and providing read and write subroutines to create a private, non-shared buffer header with the appropriate information. Separate open and close subroutines can be provided, and a special-function subroutine can be called for magnetic tape.

Stream I/O interfaces

Stream I/O interfaces provide data as a stream of bytes that is not interpreted by the system, which offers more efficient implementation for networking protocols than character I/O processing. No record boundaries exist when reading and writing using stream I/O. For example, a process reading 100 bytes from a pipe cannot determine if the process that wrote the data into the pipe did a single write of 100 bytes, or two writes of 50 bytes, or even if the 100 bytes came from two different processes.

Stream I/Os can be pipes or FIFOs (first-in, first-out files). FIFOs are similar to pipes because they allow the data to flow only one way (left to right). However, a FIFO can be given a name and can be accessed by unrelated processes, unlike a pipe. FIFOs are sometimes referred to as *named pipes*. Because it has a name, a FIFO can be opened using the standard I/O **fopen** subroutine. To open a pipe, you must call the **pipe** subroutine, which returns a file descriptor, and the standard I/O **fdopen** subroutine to associate an open file descriptor with a standard I/O stream.

Note: Stream I/O interfaces buffer data at the user level and cannot write the data until the **fclose** or **fflush** Subroutine is performed, which might lead to unexpected results when mixed with file I/O such as `read()` or `write()`.

Stream I/O interfaces are accessed through the following subroutines and macros:

fclose

Closes a stream

feof, ferror, clearerr, or fileno

Check the status of a stream

fflush

Write all currently buffered characters from a stream

fopen, freopen, or fdopen

Open a stream

fread or fwrite

Perform binary input

fseek, rewind, ftell, fgetpos, or fsetpos

Reposition the file pointer of a stream

getc, fgetc, getchar, or getw

Get a character or word from an input stream

gets or fgets

Get a string from a stream

getwc, fgetwc, or getwchar

Get a wide character from an input stream

getws or fgets

Get a string from a stream

printf, fprintf, sprintf, vsprintf, vprintf, vsprintf, or vwsprintf

Print formatted output

putc, putchar, fputc, or putw

Write a character or a word to a stream

puts or fputs

Write a string to a stream

putwc, putwchar, or fputwc

Write a character or a word to a stream

putws or fputws

Write a wide character string to a stream

scanf, fscanf, sscanf, or wscanf

Convert formatted input

setbuf, setvbuf, setbuffer, or setlinebuf

Assign buffering to a stream

ungetc or ungetwc

Push a character back into the input stream

Terminal I/O interfaces

Terminal I/O interfaces operate between a process and the kernel, providing functions such as buffering and formatted output. Every terminal and pseudo-terminal has a `tty` structure that contains the current process group ID. This field identifies the process group to receive the signals associated with the terminal. Terminal I/O interfaces can be accessed through the **iosat** command, which monitors I/O system device loading, and the **uprintfd** daemon, which allows kernel messages to be written to the system console.

Terminal characteristics can be enabled or disabled through the following subroutines:

cfgetospeed, cfsetospeed, cfgetispeed, or cfsetispeed

Get and set input and output baud rates

ioctl

Performs control functions associated with the terminal

termdef

Queries terminal characteristics

tcdrain

Waits for output to complete

tcflow

Performs flow control functions

tcflush

Discards data from the specified queue

tcgetpgrp

Gets foreground process group ID

tcsendbreak

Sends a break on an asynchronous serial data line

tcsetattr

Sets terminal state

ttylock, ttywait, ttyunlock, or ttylocked

Control tty locking functions

ttyname or isatty

Get the name of a terminal

ttyslot

Finds the slot in the **utmp** file for the current user

Asynchronous I/O interfaces

Asynchronous I/O subroutines allow a process to start an I/O operation and have the subroutine return immediately after the operation is started or queued. Another subroutine is required to wait for the operation to complete (or return immediately if the operation is already finished). This means that a process can overlap its execution with its I/O or overlap I/O between different devices. Although asynchronous I/O does not significantly improve performance for a process that is reading from a disk file and writing to another disk file, asynchronous I/O can provide significant performance improvements for other types of I/O driven programs, such as programs that dump a disk to a magnetic tape or display an image on an image display.

Although not required, a process that is performing asynchronous I/O can tell the kernel to notify it when a specified descriptor is ready for I/O (also called *signal-driven I/O*). When using LEGACY AIO, the kernel notifies the user process with the SIGIO signal. When using POSIX AIO, the **sigevent** structure is used by the programmer to determine which signal for the kernel to use to notify the user process. Signals include **SIGIO**, **SIGUSR1**, and **SIGUSR2**.

To use asynchronous I/O, a process must perform the following steps:

1. Establish a handler for the **SIGIO** signal. This step is necessary only if notification by the signal is requested.
2. Set the process ID or the process group ID to receive the **SIGIO** signals. This step is necessary only if notification by the signal is requested.
3. Enable asynchronous I/O. The system administrator usually determines whether asynchronous I/O is loaded (enabled). Enabling occurs at system startup.

The following asynchronous I/O subroutines are provided:

aio_cancel

Cancels one or more outstanding asynchronous I/O requests

aio_error

Retrieves the error status of an asynchronous I/O request

aio_fsync

Synchronizes asynchronous files.

aio_nwait

Suspends the calling process until a certain number of asynchronous I/O requests are completed.

aio_read

Reads asynchronously from a file descriptor

aio_return

Retrieves the return status of an asynchronous I/O request

aio_suspend

Suspends the calling process until one or more asynchronous I/O requests is completed

aio_write

Writes asynchronously to a file descriptor

lio_listio

Initiates a list of asynchronous I/O requests with a single call

poll or select

Check I/O status of multiple file descriptors and message queues

For use with the **poll** subroutine, the following header files are supplied:

poll.h

Defines the structures and flags used by the **poll** subroutine

aio.h

Defines the structure and flags used by the **aio_read**, **aio_write**, and **aio_suspend** subroutines

Storage protect keys

Storage protect keys provide a mechanism for you to improve the reliability of your programs.

Protect keys apply to memory pages and work at the page level of granularity, similar to the `mprotect` subroutine, which can be used to read- or write-protect one or more pages. However, with storage keys you can mark sections of your data for specific levels of read and write access protection. Protection by storage keys is a function not only of the data page, but also of the thread attempting access. You can enable certain well-defined code paths to access data that is unavailable to your larger program, thereby encapsulating critical program data and protecting it against accidental damage.

Because access to key-protected pages is an attribute of the running thread, this mechanism extends naturally to multithreaded applications, but with the restriction that these use only 1:1 (or system scope) pthreads. The `mprotect` subroutine approach does not work reliably in a multithreaded environment, because you have to remove protection for all threads when you want to grant access to any thread. You can use both mechanisms simultaneously, and both are fully enforced; therefore, your program cannot write to a write-protected page even if a protect key would otherwise allow this.

Protect keys sample uses include:

- Encapsulate your program's private data completely, limiting access to just selected code paths.
- Protect your program's private data from accidental damage by always running with read access granted, but granting write access only when you intend to modify the data. This can be especially useful when code in a core engine allows calls out to untrusted code.
- When multiple private keys are available, additional granularity of data protection is possible.

You can simplify debugging by designing your application with key protection in mind. Setting a page's protect key and setting your active user keyset are both system calls, and therefore relatively expensive operations. You should design your program so that the frequency of these operations is not excessive.

User protect keys

The following guidelines and considerations apply when using protect keys:

- Pages that are exported read-only from the kernel will continue to be visible to your program. These pages have a protect key of **UKEY_SYSTEM**. This protect key is not a protect key that is under your program's control, but is always accessible by your program.
- All of your program's memory pages initially have the user public key assigned to them. As noted above, access to key 0 storage is always granted, making this the *user public key*.

- You can set protect keys only for your normal and shared data. You cannot, for example, protect library data, low memory shared with the kernel, or program text.
- Depending on the underlying hardware and administrative choice, only a limited number of user private keys (typically just one) are available. When your program assigns a private key to one or more of its pages, the data in those pages is no longer available by default. You must explicitly grant read or write access to this data by surrounding code paths that require access with calls to a new service to manage your active user keyset.
- The physical hardware likely supports additional protect keys that are not available for use as user protect keys.
- No special privilege is needed to assign protect keys to a page. The only requirement is current write access to the page.
- There is no control of execute authority with protect keys.

If your program accesses key protected data in violation of the access rights expressed in its active user keyset, it receives a **SIGSEGV** signal, as is already the case for violating read- or write-protected pages. If you choose to handle this signal, be aware that signal handlers are invoked without access to private keys. Signal handling code must add any needed access rights to the active user keyset before referencing key-protected data.

Child processes, created by the `fork` system call, logically inherit their parent's memory and running state. This includes the protect key associated with each page, as well as the parent thread's active user keyset at the time of fork.

Regions protected by user keys

User protect keys can protect pages in the following regions:

- Data region
- Default stack region
- `mmap`ed regions
- Shared memory attached with the `shmat()` subroutine, except as below
- These categories of pages cannot use protect keys:
 - `shmat`ed files and pinned shared memory
 - Large (nonpageable) pages
 - Program text
 - Low memory shared read-only with the kernel

System prerequisites for key protection

Storage key protection is a hardware-specific privileged page protection mechanism that is made available by the AIX kernel for use in application programs. To use this feature, your system must:

- Be running on physical hardware that provides storage key protection
- Be running the 64-bit kernel
- Enable the use of user protect keys

Program prerequisites for key protection

To use user keys, your program must:

- Declare itself user-key aware and determine how many user protect keys are available, if any, with the `ukey_enable` subroutine.
- Organize its protected data on page boundaries.
- Assign a private key to each page you want to protect with the `ukey_protect` subroutine.
- If you protect the `malloc`'d data, remember to unprotect it before you free it.
- Prepare one or more keysets with the `ukeyset_init` subroutine.

- Possibly add the required keys to your keyset with the `ukeyset_add_key` subroutine, to enable future read or write accesses as required.
- Make a keyset active with the `ukeyset_activate` subroutine to grant the access rights defined in a keyset.

Your program must not:

- Include any M:N (process scope) pthreads
- Be able to have a checkpoint performed on it (for example have `CHECKPOINT=yes` in the environment)

Note: When a program is user-key aware, it has additional context associated with it to represent its active user keyset. This can be seen in:

- Signal handlers receiving a `ucontext_t` structure. The previously active user keyset is in `ucontext_t.__extctx.__ukeys`, an array of two ints containing a 64-bit user keyset value
- User context structures compiled with `__EXTABI__` defined (used by `setcontext`, `getcontext`, `makecontext`, `swapcontext`)

Subroutines

The following new AIX kernel subroutines are provided for using protect keys:

Subroutine	Description
<code>sysconf</code>	Use with <code>_SC_AIX_UKEYS</code> to determine the number of user keys supported (can be called on older versions of AIX)
<code>ukey_enable</code>	Enable the user-key aware programming environment for your process, and report how many user keys are available
<code>ukeyset_init</code>	Initialize a user keyset, which will represent a set of access rights to your private key or keys
<code>ukeyset_add_key</code>	Add read or write access, or both for a specified key to a keyset
<code>ukeyset_remove_key</code>	Remove or write access, or both for a specified key from a keyset
<code>ukeyset_add_set</code>	Add all the access rights in one keyset to another
<code>ukeyset_remove_set</code>	Remove all the access rights in one keyset from another
<code>ukeyset_activate</code>	Apply the access rights in a keyset to the running thread
<code>ukeyset_ismember</code>	Test if a given access right is contained in a keyset
<code>ukey_setjmp</code>	Extended form of <code>setjmp</code> that preserves the active keyset (uses a <code>ukey_jmp_buf</code> structure)
<code>pthread_attr_getukeyset_np</code>	Get the keyset attribute of a pthread
<code>pthread_attr_setukeyset_np</code>	Set the keyset attribute for a pthread
<code>ukey_protect</code>	Set a user protect key for a page-aligned range of user memory
<code>ukey_getkey</code>	Retrieve the user protect key for a specified address

Debugging

The `dbx` command adds limited support for protect keys:

- When debugging a running program:
 - The `ukeyset` subcommand displays the active keyset.
 - The `ukeyvalue` subcommand displays the protect key associated with a given memory location.
- When debugging a core file, the `ukeyexcept` subcommand reports the active keyset, effective address of the key exception, and the storage key involved.

Hardware details

The active user keyset in the running context of a key-aware thread parallels the actual hardware authority mask register (AMR) in format, represented by the `ukeyset_t` abstract data type. This information is provided for debugging purposes only. Use only the defined programming services to set up your active user keyset.

- The AMR is a 64-bit register comprising 32-bit pairs, one pair per key, for a maximum of 32 keys numbered 0 through 31.
 - The first bit of each pair represents write access to the corresponding numbered key.
 - Similarly, the second bit of each pair represents read access to the corresponding numbered key.
- A bit value of 0 grants the corresponding access, and a bit value of 1 denies it.
- The bit pair granting access to key 0 is not controlled by your program. User key 0 is the *user public key*, and all threads always have full access to data in this key, without regard to your settings in the active user keyset.
- All the other bit pairs represent *user private keys*, which, subject to availability, you can use to protect your data as you see fit.

Sample program

The following is a sample user-key aware program:

```
#include <stdio.h>
#include <stdlib.h>
#include <malloc.h>
#include <sys/ukeys.h>
#include <sys/sysexp.h>
#include <sys/signal.h>
#include <sys/vminfo.h>

#define ROUND_UP(size,psize)    ((size)+(psize)-1 & ~((psize)-1))

/*
 * This is an example skeleton for a user key aware program.
 *
 * The private_data_1 structure will map a malloc'd key protected area
 * which the main program can access freely, while the "untrusted"
 * subroutine will only have read access.
 */
struct private_data_1 {
    int some_data;
};
struct private_data_1 *p1;      /* pointer to structure for protected data */
ukeyset_t keyset1RW;          /* keyset to give signal handler access */

/*
 * The untrusted function here should successfully read protected data.
 *
 * When the count is 0, it will just return so the caller can write
 * the incremented value back to the protected field.
 *
 * When the count is 1, it will try to update the protected field itself.
 * This should result in a SIGSEGV.
 */
int untrusted(struct private_data_1 *p1) {
    int count = p1->some_data;    /* We can read protected data */
    if (count == 1)
        p1->some_data = count; /* But should not be able to write it */
    return count + 1;
}

/*
 * Signal handler to catch the deliberate protection violation in the
 * untrusted function above when count == 1.
 * Note that the handler is entered with NO access to our private data.
 */
void handler(int signo, siginfo_t *sip, void *ucp) {
    printf("siginfo: signo %d code %d\n", sip->si_signo, sip->si_code);
    (void)ukeyset_activate(keyset1RW, UKA_REPLACE_KEYS);
    exit(1);
}
```

```

main() {
    int nkeys;
    int pagesize = 4096;           /* hardware data page size */
    int padded_protsize_1;       /* page padded size of protected data */
    struct vm_page_info page_info;

    ukey_t key1 = UKEY_PRIVATE1;
    ukeyset_t keyset1W, oldset;
    int rc;
    int count = 0;

    struct sigaction sa;

    /*
     * Attempt to become user key aware.
     */
    nkeys = ukey_enable();
    if (nkeys == -1) {
        perror("ukey_enable");
        exit(1);
    }
    assert(nkeys >= 2);

    /*
     * Determine the data region page size.
     */
    page_info.addr = (long)&p1;    /* address in data region */
    rc = vmgetinfo(&page_info, VM_PAGE_INFO, sizeof(struct vm_page_info));

    if (rc)
        perror("vmgetinfo");
    else
        pagesize = page_info.pagesize; /* pick up actual page size */

    /*
     * We need to allocate page aligned, page padded storage
     * for any area that is going to be key protected.
     */
    padded_protsize_1 = ROUND_UP(sizeof(struct private_data_1), pagesize);
    rc = posix_memalign((void **)&p1, pagesize, padded_protsize_1);
    if (rc) {
        perror("posix_memalign");
        exit(1);
    }

    /*
     * Initialize the private data.
     * We can do this before protecting it if we want.
     *
     * Note that the pointer to the private data is in public storage.
     * We only protect the data itself.
     */
    p1->some_data = count;

    /*
     * Construct keysets to use to access the protected structure.
     * Note that these keysets will be in public storage.
     */
    rc = ukeyset_init(&keyset1W, 0);
    if (rc) {
        perror("ukeyset_init");
        exit(1);
    }

    rc = ukeyset_add_key(&keyset1W, key1, UK_WRITE); /* WRITE */
    if (rc) {
        perror("ukeyset_add_key 1W");
        exit(1);
    }

    keyset1RW = keyset1W;
    rc = ukeyset_add_key(&keyset1RW, key1, UK_READ); /* R/W */
    if (rc) {
        perror("ukeyset_add_key 1R");
        exit(1);
    }

    /*
     * Restrict access to the private data by applying a private key
     * to the page(s) containing it.
     */
    rc = ukey_protect(p1, padded_protsize_1, key1);
}

```

```

if (rc) {
    perror("ukey_protect");
    exit(1);
}

/*
 * Allow our general code to reference the private data R/W.
 */
oldset = ukeyset_activate(keyset1RW, UKA_ADD_KEYS);
if (oldset == UKSET_INVALID) {
    printf("ukeyset_activate failed\n");
    exit(1);
}

/*
 * Set up a signal handler for SIGSEGV, to catch the deliberate
 * key violation in the untrusted code.
 */
sa.sa_sigaction = handler;
SIGINITSET(sa.sa_mask);
sa.sa_flags = SA_SIGINFO;
rc = sigaction(SIGSEGV, &sa, 0);
if (rc) {
    perror("sigaction");
    exit(1);
}

/*
 * Program's main processing loop.
 */
while (count < 2) {
    /*
     * When we need to run "untrusted" code, change access
     * to the private data to R/O by removing write access.
     */
    (void)ukeyset_activate(keyset1W, UKA_REMOVE_KEYS);

    /*
     * Call untrusted subroutine here. It can only read
     * the protected data passed to it.
     */
    count = untrusted(p1);

    /*
     * Restore our full access to private data.
     */
    (void)ukeyset_activate(keyset1W, UKA_ADD_KEYS);

    p1->some_data = count;
}
ukey_protect(p1, padded_protsize_1, UKEY_PUBLIC);
free(p1);
exit(0);
}

```

Large program support

This topic provides information about using the large and very large address-space models to accommodate programs requiring data areas that are larger than those provided by the default address-space model.

The large address-space model is available on AIX 4.3 and later. The very large address-space model is available on AIX 5.1 and later.

Note: This discussion applies only to 32-bit processes.

The virtual address space of a 32-bit process is divided into 16 256-megabyte areas (or *segments*), each addressed by a separate hardware register. The operating system refers to segment 2 (virtual addresses 0x20000000-0x2FFFFFFF) as the *process-private* segment. By default, this segment contains the user stack and data, including the heap. The process-private segment also contains the u-block of the process, which is used by the operating system and is not readable by an application.

Because a single segment is used for both user data and stack, their maximum aggregate size is slightly less than 256 MB. Certain programs, however, require large data areas (initialized or uninitialized), or they

need to allocate large amounts of memory with the **malloc** or **sbrk** subroutine. Programs can be built to use the large or very large address-space model, allowing them to use up to 2 GB of data.

It is possible to use either the large or very large address-space model with an existing program, by providing a non-zero *maxdata* value. The **maxdata** value is obtained either from the **LDR_CNTRL** environment variable or from a field in the executable file. Some programs have dependencies on the default address-space model, and they will break if they are run using the large address-space model.

Understanding the large address-space model

The large address-space model allows specified programs to use more than 256 MB of data. Other programs continue to use the default address-space model. To allow a program to use the large address-space model, specify a non-zero **maxdata** value. You can specify a non-zero **maxdata** value either by using the **ld** command when you're building the program, or by exporting the **LDR_CNTRL** environment variable before executing the program.

When a program using the large address-space model is executed, the operating system reserves as many 256 MB segments as needed to hold the amount of data specified by the **maxdata** value. Then, beginning with segment 3, the program's initialized data is read from the executable file into memory. The data read begins in segment 3, even if the **maxdata** value is smaller than 256 MB. With the large address-space model, a program can have a maximum of 8 segments or 2 GB or 3.25 GB of data respectively.

In the default address-space model, 12 segments are available for use by the **shmat** or **mmap** subroutines. When the large address-space model is used, the number of segments reserved for data reduces the number of segments available for the **shmat** and **mmap** subroutines. Because the maximum size of data is 2 GB, at least two segments are always available for the **shmat** and **mmap** subroutines.

The user stack remains in segment 2 when the large address-space model is used. As a result, the size of the stack is limited to slightly less than 256 MB. However, an application can relocate its user stack into a shared memory segment or into allocated memory.

While the size of initialized data in a program can be large, there is still a restriction on the size of text. In the executable file for a program, the size of the text section plus the size of the loader section must be less than 256 MB. This is required so that these sections will fit into a single, read-only segment (segment 1, the TEXT segment). You can use the **dump** command to examine section sizes.

Understanding the very large address-space model

The very large address-space model enables large data programs in much the same way as the large address-space model, although there are several differences between them. To allow a program to use the very large address-space model, you must specify a **maxdata** value and the dynamic segment allocation (*dsa*) property. Use either the **ld** command or the **LDR_CNTRL** environment variable to specify a **maxdata** value and the **DSA** option.

If a **maxdata** value is specified, the very large address-space model follows the large-address space model in that a program's data is read into memory starting with segment 3, and occupies as many segments as needed. The remaining data segments, however, are not reserved for the data area at execution time, but are obtained dynamically. Until a segment is needed for a program's data area, it can be used by the **shmat** or **mmap** subroutines. With the very large address-space model, a program can a maximum of 13 segments or 3.25 GB of data. Of these 13 segments, 12 segments or 3 GB, are available for **shmat** and **mmap** subroutine purposes.

When a process tries to expand its data area into a new segment, the operation succeeds as long as the segment is not being used by the **shmat** or **mmap** subroutines. A program can call the **shmdt** or **munmap** subroutine to stop using a segment so that the segment can be used for the data area. After a segment has been used for the data area, however, it can no longer be used for any other purpose, even if the size of the data area is reduced.

If a **maxdata** value is not specified (*maxdata* = 0) with the *dsa* property, a slight variation from the above behaviour is achieved. The process will have its data and stack in segment 2, similar to a regular process. The process will not have access to the global shared libraries, so all shared libraries used by the process will be loaded privately. The advantage to running this way is that a process will have all 13 segments (3.25 GB) available for use by the **shmat** and **mmap** subroutines.

To reduce the chances that the **shmat** or **mmap** subroutines will use a segment that could be used for the data area, the operating system uses a different rule for choosing an address to be returned (if a specific address is not requested). Normally, the **shmat** or **mmap** subroutines return an address in the lowest available segment. When the very large address-space model is used, these subroutines will return an address in the highest available segment. A request for a specific address will succeed, as long as the address is not in a segment that has already been used for the data area. This behaviour is followed for all process that specify the `dsa` property.

With the very large address-space model, a **maxdata** value of zero or a value of up to `0xD0000000` can be specified. If a **maxdata** value larger than `0xFFFFFFFF` is specified, a program will not use globally loaded shared libraries. Instead, all shared libraries will be loaded privately. This can affect program performance.

Enabling the large and very large address-space models

The large address space model is used if any non-zero value is specified for the **maxdata** value, and the dynamic segment allocation (`dsa`) property is not specified. The very large address-space model is used if any **maxdata** value is given and the `dsa` property is specified. Use the **ld** command with the **-bmaxdata** flag to specify a **maxdata** value and to set the `dsa` property.

Use the following command to link a program that will have the maximum 8 segments reserved for its data:

```
cc -bmaxdata:0x80000000 sample.o
```

To link a program with the very large-address space model enabled on the POWER processor-based platform, use the following command:

```
cc -bmaxdata:0xD0000000/dsa sample.o
```

To link a program with the very large-address space model enabled, use the following command:

```
cc -bmaxdata:0xD0000000/dsa sample.o
```

You can cause existing programs to use the large or very large address-space models by specifying the **maxdata** value with the **LDR_CNTRL** environment variable. For example, use the following command to run the **a.out** program with 8 segments reserved for the data area:

```
LDR_CNTRL=MAXDATA=0x80000000 a.out
```

The following command runs the **a.out** program using the very large address-space model, allowing the program's data size to use up to 8 segments for data:

```
LDR_CNTRL=MAXDATA=0x80000000@DSA a.out
```

You can also modify an existing program so that it will use the large or very large address-space model. To set the **maxdata** value in an existing 32-bit XCOFF program, **a.out**, to `0x80000000`, use the following command:

```
/usr/ccs/bin/ldedit -bmaxdata:0x80000000 a.out
```

If an existing 32-bit XCOFF program, **a.out**, with a **maxdata** value of `0x80000000` does not already have the **DSA** property, you can add the property with the following command:

```
/usr/ccs/bin/ldedit -bmaxdata:0x80000000/dsa a.out
```

You can use the **dump** command to examine the **maxdata** value, or to determine whether a program has the `dsa` property.

Some programs have dependencies on the default address-space model. These programs terminate if a non-zero **maxdata** value has been specified, either by modifying the executable file of the program or by setting the **LDR_CNTRL** environment variable.

Executing programs with large data areas

When you execute a program that uses the large address-space model, the operating system attempts to modify the soft limit on data size, if necessary, to increase it to match the **maxdata** value. If the **maxdata** value is *larger* than the current hard limit on data size, either the program will not execute if the environment variable **XPG_SUS_ENV** has the value set to ON, or the soft limit will be set to the current hard limit.

If the **maxdata** value is *smaller* than the size of the program's static data, the program will not execute.

After placing the program's initialized and uninitialized data in segments 3 and beyond, the break value is computed. The break value defines the end of the process's static data and the beginning of its dynamically allocatable data. Using the **malloc**, **brk** or **sbrk** subroutine, a process can move the break value to increase the size of the data area.

For example, if the **maxdata** value specified by a program is 0x68000000, then the maximum break value is in the middle of segment 9 (0x98000000). The **brk** subroutine extends the break value across segment boundaries, but the size of the data area cannot exceed the current soft data limit.

The **setrlimit** subroutine allows a process to set its soft data limit to any value that does not exceed the hard data limit. The maximum size of the data area, however, is limited to the original **maxdata** value, rounded up to a multiple of 256 MB.

The majority of subroutines are unaffected by large data programs. The **shmat** and **mmap** subroutines are the most affected, because they have fewer segments available for use. If a large data-address model program forks, the child process inherits the current data resource limits.

Special considerations

Programs with large data spaces require a large amount of paging space. For example, if a program with a 2-GB address space tries to access every page in its address space, the system must have 2 GB of paging space. The operating system terminates processes when paging space runs low. Programs with large data spaces are terminated first because they typically consume a large amount of paging space.

Debugging programs using the large data model is no different than debugging other programs. The **dbx** command can debug these large programs actively or from a core dump. A full core dump from a large-data program can be quite large. To avoid truncated core files, be sure the **coredump** resource limit is large enough, and make sure that there is enough free space in the file system where your program is running.

Some application programs might be written in such a way that they rely on characteristics of the default address space model. These programs might not work if they execute using the large or very large address-space model. Do not set the **LDR_CNTRL** environment variable when you run these programs.

Processes using the very large address-space model must make code changes to their programs in order to move the break value of the address-space in chunks larger than 2 GB. This is a limitation of the **sbrk** system call which takes a signed value as the parameter. As a workaround, a program can call **sbrk** more than one time to move the break value to the desired position.

Related information

[brk](#)

[exec](#)

[fork](#)

[malloc](#)

[setrlimit](#)

[dd Command](#)

[ld Command](#)

[XCOFF Object \(a.out\) File Format](#)

Programming on multiprocessor systems

On a uniprocessor system, threads execute one after another in a time-sliced manner. This contrasts with a multiprocessor system, where several threads execute at the same time, one on each available processor. Overall performance can be improved by running different process threads on different processors. However, an individual program cannot take advantage of multiprocessing, unless it has multiple threads.

Multiprocessing is not apparent to most users because it is handled completely by the operating system and the programs it runs. Users can bind their processes (force them to run on a certain processor); however, this is not required, nor recommended for ordinary use. Even for most programmers, taking advantage of multiprocessing simply amounts to using multiple threads. On the other hand, kernel programmers have to deal with several issues when porting or creating code for multiprocessor systems.

Identifying processors

Symmetric multiprocessor (SMP) machines have one or more CPU boards, each of which can accommodate two processors.

For example, a four-processor machine has two CPU boards, each having two processors. Commands, subroutines, or messages that refer to processors must use an identification scheme. Processors are identified by physical and logical numbers, and by Object Data Manager (ODM) processor names and location codes.

ODM processor names

ODM is a system used to identify various parts throughout a machine, including bus adapters, peripheral devices such as printers or terminals, disks, memory boards, and processor boards.

ODM assigns numbers to processor boards and processors in order, starting from 0 (zero), and creates names based on these numbers by adding a prefix `cpucard` or `proc`. Thus, the first processor board is called `cpucard0`, and the second processor on it is called `proc1`.

ODM location codes for processors consist of four 2-digit fields, in the form *AA-BB-CC-DD*, as follows:

Code	Description
-------------	--------------------

<i>AA</i>	Always 00. It indicates the main unit.
-----------	--

<i>BB</i>	Indicates the processor board number. It can be 0P, 0Q, 0R, or 0S, indicating respectively the first, second, third or fourth processor card.
-----------	---

<i>CC</i>	Always 00.
-----------	------------

<i>DD</i>	Indicates the processor position on the processor board. It can be 00 or 01.
-----------	--

Logical processor numbers

Processors can also be identified using logical numbers, which start with 0 (zero). Only enabled processors have a logical number.

The logical processor number 0 (zero) identifies the first physical processor in the enabled state; the logical processor number 1 (one) identifies the second enabled physical processor, and so on. Generally, all operating system commands and library subroutines use logical numbers to identify processors.

ODM processor states

If a processor functions correctly, it can be enabled or disabled using a software command. A processor is marked **faulty** if it has a detected hardware problem. ODM classifies processors using the following states:

State	Description
enabled	Processor works and can be used by AIX.

State	Description
disabled	Processor works, but cannot be used by AIX.
faulty	Processor does not work (a hardware fault was detected).

Controlling processor use

This section describes how to control the use of processors on the multiprocessor system.

On a multiprocessor system, the use of processors can be controlled in the following way:

- A user can force a process or kernel threads to run on a specific processor.

Binding processes and kernel threads

Users may force their processes to run on a given processor; this action is called *binding*. A system administrator may bind any process. From the command line, binding is controlled with the **bindprocessor** command.

The process itself is not bound, but rather its kernel threads are bound. After kernel threads are bound, they are always scheduled to run on the chosen processor, unless they are later unbound. When a new kernel thread is created, it has the same bind properties as its creator.

This situation applies to the initial thread in the new process created by the **fork** subroutine; the new thread inherits the bind properties of the thread that called the **fork** subroutine. When the **exec** subroutine is called, bind properties are left unchanged. After a process is bound to a processor, if no other binding or unbinding action is performed, all child processes will be bound to the same processor.

It is only possible to bind processes to enabled processors using logical processor numbers. To list available logical processor numbers, use the **bindprocessor -q** command. For a system with four enabled processors, this command produces output similar to the following:

```
The available processors are: 0 1 2 3
```

Binding may also be controlled within a program using the **bindprocessor** subroutine, which allows the programmer to bind a single kernel thread or all kernel threads in a process. The programmer can also unbind either a single kernel thread or all kernel threads in a process.

Using Dynamic Processor Deallocation

Starting with machine type 7044 model 270, the hardware of all systems with more than two processors can detect correctable errors, which are gathered by the firmware. These errors are not fatal and, as long as they remain rare occurrences, can be safely ignored. However, when a pattern of failures seems to be developing on a specific processor, this pattern may indicate that this component is likely to exhibit an unrecoverable failure in the near future. This prediction is made by the firmware based-on-failure rates and threshold analysis.

AIX implements continuous hardware surveillance and regularly polls the firmware for hardware errors. When the number of processor errors hits a threshold and the firmware recognizes the distinct probability that this system component will fail, the firmware returns an error report to AIX and logs the error in the system error log. In addition, on multiprocessor systems, depending on the type of failure, AIX attempts to stop using the untrustworthy processor and deallocate it. This feature is called *dynamic processor deallocation*.

At this point, the firmware flags the processor for persistent deallocation for subsequent reboots, until service personnel replace the processor.

Potential impact to applications

Processor deallocation is not apparent for the vast majority of applications, including drivers and kernel extensions. However, you can use AIX published interfaces to determine whether an application or kernel extension is running on a multiprocessor machine, find out how many processors there are, and bind threads to specific processors.

The bindprocessor interface for binding processes or threads to processors uses bind CPU numbers. The bind CPU numbers are in the range [0.. $N-1$] where N is the total number of CPUs. To avoid breaking applications or kernel extensions that assume no "holes" in the CPU numbering, AIX always makes it appear for applications as if the CPU is the "last" (highest numbered) bind CPU to be deallocated. For instance, on an 8-way SMP, the bind CPU numbers are [0..7]. If one processor is deallocated, the total number of available CPUs becomes 7, and they are numbered [0..6]. Externally, CPU 7 seems to have disappeared, regardless of which physical processor failed.

Note: In the rest of this description, the term *CPU* is used for the logical entity and the term *processor* for the physical entity.

Applications or kernel extensions using processes/threads binding could potentially be broken if AIX silently terminated their bound threads or forcibly moved them to another CPU when one of the processors needs to be deallocated. Dynamic processor deallocation provides programming interfaces so that those applications and kernel extensions can be notified that a processor deallocation is about to happen. When these applications and kernel extensions get this notification, they are responsible for moving their bound threads and associated resources (such as timer request blocks) away from the last bind CPU ID and adapt themselves to the new CPU configuration.

If, after notification of applications and kernel extensions, some of the threads are still bound to the last bind CPU ID, the deallocation is aborted. In this case, AIX logs the fact that the deallocation has been aborted in the error log and continues using the ailing processor. When the processor ultimately fails, it creates a total system failure. Thus, it is important for applications or kernel extensions that are binding threads to CPUs to get the notification of an impending processor deallocation, and to act on this notice.

Even in the rare cases where the deallocation cannot go through, dynamic processor deallocation still gives advanced warning to system administrators. By recording the error in the error log, it gives them a chance to schedule a maintenance operation on the system to replace the ailing component before a global system failure occurs.

Flow of events for processor deallocation

The typical flow of events for processor deallocation is as follows:

1. The firmware detects that a recoverable error threshold has been reached by one of the processors.
2. AIX logs the firmware error report in the system error log, and, when executing on a machine supporting processor deallocation, starts the deallocation process.
3. AIX notifies non-kernel processes and threads bound to the last bind CPU.
4. AIX waits for all the bound threads to move away from the last bind CPU. If threads remain bound, AIX eventually times out (after ten minutes) and aborts the deallocation. Otherwise, AIX invokes the previously registered High Availability Event Handlers (HAEHs). An HAEH may return an error that will abort the deallocation. Otherwise, AIX continues with the deallocation process and ultimately stops the failing processor.

In case of failure at any point of the deallocation, AIX logs the failure, indicating the reason why the deallocation was aborted. The system administrator can look at the error log, take corrective action (when possible) and restart the deallocation. For instance, if the deallocation was aborted because at least one application did not unbind its bound threads, the system administrator could stop the application(s), restart the deallocation (which should continue this time) and restart the application.

Programming interfaces dealing with individual processors

The following sections describe available programming interfaces:

Interfaces to determine the number of CPUs on a system

sysconf subroutine

The **sysconf** subroutine returns a number of processors using the following parameters:

- **_SC_NPROCESSORS_CONF**: Number of processors configured
- **_SC_NPROCESSORS_ONLN**: Number of processors online

The value returned by the **sysconf** subroutine for `_SC_NPROCESSORS_CONF` will remain constant between reboots. Uniprocessor (UP) machines are identified by a 1. Values greater than 1 indicate multiprocessor (MP) machines. The value returned for the `_SC_NPROCESSORS_ONLN` parameter will be the count of active CPUs and will be decremented every time a processor is deallocated.

The `_system_configuration.ncpus` field identifies the number of CPUs active on a machine. This field is analogous to the `_SC_NPROCESSOR_ONLN` parameter.

For code that must recognize how many processors were originally available at boot time, the `ncpus_cfg` field is added to the `_system_configuration` table, which remains constant between reboots.

The CPUs are identified by bind CPU IDs in the range `[0..(ncpus-1)]`. The processors also have a physical CPU number that depends on which CPU board they are on, in which order, and so on. The commands and subroutines dealing with CPU numbers always use bind CPU numbers. To ease the transition to varying numbers of CPUs, the bind CPU numbers are contiguous numbers in the range `[0..(ncpus-1)]`. The effect of this is that from a user point of view, when a processor deallocation takes place, it always looks like the highest-numbered ("last") bind CPU is disappearing, regardless of which physical processor failed.

Note: To avoid problems, use the `ncpus_cfg` variable to determine what the highest possible bind CPU number is for a particular system.

Interfaces to bind threads to a specific processor

The `bindprocessor` Command and the `bindprocessor` programming interface allow you to bind a thread or a process to a specific CPU, designated by its bind CPU number. Both interfaces will allow you to bind threads or processes only to active CPUs. Those programs that directly use the `bindprocessor` programming interface or are bound externally by a `bindprocessor` command must be able to handle the processor deallocation.

The primary problem seen by programs that bind to a processor when a CPU has been deallocated is that requests to bind to a deallocated processor will fail. Code that issues `bindprocessor` requests should always check the return value from those requests.

Interfaces for processor deallocation notification

The notification mechanism is different for user-mode applications having threads bound to the last bind CPU than it is for kernel extensions.

Notification in user mode

Each thread of a user mode application that is bound to the last bind CPU is sent the `SIGCPUFAIL` and `SIGRECONFIG` signals. These applications need to be modified to catch these signals and dispose of the threads bound to the last bind CPU (either by unbinding them or by binding them to a different CPU).

Notification in kernel mode

The drivers and kernel extensions that must be notified of an impending processor deallocation must register a High-Availability Event Handler (HAEH) routine with the kernel. This routine will be called when a processor deallocation is imminent. An interface is also provided to unregister the HAEH before the kernel extension is unconfigured or unloaded.

Registering a high-availability event handler

The kernel exports a new function to allow notification of the kernel extensions in case of events that affect the availability of the system.

The system call is:

```
int register_HA_handler(ha_handler_ext_t *)
```

For more information on this system call, see `register_HA_handler` in *Operating system and device management*.

The return value is equal to 0 in case of success. A non-zero value indicates a failure.

The system call argument is a pointer to a structure describing the kernel extension's HAEH. This structure is defined in a header file, named **sys/high_avail.h**, as follows:

```
typedef struct _ha_handler_ext_ {
    int (*_fun)();          /* Function to be invoked */
    long long _data;       /* Private data for (*_fun)() */
    char _name[sizeof(long long) + 1];
} ha_handler_ext_t;
```

The private **_data** field is provided for the use of the kernel extension if it is needed. Whatever value given in this field at the time of registration will be passed as a parameter to the registered function when the field is called due to a CPU predictive failure event.

The **_name** field is a null-terminated string with a maximum length of 8 characters (not including the null character terminator) which is used to uniquely identify the kernel extension with the kernel. This name must be unique among all the registered kernel extensions. This name is listed in the detailed data area of the CPU_DEALLOC_ABORTED error log entry if the kernel extension returns an error when the HAEH routine is called by the kernel.

Kernel extensions should register their HAEH only once.

Invocation of the high-availability event handler

The following parameters call the HAEH routine:

- The value of the **_data** field of the **ha_handler_ext_t** structure passed to **register_HA_handler**.
- A pointer to a **ha_event_t** structure defined in the **sys/high_avail.h** file as:

```
typedef struct {
    uint _magic;          /* High-availability related event */
    #define HA_CPU_FAIL 0x40505546 /* Identifies the kind of the event */
    union {
        struct {
            cpu_t dealloc_cpu; /* Predictive processor failure */
            ushort domain;    /* CPU bind ID of failing processor */
            ushort nodeid;    /* future extension */
            ushort reserved3; /* future extension */
            uint reserved[4];  /* future extension */
        } _cpu;
        /* ... */
        /* Additional kind of events -- */
        /* future extension */
    } _u;
} haeh_event_t;
```

The function returns one of the following codes, also defined in the **sys/high_avail.h** file:

```
#define HA_ACCEPTED 0 /* Positive acknowledgement */
#define HA_REFUSED -1 /* Negative acknowledgement */
```

If any of the registered extensions does not return HA_ACCEPTED, the deallocation is aborted. The HAEH routines are called in the process environment and do not need to be pinned.

If a kernel extension depends on the CPU configuration, its HAEH routine must react to the upcoming CPU deallocation. This reaction is highly application-dependent. To allow AIX to proceed with the deconfiguration, they must move the threads that are bound to the last bind CPU, if any. Also, if they have been using timers started from bound threads, those timers will be moved to another CPU as part of the CPU deallocation. If they have any dependency on these timers being delivered to a specific CPU, they must take action (such as stopping them) and restart their timer requests when the threads are bound to a new CPU, for instance.

Canceling the registration of a high-availability event handler

To keep the system coherent and prevent system crashes, the kernel extensions that register an HAEH must cancel the registration when they are unconfigured and are going to be unloaded. The interface is as follows:

```
int unregister_HA_handler(ha_handler_ext_t *)
```

This interface returns 0 in case of success. Any non-zero return value indicates an error.

Deallocating a processor in the test environment

To test any of the modifications made in applications or kernel extensions to support this processor deallocation, use the following command to trigger the deallocation of a CPU designated by its logical CPU number. The syntax is:

```
cpu_deallocate cpunum
```

where:

cpunum is a valid logical CPU number.

You must reboot the system to get the target processor back online. Hence, this command is provided for test purposes only and is *not* intended as a system administration tool.

Dynamic memory guarding

AIX systems are designed to be resilient in regards to memory errors. Memory error resilience is the result of both hardware and operating system-level recoveries.

There are multiple ways to categorize memory errors, but for the purposes of this discussion, memory errors are classified as recoverable and non-recoverable errors.

Recoverable errors result in data located in specific locations being retrievable, and unrecoverable errors result in a loss of data from the specific location in question. Unrecoverable errors are typically resolved by using hardware redundancy in the memory subsystem, or by masking the area in question from use during boot time of the operating system.

AIX supports resilience as a means of preventing recoverable memory errors from becoming unrecoverable errors through a technique known as Dynamic Memory Guarding. Dynamic Memory Guarding is based on support provided by the hardware. Hardware provides mechanisms for the detection of and recovery from errors (such as memory scrubbing and error correcting circuits (ECC)). Hardware can provide mechanisms for avoiding future unrecoverable errors as well, including redundant bit steering.

As a complement to these hardware mechanisms, the hardware can inform the operating system about errors best handled through Dynamic Memory Guarding. This is done by identifying areas of memory to be deallocated. The AIX operating system uses this information to mask off the memory area in question and to stop using it. The operating system will move any data currently contained in the memory area in error to another memory area, and then stop using the memory page that contains the memory location in error. This memory guarding is done by the operating system without any user intervention and is transparent to end users and applications.

Creating locking services

Some programmers may want to implement their own high-level locking services instead of using the standard locking services (mutexes) provided in the threads library.

For example, a database product may already use a set of internally defined services; it can be easier to adapt these locking services to a new system than to adapt all the internal modules that use these services.

For this reason, AIX provides atomic locking service primitives that can be used to build higher-level locking services. To create services that are multiprocessor-safe (like the standard mutex services), programmers must use the atomic locking services described in this section and not atomic operations services, such as the **`compare_and_swap`** subroutine.

Multiprocessor-safe locking services

Locking services are used to serialize access to resources that may be used concurrently. For example, locking services can be used for insertions in a linked list, which require several pointer updates. If the update sequence by one process is interrupted by a second process that tries to access the same list, an error can occur. A sequence of operations that should not be interrupted is called a *critical section*.

Locking services use a lock word to indicate the lock status: 0 (zero) can be used for free, and 1 (one) for busy. Therefore, a service to acquire a lock would do the following:

```
test the lock word
if the lock is free
    set the lock word to busy
return SUCCESS
...
```

Because this sequence of operations (read, test, set) is itself a critical section, special handling is required. On a uniprocessor system, disabling interrupts during the critical section prevents interruption by a context switch. But on a multiprocessor system, the hardware must provide a test-and-set primitive, usually with a special machine instruction. In addition, special processor-dependent synchronization instructions called *import and export fences* are used to temporarily block other reads or writes. They protect against concurrent access by several processors and against the read and write reordering performed by modern processors and are defined as follows:

Import fences

The *import fence* is a special machine instruction that delays until all previously issued instructions are complete. When used in conjunction with a lock, this prevents speculative execution of instructions until the lock is obtained.

Export fences

The *export fence* guarantees that the data being protected is visible to all other processors prior to the lock being released.

To mask this complexity and provide independence from these machine-dependent instructions, the following subroutines are defined:

_check_lock

Conditionally updates a single word variable atomically, issuing an *import fence* for multiprocessor systems. The **compare_and_swap** subroutine is similar, but it does not issue an import fence and, therefore, is not usable to implement a lock.

_clear_lock

Atomically writes a single word variable, issuing an *export fence* for multiprocessor systems.

Kernel programming

For complete details about kernel programming, see *Kernel Extensions and Device Support Programming Concepts*. This section highlights the major differences required for multiprocessor systems.

Serialization is often required when accessing certain critical resources. Locking services can be used to serialize thread access in the process environment, but they will not protect against an access occurring in the interrupt environment. New or ported code should use the **disable_lock** and **unlock_enable** kernel services, which use simple locks in addition to interrupt control, instead of the **i_disable** kernel service. These kernel services can also be used for uniprocessor systems, on which they simply use interrupt services without locking. For detailed information, see [Locking Kernel Services](#) in *Kernel Extensions and Device Support Programming Concepts*.

Device drivers by default run in a logical uniprocessor environment, in what is called *funneled* mode. Most well-written drivers for uniprocessor systems will work without modification in this mode, but must be carefully examined and modified to benefit from multiprocessing. Finally, kernel services for timers now have return values because they will not always succeed in a multiprocessor environment. Therefore, new or ported code must check these return values. For detailed information, see [Using Multiprocessor-Safe Timer Services](#) in *Kernel Extensions and Device Support Programming Concepts*.

Example of locking services

The multiprocessor-safe locking subroutines can be used to create custom high-level routines independent of the threads library. The example that follows shows partial implementations of subroutines similar to the **pthread_mutex_lock** and **pthread_mutex_unlock** subroutines in the threads library:

```
#include <sys/atomic_op.h>      /* for locking primitives */
#define SUCCESS 0
```



```

#define FAILURE          -1
#define LOCK_FREE       0
#define LOCK_TAKEN     1

typedef struct {
    atomic_p      lock; /* lock word */
    tid_t        owner; /* identifies the lock owner */
    ...          /* implementation dependent fields */
} my_mutex_t;

...

int my_mutex_lock(my_mutex_t *mutex)
{
    tid_t  self; /* caller's identifier */

    /*
     * Perform various checks:
     *   is mutex a valid pointer?
     *   has the mutex been initialized?
     */
    ...

    /* test that the caller does not have the mutex */
    self = thread_self();
    if (mutex->owner == self)
        return FAILURE;

    /*
     * Perform a test-and-set primitive in a loop.
     * In this implementation, yield the processor if failure.
     * Other solutions include: spin (continuously check);
     *   or yield after a fixed number of checks.
     */
    while (_check_lock(mutex->lock, LOCK_FREE, LOCK_TAKEN))
        yield();

    mutex->owner = self;
    return SUCCESS;
} /* end of my_mutex_lock */

int my_mutex_unlock(my_mutex_t *mutex)
{
    /*
     * Perform various checks:
     *   is mutex a valid pointer?
     *   has the mutex been initialized?
     */
    ...

    /* test that the caller owns the mutex */
    if (mutex->owner != thread_self())
        return FAILURE;

    _clear_lock(mutex->lock, LOCK_FREE);
    return SUCCESS;
} /* end of my_mutex_unlock */

```

Related information

[Locking Kernel Services](#)

[Using Multiprocessor-Safe Timer Services](#)

[bindprocessor](#)

[compare_and_swap](#)

[pthread_mutex_unlock](#)

[disable_lock](#)

[i_disable](#)

[unlock_enable](#)

[bindprocessor command](#)

ProbeVue dynamic tracing facility

You can use the ProbeVue dynamic tracing facility for both performance analysis and problem debugging. ProbeVue uses the Vue programming language to dynamically specify trace points and provide the actions to run at the specified trace points.

ProbeVue includes the following features:

- No pre-compiled trace hooks. ProbeVue works on unmodified kernel or user applications.
- Trace hooks do not have to be pre-compiled. They are compiled as part of your source programs. ProbeVue works on unmodified kernel or user applications.
- Trace hooks have no effect (do not exist) until they are dynamically enabled.
- Tracing actions (specified by the instrumentation code) to be issued at a trace hook are provided dynamically at the time the trace hook is enabled.
- Trace data captured as part of the tracing actions are available for viewing immediately and can be displayed as terminal output or saved to a file for later viewing.

Note: dbx and ProbeVue cannot debug a process simultaneously. Sometimes trying to debug a executable started by ProbeVue might result in dbx process waiting to attach to the process.

ProbeVue concepts

ProbeVue uses tracing actions or probe actions to capture information by writing the current values of global and context-specific information to a trace buffer.

The information captured is called trace data. The system usually provides facilities to read the data out of the trace buffer and make it available to the users of the system.

A *probe point* identifies a point during normal system activity that is capable of being probed. With dynamic tracing, probe points do not have any probes attached to them unless they are being traced. *Enabling* a probe is the operation of attaching a probe to a probe point and *disabling* a probe is the operation of removing a probe from a probe point. A probe is *triggered* or *fired* when an enabled probe point is reached during system activity and the tracing actions are performed.

ProbeVue supports the following two broad categories of probe points:

Probe location

A location in user or kernel code where some tracing action is to be performed. Enabled probes at a probe location fire when any thread running code reaches that location.

Probe event

An event of interest at whose occurrence some tracing action is performed. Probe events do not easily map to a specific code location. Enabled probes that indicate a probe event are configured to fire when that event occurs.

ProbeVue also distinguishes probe points by their type. A probe type identifies a set of probe points that share some common characteristics, for example, probe points that identify the entry and exit of system calls, or probe points that indicate updates to system statistics. Distinguishing probes by types gives structure to the wide variety of probe points.

ProbeVue command

The `probevue` command starts a dynamic tracing session or a ProbeVue session.

The **probevue** command starts a dynamic tracing session or a ProbeVue session. The **probevue** command takes a Vue script as input reading from a file or from the command line and activates a ProbeVue session. Any trace data that is captured by the ProbeVue session can be printed to the terminal or saved to a user-specified file as per options passed in the command line.

The ProbeVue session stays active until you type a Ctrl-C on the terminal or an exit action is issued from the Vue script.

Each invocation of the **probevue** command activates a separate dynamic tracing session. Multiple tracing sessions can be active at one time, but each session presents only the trace data that is captured in that session. Concurrent sessions are generally unaware of each other.

Running the **probevue** command is a privileged operation and non-root users needs privileges to initiate a dynamic tracing session.

Vue programming language

The Vue programming language provides your tracing specifications to ProbeVue.

The Vue programming language provides your tracing specifications to ProbeVue. A Vue script or Vue program is a program written in Vue. You can use a Vue script for the following purposes:

- Identify the probe points where a probe is to be dynamically enabled.
- Specify the conditions, if any, that must be satisfied to cause a probe to fire.
- Specify the actions issued, including the trace data that you want to capture.

In short, a Vue script tells ProbeVue where to trace, when to trace, and what to trace. Vue scripts should have a file suffix of ".e" to distinguish them from other file types.

Note that Vue is both a programming and a script language: it is a dedicated dynamic tracing language. Vue supports a subset of C and scripting syntax that is most beneficial for dynamic tracing purposes.

Elements of C

Vue supports a subset of C.

The following table describes how the support the ProbeVue compiler provides for the specified C keywords. All C keywords remain restricted in Vue. Use of any of these keywords as variable names or other symbols is not flagged as a syntax error, however, the behavior of such usage is undefined.

Note: The keywords in the second column can be present in type or structure definitions or function declarations. The Vue compiler ignores them. However, these keywords must not be applied when declaring Vue script variables.

Supported	Allowed in header files or in declaration section	Unsupported
char	auto	break
double	const	case
else	extern	continue
enum	register	default
float	static	do
if	typedef	for
int	volatile	goto
long		switch
return		while
short		
signed		
sizeof		
struct		
union		

Supported	Allowed in header files or in declaration section	Unsupported
unsigned		
void		

The following list describes the set of C features that Vue supports:

Statements

All C statements except for loop and some control flow statements.

Operators

All C unary, binary, and ternary operators except for the comma operator. Operator precedence and associativity follows C language rules.

Data types

Most legal C-89 defined variable types including all statements and keywords (struct, union, enum, typedef, and so on) for declaring types subject to the restrictions. This includes the types for kernel or application variables and parameters.

Note: Vue has its own rules on scope and storage classes.

Type conversions

Implicit type conversions as well as explicit conversions with type casts.

Subroutine

The syntax for calling subroutines and passing parameters to functions. There are, however, restrictions on what functions can be called.

Variable names

The naming conventions for variables follow C identifier rules. The complete variable specification can include colons if a variable class name is prefixed to the variable name.

Header files

Header files can be included for explicitly declaring the types of kernel global variables or the prototypes of functions in applications and the kernel. There are some restrictions on how the header files are included.

Punctuators

All C punctuators are supported and their rules are enforced. Thus, statements must be separated by the semicolon (;) character. All C white space rules are followed.

Literals

Representation of strings (using double quotes characters(")), character literals (using single quotes characters (')), octal and hexadecimal integers, and special characters like the `\n` and `\t` escape sequences.

Comments

C-style and C++-style comments. Comments can appear both inside and outside of a clause. Any line starting with a # character is ignored. Avoid using the character to indicate a comment line.

Differences from the C language

Vue has a different behavior for some C features. Some restrictions are imposed for maintaining efficiency or for ensuring that a Vue script can be issued safely inside the kernel and that it does not affect the probed process.

Loop statements

Loop statements are not accepted by a Vue script. This is a precaution to prevent any Vue probe from never completing.

Conditional flow statements

Only the "if-else" style control flow statement is accepted by a Vue script. Most conditional logic flow can be achieved through the proper use of "if" statements. Predicates are a more efficient way to do high-level conditional logic in a Vue script.

Return statement

The return statement is accepted by Vue to signal that the execution of the action block is to be ended immediately. However, the return statement does not take any expressions in Vue, as the Vue action block has no return values.

Subroutines

Vue scripts do not have access to the functions provided by the AIX system or general user libraries. There is no support for creating archives (libraries of functions) or user functions callable from within the probes. Instead, a special internal library is available for you that provide a set of functions generally useful for dynamic tracing programs.

Floating point

Floating-point variables are not accepted by any clause associated with a kernel probe point. You can only use floating-point variables in simple assignment statements and as parameters to Vue functions that print data. Vue language support for floating-point variables is restricted to its capture.

Variable modifications

External variables are not accepted by the left-hand side of an assignment statement, that is, they cannot be modified in a Vue script.

Header files

Vue does not support an explicit inclusion of the header file in the Vue script itself. Instead, the name of a header file to be included must be passed through command line arguments to the **probevue** command. Any C-preprocessor operators or directives in the header file are ignored. This might cause unexpected behavior. To avoid this, hand-code the header file or explicitly run the C-preprocessor directly on the set of relevant header files and generate a post-processed header file for inclusion. You can include the function prototypes and structure or union definitions in the Vue script itself if they are inserted at the very beginning before starting any of the probe clauses.

C-preprocessor

C-preprocessor operators, macro definitions, line or pragma directives and pre-defined macro names are ignored.

Pointer operations

Vue does not accept pointers to script variables. For example, the address of a script variable cannot be taken. However, the address of a kernel variable can be taken and assigned to a Vue pointer variable and pointer operations supported using the pointer variable.

Miscellaneous

- Trigraphs are not accepted.
- The comma operator is not accepted.
- Declaration statements cannot contain any initialization.

Vue scripts

Unlike procedures in procedural languages, an action block in Vue does not have an output or return value.

Unlike procedures in procedural languages, an action block in Vue does not have an output or return value. It also does not have inherent support for a set of input parameters. On the other hand, the context data at the point where a probe is entered can be accessed within the action block. For example, the parameters passed to a function can be referenced within the action block of a Vue clause if the probe point is at the function's entry point.

Predicates

You must use predicates when execution of clauses at probe points must be performed conditionally. The predicate section is identified by the presence of the **when** keyword immediately after the probe specification section. The predicate itself consists of regular C-style conditional expressions with the enclosing parentheses.

A predicate has the following format:

```
when ( <condition> )
```

For example:

```
when ( __pid == 1678 )
```

Vue script example

The following script is an example of a Vue script:

```
/* Global variables are auto-initialized to zero */ [1]

int count; [2]
/*
 * File: count.e
 *
 * Count number of times the read or write system call is entered
 * by process with Id 400
 */

@@BEGIN
{
    printf("Start probing\n");
}

@@syscall:*:read:entry, @@syscall:*:write:entry [3]
    when ( __pid == 400) [4]

{ [5]

    count++;
    /* Print a message for every 20 system calls */
    if (count % 20 == 0)
        printf("Total read/writes so far: %d\n", count);
    /* Exit when we exceed 100 system calls */
    if (count > 100)
        exit();
} [6]

/* print some statistics at exit */
@@END
{
    printf("Terminating probe after %d system calls.\n", count);
}
```

The following superscripts used in the above example identify the different elements of a Vue script:

1. Comments
2. (Optional) Declaration section
3. Probe specification
4. (Optional) Predicate
5. Start of action block
6. End of action block

You can start this simple script by issuing the following command. Note that this example displays some sample output.

```
# probevue count.e
Total read/writes so far: 20
Total read/writes so far: 40
Total read/writes so far: 60
...
...
```

Running the **probevue** command requires privileges. To issue the above command successfully, you must have logged in as the superuser or been granted privileges to probe system calls made by any process in the system.

Probe point specification

A probe point specification consists of one or more probe point tuples.

Each probe point tuple identifies a code location whose execution or an event whose occurrence must trigger the probe actions. Multiple probe points can be associated with the same set of probe actions and the predicate, if any, by providing a comma-separated list of probe tuples at the top of the Vue clause.

The following are some of the probe types that are supported:

- User function entry probes (or **uft** probes)
- System call entry or exit probes (or **syscall** probes)
- Probes that fire at specific time intervals (or **interval** probes)

For a complete list of probe types that are supported, refer the probe manager section.

The probe point tuple is an ordered list of fields that are separated by colons that uniquely identify a probe point. It has the following general format, although the location field is usually present only if the probe point is a probe location.

```
@@ <probetype>:  
<one or more probetype-specific fields separated by colons>:<location>
```

The probe manager defines the acceptable values for the probe-type specific fields in the probe point tuple and the length of the probe tuple. However, the following general rules are followed by all probe managers when defining probe point tuples:

- Each probe point tuple is at least a 3-tuple, that is, it has a minimum of 3 fields.
- The first field always identifies the probe type and thus its probe manager.
- For probe managers that support process-specific tracing, the second field must be a process ID.
- For probe managers that support function entry or exit probes, the location field (the last field) must use the **entry** or **exit** keyword.
- Fields are separated by the colon (:) symbol.
- An asterisk or the "*" symbol for a field in the probe point tuple indicates that it matches any possible value for that field. For example, the **syscall** probe manager allows for system calls of a specific process or for all processes to be probed. In the first case, the second field must be the process ID of the process to be probed. In the latter case, when all processes are to be probed, the second field must be the "*" symbol. A second use of the asterisk symbol for a field is to allow for finer-grained probes in the future while maintaining binary compatibility with existing scripts. For example, the **uft** probe manager currently requires the third field to be an asterisk. In the future, it can support a module name as the third field to limit probes to only functions defined in that module.
- Maximum length of the probe specification is 1023 characters.

For example:

@@uft:34568*:foo:entry

Probe at entry into any function called **foo** in process with ID = 34568. The asterisk in the third field indicates that the **foo** function is to be probed if it exists in any module of the process.

@@syscall*:read:exit

Probe at exit of the **read** system call. The asterisk indicates that the **read** system call for all processes are to be probed.

@@interval*:clock:500

Probe to fire every 500 milliseconds (wall clock time). The asterisk is a placeholder for supporting finer probe points in the future.

The process ID for a process is often not known at the time the Vue script is written. Vue provides a simple method to avoid having to hard-code a process ID in the second field of the probe specification or anywhere in a Vue script (for example, in the predicate section).

A single Vue script can contain probe points from multiple processes in user space and in the kernel. Any trace output generated is always displayed in time order.

In addition to regular probe points defined by probe managers, Vue supports two special probe points. Each Vue script can contain a `@@BEGIN` probe point to indicate any action that needs to be issued before enabling any probes and an `@@END` probe point to indicate any action to be issued after the tracing has been terminated.

Action block

The action block identifies the trace actions to be performed when the associated probe point is triggered. Supported actions are not restricted to capturing and formatting of trace data but the full power of the Vue language can be employed.

An action block in Vue is similar to a procedure in procedural languages. It consists of a sequence of statements that are issued in order. The flow of execution is essentially sequential. The only exceptions are that conditional execution is possible using the "if-else" statement and control can be returned from the action block using the "return" statement. Vue also supports an **exit** function that terminates the entire script and ends the tracing session. There are no constructs for looping in Vue and C language statements, so "for", "do", "goto", and so on, are not supported.

Unlike procedures in procedural languages, an action block in Vue does not have an output or return value. It also does not have inherent support for a set of input parameters. On the other hand, the context data at the point where a probe is entered can be accessed within the action block. For example, the parameters passed to a function can be referenced within the action block of a Vue clause if the probe point is at the function's entry point.

Predicates

You must use predicates when execution of clauses at probe points must be performed conditionally. The predicate section is identified by the presence of the **when** keyword immediately after the probe specification section. The predicate itself consists of regular C-style conditional expressions with the enclosing parentheses.

A predicate has the following format:

```
when ( <condition> )
```

For example:

```
when ( __pid == 1678 )
```

ProveVue variables

The Vue language supports most of the traditional C data types, namely those recognized by the C-89 specification. In addition, Vue includes some extensions to make powerful dynamic tracing programs be written easily.

Vue supports variables with three different scope rules:

- Variables that are local to one action block only
- Variables that have global scope
- Variables that have thread-local scope

In addition, Vue can access variables with external scope like global variables in the kernel or user data in an application being probed.

In general, variables need to be declared before their first use in the script, although Vue also supports a very limited form of implicit type recognition. Variable declaration statements inside an action block must appear before any of the executable statements. They cannot be inside nested blocks like within an if statement. In some cases, you can declare variables outside any of the action blocks, but in this case, all such declarations must appear before the first action block.

Variable classes

Vue supports several classes of variables with varying rules on scope, on how they are initialized, on whether they can be updated or not and on how their types are determined. As in the C language, any declaration statement for a variable must textually precede its first use in the script.

Vue provides special type qualifiers that are added to the declaration statement to indicate the class of the variables being declared. For example, the **__global** keyword is a class qualifier that you can include in the declaration statement to specify that the variables being declared have "global" class.

In the following example, both **foo** and **bar** are declared to be variables of global class:

```
__global int foo, bar;
```

Vue also supports implicit recognition of the type of a variable based on its first usage in the script. In this case, there is no declaration statement, but the class of the variable can still be provided by directly attaching a class qualifier to the variable as follows on its first textual reference in the script:

```
global:count = 5; /* First reference to variable count in the script */
```

In the preceding example, the **global:** keyword is a qualifier that specifies the **count** variable to be a variable of global class. This variable will also implicitly be assigned the int type because the first reference to it is an assignment expression whose right hand side is an integer constant.

Note: You need to use the **__global** keyword when specifying the class qualifier with the declaration statement, but the **global:** keyword when defining it at the first use of the variable in the script. The syntax rules are similar for the other class qualifiers supported by Vue.

Automatic class variables

An automatic variable is clause-specific and is similar to an automatic or stack variable in C. It has scope only within the action block portion of the clause where it is defined or used and is recreated for each invocation of the action block. Automatic variables are always undefined at the start of an action block and must be initialized through an assignment statement before you can use them in an expression or in any other executable statement.

An automatic variable is identified by using the **auto:** prefix, for example `auto:lticks` indicates an automatic variable. You can also declare automatic variables using the **__auto** declaration statement in which case the **auto:** prefix can be omitted.

You cannot use automatic class variables in the predicate section of a Vue clause.

The following script is an example of the **__auto** declaration statement:

```
__auto int i; /* Explicit declaration */  
auto:j = 0; /* Implicit declaration */
```

Thread-local class variables

A thread-local variable is instantiated per traced thread the first time it issues an action block that assigns a value to the variable. Once created, the thread-local variable exists as long as the Vue script is active and the traced thread does not exit. The value of the thread-local variable is thread-specific and retained across executions of any of the clauses of the same program. In other words, variables of this class are visible everywhere within the Vue script. However, each thread that issues the Vue script obtains its own copy of these variables and the variables in each such copy are accessible and modifiable anywhere within the script only by the thread that instantiated them.

A thread-local variable is distinguished by using the **thread:** prefix. For example, `thread:count` indicates a thread-local variable. You can also declare thread-local variables using the **__thread** declaration statement in which case the **thread:** prefix can be omitted with the following one exception.

You can use a thread-local variable in the predicate section of a Vue clause even before it is instantiated. Predicates with un-instantiated thread-local variables are always evaluated to a value of FALSE. When used in the predicate section, the **thread:** prefix must always be included to identify it as a thread-local variable.

The following script is an example of the `__thread` declaration statement:

```
__thread int i;      /* Explicit declaration */
__thread j = 0;     /* Implicit declaration */
```

Note: Although you can declare thread-locals inside the `@@BEGIN` and `@@END` probes, any other references to them in these special probes can produce undefined behavior. A declaration statement by itself does not cause the thread-local variable to be instantiated.

Global class variables

Variables of global class have global scope and are visible everywhere within a Vue script. You can use a global variable in one or more clauses of a Vue script. They can also be declared at the beginning textually before the first clause for clarity. Global variables are initialized to zero or NULL as appropriate.

All variables in a Vue script are by default assigned global class, unless an explicit non-global class specifier is prefixed to the declaration. You can also explicitly declare global variables by using the `__global` class specifier when declaring a variable. List variables are, by definition, always created as variables of global class.

Reads and updates of global variables are not serialized unless they are of the list type. There are no guarantees on data races when probes are issued simultaneously. Global variables, which are not of the list type, are useful for collecting profiling and other statistics.

You can use global variables in the predicate section of a Vue clause.

The following scripts are examples for initializing and using global variables:

```
int wcount;          /* Global variable declared before first clause */

@@BEGIN
{
    int f_count;      /* Global variable declared inside @@BEGIN */
    __global int z_count; /* Global variable declared with __global prefix */

    f_count = 12;
}

@@syscall:*:read:entry
when (z_count == 0)
{
    int m_count;      /* Global variable declared inside a probe */
    m_count += f_count; /* f_count already declared in earlier probe */
    printf("m_count = %d\n", m_count);
    if (wcount == 1)
        exit();
}

@@syscall:*:write:entry
{
    m_count++;        /* m_count already declared in earlier probe */
}

@@syscall:*:write:exit
{
    wcount = 1;       /* w_count declared globally */
}
```

Kernel global class variables

In ProbeVue, a privileged user can access kernel global variables inside the action block of any Vue clause, even for probe points in user space like the `uft` probe points. Before using or referring to the kernel variable in the Vue script, you must explicitly declare it using the `__kernel` declaration statement. Only variables exported by the kernel, that is, only those are present in the export list of `/unix` are accessible.

Kernel variables cannot appear in the predicate section of a clause. Kernel variables are always treated as read-only variables in a Vue script. Any attempt to write to a kernel variable either causes a syntax error or fail later with a script abort message.

For an example of how kernel variables can be declared and used in a Vue script.

Access only pinned kernel variables. If the page containing the kernel variable is not in memory (has been paged out), ProbeVue returns a value of zero for that variable.

You can access integral type kernel variables and kernel variables that are structures or unions and even pointers. Further, you can also refer to the member names of kernel structures and unions in a Vue script. Kernel arrays can also be accessed but there is no support for copying kernel character data into a ProbeVue string.

Useful kernel variables

The following table lists a few examples of useful kernel variables that can be accessed from within a Vue script. Be careful when using them in a Vue script because although not likely, the names of these variables or their meanings can change across different releases of AIX. All these kernel variables are pinned in memory and exported from the kernel.

Kernel variable	Description	Associated header files
struct system_configuration _system_configuration	System configuration structure.	sys/systemcfg.h
struct var v	Base kernel tunable (and other) parameters.	sys/var.h
struct timestruc_t tod	Memory-mapped time of day clock. Seconds and nanoseconds since Epoch.	sys/time.h
cpu_t high_cpuid	Highest logical CPU ID ever online.	sys/encap.h
struct vminfo vmminfo	Data structure that contains the information shown by the vmstat command.	sys/vminfo.h
time_t lbolt	Number of ticks since last boot.	sys/time.h
char spurr_version	Identifies if current system supports the SPURR register 0=No SPURR, 1=CPUs have SPURR.	sys/sysinfo.h
struct utsname utsname	System name structure that includes the operating system name, node name, release level, and so on.	sys/utsname.h

Entry class variables

Clauses associated with a probe point that is at the entry location point of a system call or a user function can access the arguments passed to the system call or function being probed.

Probes at entry location points are supported by the system call and the user function tracing probe managers. For example, the read system call takes three arguments: a file descriptor ID, a pointer to a user buffer, and a value for the number of bytes of data to be read. The values of these three arguments can be accessed if the probe specification is **@@syscall:*:read:entry**, which specifies a probe at the read system call entry point.

Parameters to functions are referenced using the special built-in entry class variable names **__arg1**, **__arg2**, **__arg3**, ... up to the number of arguments passed to the function. For example, in the clause

associated with the read system call entry point, `__arg1` refers to the value of file descriptor id parameter, `__arg2` refers to the value of the buffer pointer parameter, and `__arg3` to the size of the data to be read.

Note: When one or more probe point tuples are specified, then `__arg <x>` variables are not allowed in the Action Block and will result in error as shown in the example below.

```
@@syscall:*:read:entry,@@syscall:*:write:entry
{
    char *argument;
    argument=__arg2;  -> Not Allowed.
}
```

ProbeVue would exit with the following error message: *arg builtin cannot be used. No defined function.*

Use of entry class variables in a Vue clause is legal only if the C-style declaration of the function being probed, specifically the data type of the parameters being passed to the function, are also provided in the Vue script. This must appear textually before the first Vue clause that references the entry clause. Place the declaration textually before any Vue clause at the top of the Vue script.

The following script is an example of using entry class variables:

```
int read(int fd, char *buf, unsigned long size);
@@syscall:*:read:entry
{
    printf("Number of bytes to read = %d\n", __arg3);
}
```

Note: In the preceding example, the definition of the read system call function specified in the script does not exactly match what is given in the `/usr/include/unistd.h` file, but it works just as well.

A second requirement is that the probe specification associated with the clause identify a unique probe point. Entry class variables cannot be used in a Vue clause that has multiple probe points specified in the probe specification irrespective of whether the functions being probed are the same or have similar function prototypes. The following script is an illegal script and will cause the ProbeVue compiler to fail with a syntax error since the probe specification includes two probe points:

```
int read(int fd, char *buf, unsigned long size);
int write(int fd, char *buf, unsigned long size);
@@syscall:*:read:entry, @@syscall:*:write:entry
{
    /* Cannot use __arg3 in here, as this clause has multiple probe
     * points associated with it. This script will fail with a
     * syntax error in the compilation phase of the probevue command.
     */
    printf("Number of bytes to read/write = %d\n", __arg3);
}
```

The following modified script can work:

```
int read(int fd, char *buf, unsigned long size);
int write(int fd, char *buf, unsigned long size);
@@syscall:*:read:entry
{
    printf("Number of bytes to read = %d\n", __arg3);
}
@@syscall:*:write:entry
{
    printf("Number of bytes to write = %d\n", __arg3);
}
```

Exit class variables

Clauses associated with a probe point that is at the exit location points of a system call or user function can access the return value of the system call or user function.

There is only one exit class variable that is defined by the Vue language. This is the return value from a function or a system call, which can be accessed by using the special built-in variable name `__rv`.

Probes at exit location points are supported by the system call probe manager. For example, the read system call returns the actual number of bytes read or an error return code of -1. This returned value can be accessed at the `@@syscall:*:read:exit` probe point, which identifies all exit points from the read system call.

Similar to entry class variables, the use of exit class variables in a Vue clause is legal only if the probe specification associated with the clause identifies a unique probe point. Thus, `__rv` cannot be used in a Vue clause that has multiple probe points specified in the probe specification. Furthermore, the C-style declaration of the function being probed, specifically the data type of the return value, must be explicitly provided in the Vue script. In fact, it is an error to specify a function declaration without providing its return type.

You can use exit class variables in the predicate section of a clause.

The following script is an illegal script and will cause the ProbeVue compiler to fail with a syntax error since the return type of the `read` function is not specified:

```
/* Bad example. */
int read(int fd, char *buf, unsigned long size);
@@syscall:*:read:exit
when (__rv > 0)
{
    /* Entered on read success: return value = # of bytes read */
    printf("Number of bytes read = %d\n", __rv);
}
```

The following modified script can work:

```
/* Good example. */
int read(int fd, char *buf, unsigned long size);
@@syscall:*:read:exit
when (__rv > 0)
{
    /* Entered on read success: return value = # of bytes read */
    printf("Number of bytes read = %d\n", __rv);
}
```

Built-in class variables

In addition to the special built-in variables, `__arg1` through `__arg32` and `__rv`, Vue also defines a set of general-purpose built-in variables. These general-purpose built-in variables are discussed in more detail in this section and some probe manager specific built-in variables are discussed in their respective probe manager section. Built-in class variables are functions, but are treated as variables by ProbeVue. Therefore, you can use these built-in variables in the predicate section of a Vue clause.

The following built-in variables are supported in Vue:

- __tid**
Thread ID of traced thread.
- __pid**
Process ID of traced thread.
- __ppid**
Parent process ID of traced thread.
- __pgid**
Process group ID of traced thread.
- __pname**
Process name of traced thread.

__uid, __euid

Real and effective user ID of traced thread.

__trcid

Process ID of tracing process (that is, of the **probevue** command)

__errno

Current errno value for the traced thread.

__kernelmode

Current executable mode: is either 1 (in kernel mode) or 0 (in user mode).

__r3, ..., __r10

General purpose register values (for function parameters or return values).

__curthread

Current thread.

__curproc

Current process.

__ublock

User area of the current process.

__mst

Built-in variable to access the hardware register content of the current thread's machine state save area (MST).

__stat

Built-in variable to provides access to the system statistics for various AIX® kernel components.

The following script is an example using built-in variables:

```
@syscall:*:read:entry
{
    printf("Thread ID:%d, Process ID:%d, Parent Process ID:%d\n",
        __tid, __pid, __ppid);
    printf("Process Group ID: %d\n", __pgid);
    printf("Process name = %s\n", __pname);

    printf("Real UID=%d, Effective UID=%d\n", __uid, __euid);S
    printf("probevue command process ID = %d\n", __trcid);

    printf("Errno = %d\n", __errno);
    printf("Mode = %s\n", __kernelmode == 1 ? "kernel" : "user");

    printf("Current values of GPRs: r3=0x%016llx, r4=0x%016llx, r5=0x%016llx\n",
        __r3, __r4, __r5);
    printf("                r6=0x%016llx, r7=0x%016llx, r8=0x%016llx\n",
        __r6, __r7, __r8);
    printf("                r9=0x%016llx, r10=0x%016llx\n",
        __r9, __r10);
}
```

__curthread built-in variable

__curthread is a special built-in using which the user can access some of the thread related information for the current thread. The information can be accessed using the **->** operator on the **__curthread** built-in. This built-in cannot be used in **systrace**, **BEGIN** and **END** probes. Also it can be used in interval probes only if **PID** is mentioned. This built-in will basically provide functionality similar to **getthrds/getthrds64** but only limited to the current thread. The data that can be accessed are

tid

Thread ID

threadstate

State of the thread

pid

Process ID

policy

Scheduling Policy

pri

Priority

cpuusage

CPU Usage

cpuid

Processor to which the current thread is bound to

sigmask

Signal blocked on the thread

lockcount

Number of kernel lock taken by the thread

ptid

The pthread identifier of this thread (0 if it is a kernel thread, 1 if it is a single-threaded application)

homecpu

Home CPU of a thread.

homesrad

Home srad of a thread

Usage Example

```
Tid of the current thread can be accessed using __curthread->tid.
```

__curproc built-in variable

__curproc is a special built-in using which the user can access some of the process related information for the current process. The information can be accessed using the **->** operator on the **__curproc** built-in. This built-in cannot be used in **systrace**, **BEGIN** and **END** probes. Also, it can be used in interval probes only if **PID** is mentioned. This built-in will basically provide functionality similar to **getproc** but only limited to the current process. The data that can be accessed are

pid

Process ID.

ppid

Parent Process ID

pgid

Process Group ID

uid

Real user ID

suid

Saved user ID

pri

Priority

nice

Nice value

cpu

Processor usage

adspace

Process Address Space

majflt

I/O Page Fault

minflt

Non I/O Page Fault

size

Size of image in pages

sigpend

Signals pending on the process

sigignore

Signals ignored by the process

sigcatch

Signals being caught by the process

forktime

Creation time of the process

threadcount

No of threads in the process

cwd

Current working directory. If a free page fault context is not available or the per-CPU computation stack size is less than 96 KB, or in a probe where a page fault is not allowed (for example, interval probe), then this built-in returns a null string

Usage Example

```
Parent process id of the current process can be accessed using __curproc->ppid.
```

__ublock built-in variable

__ublock is a special built-in using which the user can access some of the process related information for the current process. This built-in cannot be used in **systrace**, **BEGIN** and **END** probes. Also it can be used in interval probes only if **PID** is mentioned. The information can be accessed using the **->** operator on the **__ublock** built-in. The data that can be accessed are

text

Start of text

tsize

Text size (bytes)

data

Start of Data

sdata

Current Data Size (bytes)

mdata

Maximum data size (bytes)

stack

Start of stack

stkmax

Stack Max (bytes)

euid

Effective user id

uid

Real user id

egid

Effective group id

gid

Real group id

utime

Process User resource usage time in seconds

stime

Process System resource usage time in seconds

maxfd

Max fd value in user

is64u

Set to 1, if in context of a 64-bit process

Usage Example

```
Start of the text for the current process can be accessed using __ublock->text.
```

__mst built-in variable

__mst is a special built-in variable in which you can access hardware register content of the current thread. This built-in variable cannot be used in the **systrace**, **BEGIN**, and **END** probes. Also, this built-in variable can be used in the interval probes only if **PID** is mentioned. The information can be accessed using the **->** operator on the **__ublock** built-in. The registers that can be accessed follow:

r1-r10

General purpose register r1 to r10

r14-r31

General purpose register r14 to r31

iar

Instruction address register

lr

Link register

isisr

Set if in interrupt or exception context.

Usage Example

To access the **lr** value in a probe, use the following command:

```
__mst->lr
```

__stat built-in variable

This built-in variable provides access to the system statistics for various AIX kernel components by using the Vue script. The system statistics are provided as running counters that can be accessed from any probe point of any ProbeVue script. A new probe point is not added to support the system statistics. To access system statistics, you must have `aix.ras.probevue.trace` kernel tracing privileges.

System statistics are useful for the following reasons:

- Statistics can be accessed without turning on a system trace or a component trace.
- Only required fields can be displayed. This is not possible by using the current statistics commands. You can avoid copying a large amount of data by accessing fields directly from the kernel structures.
- The system statistics is now available in the Vue scripts to easily complete the arithmetic logical operation. For example, you can use ProbeVue to add the number of operations across two disks.

The Vue built-in **__stat** variable allows statistics to be retrieved both at the global level and the individual component level. Data is provided as running counters and Vue scripts require access to the counters periodically. You can save the values that you need, and compute delta values to get the required results. You can access system statistics for the following reasons:

- To write a simple statistic tool by using the Vue script without calling the C/C++ API to print out the delta value for the counters every second or at an interval that is requested by the user.
 - To monitor counter value in case if the delta value exceeds threshold. When the value exceeds the threshold, a script logs the message.

- To use the Vue script to log a message when the actual value of the counter, which is not a delta value, exceeds the threshold value (for example, the maximum service time of a disk).

There are different modes that are provided by the built in **__stat** variable to access statistics from the source. The source of the system statistic is obtained from different system components. The statistics access modes follow.

Synchronous mode access

ProbeVue provides a direct access to the system statistics when the Vue script is run. This type of direct access is a live data access. By default, ProbeVue selects a live data access, if it is available. Every statistic might not be accessed in this mode because components do not provide the direct access or direct access is not possible in the current thread.

Asynchronous or Cached mode

Data is periodically gathered from source and cached in the ProbeVue. The Vue script uses a cache to access the data. A cache refresh interval is tuned at the single session or at all session level. Every source component provides a method to access statistics in asynchronous mode. In such cases, you can provide access to all statistics in cached mode by using the **fetch_stats_async_only** tunable. Cached mode is used when you cannot live data or you want to monitor the system where live data is not required.

Examples

1. The following example prints the number of I/O transfers for a disk named hdisk9 every second:

```
@@interval:*:clock:1000
{
  printf("Number of transfers = %lld\n", __stat.io.disk.hdisk9->transfers);
}
```

2. The following example displays the units of service time is microseconds:

```
@@syscall:*:read:exit
{
  rdservtime = __stat.io.disk.hdisk10->rd_service_time;
  printf("rdservtime=%lld microseconds\n", rdservtime);
}
```

__stat built-in syntax

The general syntax for using the **__stat** built-in expressions is as follows: **__stat.<level1_keyword>[.<level2_keyword>.....][.<inst1_keyword>.....]-><fieldname>**

Pre-defined levels and instances for the Storage IO Statistics are shown in the following table:

Built-in	Level 0	Level 1	Instance 0	Instance 1	Fields Names
__stat	io	disk	hdisk[0...n]		• See Table 4 on page 201
			hdisk[0...n]	path[0...n]	• See Table 5 on page 202
__stat	io	adapter	vscsi[0....n]		• See Table 6 on page 202 • See Table 7 on page 203
__stat	io	adapter	fcs[0...n]		• See Table 8 on page 204

Pre-defined levels and instances for the Network Statistics are shown in the following table:

Built-in	Level 0	Level 1	Level 2	Instance 0	Fields Names
__stat	net	adapter		ent[0...n]	• See Table 9 on page 206
__stat	net	interface		en[0....n]	• See Table 10 on page 209

Table 3. Predefined level and instances (continued)

Built-in	Level 0	Level 1	Level 2	Instance 0	Fields Names
__stat	net	protocol	ip		• See Table 11 on page 211
__stat	net	protocol	ipv6		• See Table 12 on page 214
__stat	net	protocol	tcp		• See Table 13 on page 216
__stat	net	protocol	udp		• See Table 14 on page 221
__stat	net	protocol	icmp		• See Table 15 on page 222
__stat	net	protocol	icmpv6		• See Table 16 on page 223
__stat	net	protocol	igmp		• See Table 17 on page 226
__stat	net	protocol	arp		• See Table 18 on page 227

SCSI Disk I/O statistics

The following table displays the supported field names for the Small Computer System Interface (SCSI) disk I/O statistics. These fields are specific to the disk instances. They can be accessed as `__stat.io.disk.<hdisk0...n>->fieldname`. The following statistics can be accessed in synchronous mode and asynchronous mode.

Table 4. SCSI disk I/O statistics

Vue script field name	Data type	Description
name	String[32]	Disk name
block_size	unsigned long long	Disk block size in bytes
transfers	unsigned long long	Number of transfers to or from disk
rd_block_count	unsigned long long	Number of disk blocks read
rd_service_time	unsigned long long	Total read service name or receive service time in microseconds.
rd_min_service_time	unsigned long long	Minimum read service time in microseconds.
rd_max_service_time	unsigned long long	Maximum read service time in microseconds.
rd_timeouts	unsigned long long	Number of read timeouts
rd_failures	unsigned long long	Number of read failures
wr_block_count	unsigned long long	Number of blocks written
wr_service_time	unsigned long long	Total write service time in microseconds.
wr_min_service_time	unsigned long long	Minimum write service time in microseconds.
wr_max_service_time	unsigned long long	Maximum write service time in microseconds.
wr_timeouts	unsigned long long	Number of write timeouts
wr_failures	unsigned long long	Number of write failures

Table 4. SCSI disk I/O statistics (continued)

Vue script field name	Data type	Description
wait_queue_depth	unsigned long long	Driver wait queue depth
accum_wait_queue_time	unsigned long long	Accumulated wait queuing time in microseconds
min_wait_queue_time	unsigned long long	Minimum wait queue time in microseconds.
max_wait_queue_time	unsigned long long	Maximum wait queue time in microseconds.
num_queue_full	unsigned long long	Number of in-flight queue full count

SCSI disk path I/O statistics

The following table displays the supported field names for the Small Computer System Interface (SCSI) disk path I/O statistics. These fields are specific to the disk and path instances. They can be accessed as `__stat.io.disk.<hdisk0...n>.path[0...n]->fieldname`. The following statistics can be accessed in synchronous mode and the asynchronous mode.

Note: The SCSI disk path I/O statistics supports only the IBM® multi-path driver (MPIO).

Table 5. SCSI disk path I/O statistics

Vue script field name	Data type	Description
name	String[32]	Disk name
block_size	unsigned long long	Disk block size in bytes
transfers	unsigned long long	Number of transfers to or from disk
rd_block_count	unsigned long long	Number of disk blocks read
wr_block_count	unsigned long long	Number of blocks written

vSCSI client I/O statistics

The following table displays the supported field names for the virtual SCSI (vSCSI) client I/O statistics. These fields are specific to the vSCSI client instances. They can be accessed as `__stat.io.adapter.<vscsi0...n>->fieldname`. The following statistics can be accessed in the synchronous mode and the asynchronous mode.

Table 6. vSCSI client I/O statistics

Vue script field name	Data type	Description
name	String[32]	Device name
transfers	unsigned long long	Number of transfers to or from device
rd_block_count	unsigned long long	Number of blocks read
rd_service_time	unsigned long long	Total read or receive service time in microseconds.
rd_min_service_time	unsigned long long	Minimum read service time in microseconds.

Table 6. vSCSI client I/O statistics (continued)

Vue script field name	Data type	Description
rd_max_service_time	unsigned long long	Maximum write service time in microseconds.
wr_block_count	unsigned long long	Number of blocks written
wr_service_time	unsigned long long	Total write service time in microseconds.
wr_min_service_time	unsigned long long	Minimum write service time in microseconds.
wr_max_service_time	unsigned long long	Maximum write service time in microseconds.
wait_queue_depth	unsigned long long	Wait queue depth for the driver
accum_wait_queue_time	unsigned long long	Accumulated wait queuing time in microseconds.
min_wait_queue_time	unsigned long long	Minimum wait queue time in microseconds.
max_wait_queue_time	unsigned long long	Maximum wait queue time in microseconds.
num_queue_full	unsigned long long	Number of in flight queue full count

vSCSI client driver statistics

The following table displays the supported field names for the virtual SCSI (vSCSI) client driver statistics. These fields are specific to the vSCSI client instances. They can be accessed as `_stat.io.adapter.<vscsi[0...n]>->fieldname`. The following statistics can be accessed in the synchronous mode and the asynchronous mode.

Table 7. vSCSI client driver statistics

Vue script field name	Data type	Description
no_dma_failures	unsigned char	Number of times system failed to send a I/O command due to insufficient Direct Memory Space (DMA). For example, DMA_NORES
no_cmd_elem_failures	unsigned char	Number of times system failed to send a I/O command because of no free command element with the client driver.
num_ping_timeouts	unsigned char	Number of times the client driver's ping request to that mapped Virtual I/O Server (VIOS) failed.
num_bad_mad	unsigned char	Number of times system failed to process a management datagram because of the adapter is not in an active state.

Table 7. vSCSI client driver statistics (continued)

Vue script field name	Data type	Description
num_hcall_drops	unsigned char	Number of times system failed to send a command to host CRQ (on VIOS), because the command-response queue (CRQ) is full. For example, the H_SEND_CRQ() parameter failed with the H_DROPPED parameter.

Fiber Channel driver statistics

The following table displays the supported field names for the Fiber Channel driver statistics. These fields are specific to the Fiber Channel device instances. The syntax of the Vue statement is `__stat.io.adapter.fcs[0..n]->fieldname`. The following statistics are accessible only in asynchronous mode.

Table 8. Fiber Channel driver statistics

Vue script field name	Data type	Description
secs_since_last_reset	unsigned long long	Time in seconds since last reset
tx_frames	unsigned long long	Number of frames transmitted
tx_words	unsigned long long	Fibre channel kbytes transmitted
rx_frames	unsigned long long	Number of frames received
rx_words	unsigned long long	Fibre channel kbytes received
lip_count	unsigned long long	Count of Loop Initiation Primitive (LIP) events on Fibre Channel Arbitrated Loop (FC-AL)
nos_count	unsigned long long	Count of No Operating System (NOS) events
error_frames	unsigned long long	Number of frames received with the cyclic redundancy check (CRC) error or discard frames. Each adapter has different field assigning this field.
lost_frames	unsigned long long	Number of lost frames
link_fail_count	unsigned long long	Count of link failures
sync_loss_count	unsigned long long	Count of loss of sync
sig_loss_count	unsigned long long	Count of loss of signal
prim_seq_proto_errcount	unsigned long long	Count of primitive sequence errors
inval_words_received	unsigned long long	Count of invalid transmission words received
inval_crc_count	unsigned long long	Count of CRC errors in a received frames
num_interrupts	unsigned integer	Total number of interrupts

Table 8. Fiber Channel driver statistics (continued)

Vue script field name	Data type	Description
num_spurious_interrupts	unsigned integer	Total number of spurious interrupts.
elastic_buf_overrun_errcount	unsigned integer	Number of times the link interface exceeds the elastic buffer overrun.
in_reqs	unsigned long long	Input requests
out_reqs	unsigned long long	Output requests
ctrl_reqs	unsigned long long	Control requests
in_bytes	unsigned long long	Input bytes
out_bytes	unsigned long long	Output bytes
no_dma_resource_count	unsigned long long	Count of DMA failures due to no DMA resources available
no_adap_elems_count	unsigned long long	Count of failure to allocate an adapter command element due to no more command element available
no_cmd_resource_count	unsigned long long	Count of failure to allocate a command due to no command resources available
adap_num_active_cmds	unsigned integer	Number of active commands in the adapter driver
adap_active_high_wmark	unsigned integer	High water mark of active requests in adapter driver
adap_num_pending_cmds	unsigned integer	Number of pending commands in the adapter driver
adap_pending_high_wmark	unsigned integer	High water mark of pending requests in adapter driver
adap_heldoff_num_cmds	unsigned integer	Number of commands in the adapter driver held off queue
adap_heldoff_high_wmark	unsigned integer	High water mark of number of commands in the adapter driver held off queue
proto_num_active_cmds	unsigned integer	Number of active commands in SCSI-FC driver
proto_active_high_wmark	unsigned integer	High water mark of active requests in SCSI-FC driver
proto_num_pending_cmds	unsigned integer	Number of pending commands in SCSI-FC driver
proto_pending_high_wmark	unsigned integer	High water mark of pending requests in SCSI-FC driver

Network device driver statistics

The following table shows the supported field name for the network device driver statistics. These fields are specific to the network device instances. Network device driver statistics can be accessed as "`__stat.net.adapter.<ent0...n>->fieldname`"

Vue script field name	Data type	Description	Field access type (Asynchronously or Both)
flags	unsigned int	Adapter flag values. This field can have following values: <ul style="list-style-type: none"> • NDD_UP • NDD_BROADCAST • NDD_DEBUG • NDD_RUNNING • NDD_SIMPLEX • NDD_DEAD • NDD_LIMBO • NDD_PROMISC • NDD_ALTADDRS • NDD_MULTICAST • NDD_DETACHED • NDD_64BIT • NDD_HIGHFUNC_QOS • NDD_MEDFUNC_QOS • NDD_MINFUNC_QOS • NDD_QOS • NDD_CHECKSUM_OFFLOAD • NDD_PSEG • NDD_ETHERCHANNEL • NDD_VLAN • NDD_SPECFLAGS These values are available as symbolic constants.	both
max_mtu	unsigned int	Maximum transmission unit.	both
min_mtu	unsigned int	Minimum transmission unit.	both
type	unsigned int	Interface types. This field can have following values: <ul style="list-style-type: none"> • NDD_ETHER • NDD_ISO88023 • NDD_ISO88024 • NDD_ISO88025 • NDD_ISO88026 These values are available as symbolic constant. Note: Not all possible interface type values are defined, and hence there might be other options present in the value.	both
physaddr	mac_addr_t	Physical or MAC address.	both

Table 9. Network device driver statistics (continued)

Vue script field name	Data type	Description	Field access type (Asynchronously or Both)
adapter_type	unsigned int	<p>Extension of the flag field. This field can have following values:</p> <ul style="list-style-type: none"> • NDD_2_SEA • NDD_2_VIOENT • NDD_2_VASI • NDD_2_HEA • NDD_2_IPV6_LSO • NDD_2_IPV6_CSO • NDD_2_IPV6_PARTIAL_CSO • NDD_2_IPV4_PARTIAL_CSO • NDD_2_LARGE_RECEIVE • NDD_2_ARPINPUT • NDD_2_ECHAN_ELEM • NDD_2_SEA_ELEM, • NDD_2_ROCE • NDD_2_VIRTUAL_PORT • NDD_2_PHYS_LINK_UP • NDD_2_VNIC <p>These values are available as symbolic constants.</p> <p>Note: Not all possible extension of flag field values are defined, and hence there might be other options present in the value.</p>	both
vlan_id	unsigned int	Virtual LAN (VLAN) identifier (Bits 0 - 11 are used for VLAN ID).	both
vlan_pri	unsigned int	VLAN priority (Bits 13 - 15 are used for VLAN priority).	both
alias	String[16]	Name of the network adapter alias.	Asynchronously
nobufs	unsigned long long	Number of times network buffers (MBUFs) were not available to the device driver.	Asynchronously
tx_packets	unsigned long long	Number of packets transmitted successfully by the network device.	Asynchronously
tx_bytes	unsigned long long	Number of bytes transmitted successfully by the network device.	Asynchronously
tx_interrupts	unsigned long long	Number of transmit interrupts received by the driver from the adapter.	Asynchronously
tx_errors	unsigned long long	Number of transmit errors on the network device. This is a counter for unsuccessful transmissions due to hardware or network errors.	Asynchronously
tx_packets_dropped	unsigned long long	Number of times packet dropped at the time of data transmission. The number of packets are accepted by device driver for transmission which were not given to the device.	Asynchronously
tx_queue_overflow	unsigned long long	The number of outgoing packets that have overflowed the software transmit queue.	Asynchronously
tx_queue_size	unsigned long long	The maximum number of outgoing packets queued to the software transmit queue.	Asynchronously
tx_queue_len	unsigned long long	The number of pending outgoing packets on the current software and hardware transmit queues.	Asynchronously
tx_broadcast_packets	unsigned long long	The number of broadcast packets transmitted.	Asynchronously
tx_multicast_packets	unsigned long long	The number of multicast packets transmitted.	Asynchronously
tx_carrier_sense	unsigned long long	The number of unsuccessful transmissions due to no carrier sense error.	Asynchronously
tx_DMA_underrun	unsigned long long	The number of unsuccessful transmissions due to direct memory access (DMA) underrun error.	Asynchronously

Table 9. Network device driver statistics (continued)

Vue script field name	Data type	Description	Field access type (Asynchronously or Both)
tx_lost_CTS_errors	unsigned long long	The number of unsuccessful transmissions due to loss of the Clear-to-Send signal error.	Asynchronously
tx_timeout_errors	unsigned long long	The number of unsuccessful transmissions due to timeout errors reported by the network adapter.	Asynchronously
tx_max_collision_errors	unsigned long long	The number of unsuccessful transmissions due to many collisions of transmitted packets. In this case, the number of collisions of transmitted packets encountered exceeds the number of retries on the network adapter.	Asynchronously
tx_late_collision_errors	unsigned long long	The number of unsuccessful transmissions due to the late collision error.	Asynchronously
tx_deferred	unsigned long long	The number of packets that are deferred for transmission.	Asynchronously
tx_hw_q_len	unsigned long long	The number of outgoing packets that currently exist on the hardware transmit queue.	Asynchronously
tx_sw_q_len	unsigned long long	The number of outgoing packets that currently exist on the software transmit queue.	Asynchronously
tx_single_collision_count	unsigned long long	Number of single collision errors at transmission.	Asynchronously
tx_multiple_collision_count	unsigned long long	Number of multiple collision errors at transmission.	Asynchronously
sqe_test	unsigned long long	The number of Signal Quality Error (SQE) tests performed successfully during transmission.	Asynchronously
ucast_pkts_reqs	unsigned long long	Number of outbound unicast packets requested by the network device.	Asynchronously
mcast_pkts_reqs	unsigned long long	Number of outbound multicast packets requested by the network device.	Asynchronously
bcast_pkts_reqs	unsigned long long	Number of outbound broadcast packets requested by the network device.	Asynchronously
rx_packets	unsigned long long	Number of packets received successfully by the network device.	Asynchronously
rx_bytes	unsigned long long	Number of bytes received successfully by the network device.	Asynchronously
rx_interrupts	unsigned long long	Number of receive interrupts received by the driver from the adapter.	Asynchronously
rx_errors	unsigned long long	The number of input errors encountered on this device. This is a counter for unsuccessful reception due to hardware or network errors.	Asynchronously
rx_packets_dropped	unsigned long long	Number of times packet dropped at the time of data reception. The number of packets received by the device driver from this device which were not given to a network demuxer.	Asynchronously
rx_bad_packets	unsigned long long	The number of bad packets received by the device driver.	Asynchronously
rx_broadcast_packets	unsigned long long	Number of broadcast packets received.	Asynchronously
rx_multicast_packets	unsigned long long	Number of multicast packets received.	Asynchronously
rx_noresource_errors	unsigned long long	The number of incoming packets dropped by the hardware due to the no resource error.	Asynchronously
rx_alignment_errors	unsigned long long	The number of incoming packets with the alignment errors.	Asynchronously
rx_DMA_overrun	unsigned long long	The number of incoming packets with the DMA overrun errors.	Asynchronously
rx_CRC_errors	unsigned long long	The number of incoming packets with the checksum error.	Asynchronously
rstart_cnt	unsigned long long	Number of times the adapter error recovery performed.	Asynchronously

Table 9. Network device driver statistics (continued)

Vue script field name	Data type	Description	Field access type (Asynchronously or Both)
rx_collision_errors	unsigned long long	The number of incoming packets with collision errors during the reception.	Asynchronously
rx_packet_tooshort_errors	unsigned long long	The number of incoming packets with the length error indicating that the packet size is less than the Ethernet minimum packet size.	Asynchronously
rx_packet_toolong_errors	unsigned long long	The number of incoming packets with the length error indicating that the packet size is bigger than the Ethernet maximum packet size.	Asynchronously
rx_packets_discardedbyadapter	unsigned long long	The number of incoming packets dropped by the adapter for any other reasons.	Asynchronously
rx_start	unsigned long long	The number of times that the receiver on the adapter has been started.	Asynchronously

Network interface based statistics

The following table shows the supported field name for the network interface based statistics. These fields are specific to the network interface instances. The network interface based statistics can be accessed in synchronous and asynchronous modes. The Network interface based statistics can be accessed as "`__stat.net.interface.<en0...n>->fieldname`"

Table 10. Network interface based statistics

Vue script field name	Data type	Description
mtu	unsigned long long	Maximum transmission unit. The maximum size of packets in bytes, that are transmitted using the interface.
flags	unsigned long long	Interface flag. This flag can have following values: <ul style="list-style-type: none"> • IFF_UP • IFF_BROADCAST • IFF_DEBUG • IFF_LOOPBACK • IFF_POINTOPOINT • IFF_VIPA • IFF_NOTRAILERS • IFF_RUNNING • IFF_PROMISC • IFF_NOARP <p>These values are available as symbolic constants.</p> <p>Note: Not all possible interface flag values are defined, and hence there might be other options present in the value.</p>

Table 10. Network interface based statistics (continued)

Vue script field name	Data type	Description
type	unsigned int	Interface type. This field can have following values: <ul style="list-style-type: none"> • IFT_ETHER • IFT_IB • IFT_LOOP • IFT_FDDI • IFT_ISO88023 • IFT_ATM • IFT_OTHER • IFT_TUNNEL These values are available as symbolic constants. Note: Not all possible interface flag values are defined, and hence there might be other options present in the value.
ipackets	unsigned long long	The number of packets received on this network interface.
ibytes	unsigned long long	The number of bytes received on this network interface.
ierrors	unsigned long long	The number of input errors. For example, malformed packets, checksum errors, or insufficient buffer space in the device driver.
opackets	unsigned long long	The number of packets transmitted on this network interface.
obytes	unsigned long long	The number of bytes transmitted on this network interface.
oerrors	unsigned long long	The number of output errors. For example, a fault in the local host connection or adapter output queue overrun.
collisions	unsigned long long	The number of packet collisions detected on carrier sense multiple access (CSMA) interfaces.
if_arpdrops	unsigned long long	Dropped because no Address Resolution Protocol (ARP) response.
if_iqdrops	unsigned long long	Number of times packet dropped at the time of data reception on this network interface.
index	unsigned int	Interface index number.
tx_mcasts	unsigned long long	The number of multicast packets transmitted on this network interface.
rx_mcasts	unsigned long long	The number of multicast packets received on this network interface.
no_proto	unsigned long long	Unsupported protocol.
bitrate	unsigned int	Bitrate as the rate at which data is sent on the wire.

Table 10. Network interface based statistics (continued)

Vue script field name	Data type	Description
dev_num	unsigned long long	Device Number.
options	unsigned int	Options field. This field can have following values: <ul style="list-style-type: none"> • IFO_FLUSH • IFO_HIGHFUNC_QOS • IFO_MEDFUNC_QOS • IFO_MINFUNC_QOS • IFO_QOS • IFO_THREAD • IFO_LARGESEND • IFO_PKTCHAIN • IFO_AACCT • IFO_MONITOR • IFO_VIRTUAL_ETHERNET • IFO_CSO_IPV6 • IFO_LSO_IPV6 • IFO_PARTIAL_CSO_IPV6 • IFO_PARTIAL_CSO_IPV4 • IFO_RNIC • IFO_FIRSTALIAS • IFO_PSEUDO_CLUSTER These values are available as symbolic constants. Note: Not all possible interface option values are defined, and hence there might be other options present in the value.

Network protocol based statistics

Network protocol based statistics

The following table shows the supported field name for the network protocol based statistics. These fields are specific to the specified network protocol. They can be accessed as `"__stat.net.protocol.<protocol name>->fieldname"`. For example, the IPv4 protocol specific statistics can be accessed as `"__stat.net.protocol.ip->fieldname"`. These statistics can be accessed in synchronous and asynchronous modes.

The IPv4 protocol based statistics can be accessed as `"__stat.net.protocol.ip->fieldname"`

Table 11. Network protocol based statistics (IPv4)

Vue script field name	Data type	Description
ipackets	unsigned long long	Total number of IP packets received.
rx_bytes	unsigned long long	Total number of bytes received in IP datagrams.
tx_bytes	unsigned long long	Total number of IP data bytes transmitted.

Table 11. Network protocol based statistics (IPv4) (continued)

Vue script field name	Data type	Description
bad_cksum	unsigned long long	The number of IP packets with bad header checksum.
shorts_pkts	unsigned long long	The buffer holding the IP packet has less bytes than the represented bytes in the total length field of the IP header (Total length includes IP header length + data).
small_pkts	unsigned long long	The buffer holding IP packet has less bytes than the represented bytes in IPv4 header length field.
bad_hdr_len	unsigned long long	The number of IP packets with bad IP headers. IP packets header length field has incorrect length (IP header length is smaller than minimum IP packet size).
bad_data_len	unsigned long long	The number of IP packets with bad length. IP packet total length field has less bytes than IP header length (Total length includes IP header length + data) or IP data size is bigger than the maximum supported packet size (IP_MAXPACKET).
bad_opts	unsigned long long	The number of IP packets with bad options.
bad_vers	unsigned long long	The number of IP packets with incorrect version number.
rx_frags	unsigned long long	The number of IP fragments received.
frag_drops	unsigned long long	The number of IP fragments dropped (duplicate or out of space).
frag_timeout	unsigned long long	The number of IP fragments dropped after timeout.
reassembled	unsigned long long	Total number of IP packets reassembled successfully.
forward	unsigned long long	The number of IP packets forwarded.
no_proto	unsigned long long	The number of unknown or unsupported protocol packets.
cant_fwd	unsigned long long	The number of packets that cannot be forwarded. The packets received from unreachable destination.
tx_redirect	unsigned long long	The number of redirects transmitted.
tx_drops	unsigned long long	The number of output packets dropped due to unavailability of network buffers (MBUF).
no_route	unsigned long long	The number of output packets discarded due to no route.
tx_frags	unsigned long long	The number of output fragments created.
cant_frag	unsigned long long	The number of datagrams that cannot be fragmented. The don't fragment flag is set.

Table 11. Network protocol based statistics (IPv4) (continued)

Vue script field name	Data type	Description
fragmented	unsigned long long	The number of output datagrams fragmented successfully.
threads_pkts	unsigned long long	The number of IP packets processed by kernel (dog) threads.
thread_drops	unsigned long long	The number of IP packets dropped by kernel threads because no more packets can be queued.
iqueueoverflow	unsigned long long	The number of IP packets dropped because the socket receive buffer is full.
pmtu_disc	unsigned long long	The number of successful path Maximum Transmission Unit (MTU) discovery cycles.
pmtu_redisc	unsigned long long	The number of path MTU rediscovery cycles attempted.
pmtu_guesses	unsigned long long	The number of path MTU discovery guesses or estimates due to no- response.
pmtu_timeouts	unsigned long long	The number of path MTU discovery response timeouts.
pmtu_decs	unsigned long long	The number of path MTU discovery decreases detected.
tx_pmtu_pkts	unsigned long long	The number of path MTU discovery packets transmitted.
pmtu_nomem	unsigned long long	The number of path MTU discovery memory allocation failures.
tx_dgd_pkts	unsigned long long	The number of dead gateway detection packets transmitted.
dgd_nomem	unsigned long long	The number of dead gateway detection (DGD) packets that are not transmitted due to allocation failures.
dgd_nogw	unsigned long long	The number of dead gateway detection (DGD) gateways that are not added due to allocation failures.
bad_src	unsigned long long	The number of packets with illegal source address.
delivered	unsigned long long	The number of IP packets consumed.
tx_local	unsigned long long	Total number of IP packets generated.
tx_raw	unsigned long long	Total number of RAW IP packets generated.
hdr_errs	unsigned long long	The number of header errors.
addr_errs	unsigned long long	The number of datagrams with IP address errors.
rx_discards	unsigned long long	The number of input datagrams discarded.
mcast_addr_errs	unsigned long long	The number of IP multicast packets dropped due to no receiver.

Table 11. Network protocol based statistics (IPv4) (continued)

Vue script field name	Data type	Description
rx_mcast_bytes	unsigned long long	The number of IP multicast bytes received.
tx_mcast_bytes	unsigned long long	The number of IP multicast bytes transmitted.
rx_mcast_pkts	unsigned long long	The number of IP multicast datagrams received.
tx_mcast_pkts	unsigned long long	The number of IP multicast datagrams transmitted.
rx_bcast_pkts	unsigned long long	The number of IP broadcast datagrams received.
tx_bcast_pkts	unsigned long long	The number of IP broadcast datagrams transmitted.
tx_mls_drops	unsigned long long	The number of outgoing IP packets dropped due to Multi Level Security (MLS) filters.
rx_mls_drops	unsigned long long	The number of IP incoming packets dropped due to MLS filters.

The Internet Protocol version 6 (IPv6) protocol based statistic can be accessed as "`__stat.net.protocol.ipv6->fieldname`"

Table 12. Network protocol based statistics (IPv6)

Vue script field name	Data type	Description
ipackets	unsigned long long	Total number of IPv6 packets received.
rx_bytes	unsigned long long	Total number of bytes received in IPv6 datagrams.
tx_bytes	unsigned long long	Total number of IPv6 data bytes transmitted.
raw_cksum	unsigned long long	The number of IPv6 packets that are not delivered due to bad raw IPv6 checksum.
shorts_pkts	unsigned long long	The MBUF does not have enough space to store IPv6 packet (IPv6 header + data).
small_pkts	unsigned long long	The MBUF does not have enough space to store the IPv6 header.
rx_nomem	unsigned long long	The number of times network buffers (MBUFs) were not available for input packets.
tx_nomen	unsigned long long	The number of times network buffers (MBUFs) were not available for output packets.
no_proto	unsigned long long	The number of unknown or unsupported protocol packets.

Table 12. Network protocol based statistics (IPv6) (continued)

Vue script field name	Data type	Description
bad_vers	unsigned long long	The number of IPv6 packets with incorrect version number.
rx_frags	unsigned long long	The number of IPv6 fragments received.
frag_drops	unsigned long long	The number of IPv6 fragments dropped (duplicate or out of space).
frag_timeout	unsigned long long	The number of IPv6 fragments dropped after timeout.
fragmented	unsigned long long	The number of output datagrams fragmented successfully.
tx_frags	unsigned long long	The number of output fragments created.
reassembled	unsigned long long	Total number of IPv6 packets reassembled successfully.
cant_frag	unsigned long long	The number of datagrams that cannot be fragmented. The don't fragment flag is set.
forward	unsigned long long	The number of IPv6 packets forwarded.
cant_fwd	unsigned long long	The number of packets that cannot be forwarded. The packets received from unreachable destination.
bad_src	unsigned long long	The number of packets with illegal source address.
tx_drops	unsigned long long	The number of output packets dropped due to unavailability of network buffers (MBUF).
no_route	unsigned long long	The number of output packets that are discarded due to no route.
delivered	unsigned long long	The number of IPv6 packets consumed.
tx_local	unsigned long long	Total number of IPv6 packets generated.
iqueueoverflow	unsigned long long	The number of IPv6 packets that are dropped due to the socket receive buffer is full.
big_pkts	unsigned long long	The number of IPv6 packets that are not forwarded because packet size is bigger than MTU.
tx_raw	unsigned long long	Total number of raw IPv6 packets generated.

Table 12. Network protocol based statistics (IPv6) (continued)

Vue script field name	Data type	Description
hdr_errs	unsigned long long	The number of header errors.
addr_errs	unsigned long long	The number of datagrams with IPv6 address errors.
rx_discards	unsigned long long	The number of input datagrams discarded.
rx_mcast_bytes	unsigned long long	The number of IPv6 multicast bytes received.
tx_mcast_bytes	unsigned long long	The number of IPv6 multicast bytes transmitted.
rx_mcast_pkts	unsigned long long	The number of IPv6 multicast datagrams received.
tx_mcast_pkts	unsigned long long	The number of IPv6 multicast datagrams transmitted.
rx_bcast_pkts	unsigned long long	The number of IPv6 broadcast datagrams received.
tx_bcast_pkts	unsigned long long	The number of IPv6 broadcast datagrams that are transmitted.
tx_mls_drops	unsigned long long	The number of IPv6 outgoing packets that are dropped due to MLS filters.
rx_mls_drops	unsigned long long	The number of IPv6 incoming packets that are dropped due to MLS filters.

The Transmission Control Protocol (TCP) based statistics can be accessed as `"__stat.net.protocol.tcp->fieldname"`

Table 13. Network protocol based statistics (TCP)

Vue script field name	Data type	Description
tx_total	unsigned long long	Total number of TCP packets transmitted. This count includes data and ack packets.
rx_total	unsigned long long	Total number of TCP packets received. This count includes data and ack packets.
opackets	unsigned long long	The number of TCP data packets transmitted.
ipackets	unsigned long long	The number of TCP data packets received.
tx_bytes	unsigned long long	The number of TCP data bytes transmitted.
rx_bytes	unsigned long long	The number of TCP data bytes received in sequence.

Table 13. Network protocol based statistics (TCP) (continued)

Vue script field name	Data type	Description
rexmit_pkts	unsigned long long	The number of TCP data packets re-transmitted.
rexmit_bytes	unsigned long long	The number of TCP data bytes re-transmitted.
tx_ack_pkts	unsigned long long	The number of TCP ACK only packets transmitted.
rx_ack_pkts	unsigned long long	The number of TCP ACK only packets received.
rx_ack_bytes	unsigned long long	The number of TCP ACK bytes received.
rx_dup_pkts	unsigned long long	The number of duplicate-only TCP packets received.
rx_dup_bytes	unsigned long long	The number of duplicate-only TCP bytes received.
rx_part_dup_pkts	unsigned long long	The number of packets with some duplicate data (part-duplicate packets) received.
rx_part_dup_bytes	unsigned long long	The number of duplicate bytes received from part-duplicate packets.
rx_dup_ack_pkts	unsigned long long	The number of TCP duplicate ACK packets received.
tx_win_probe	unsigned long long	The number of TCP window probe packets transmitted.
rx_win_probe	unsigned long long	The number of TCP window probe packets received.
tx_win_update	unsigned long long	The number of TCP window update packets transmitted.
rx_win_update	unsigned long long	The number of TCP window update packets received.
tx_delay_ack_pkts	unsigned long long	The number of TCP delayed ACK packets transmitted.
tx_urg_pkts	unsigned long long	The number of URG only packets transmitted.
tx_ctrl_pkts	unsigned long long	The number of control (SYN FIN RST) packets transmitted.
tx_large_send_pkts	unsigned long long	The number of large send packets transmitted.
tx_large_send_bytes	unsigned long long	The number of bytes transmitted using large send offload option.
tx_large_send_max	unsigned long long	Maximum number of bytes that can be transmitted using large send offload option.

Table 13. Network protocol based statistics (TCP) (continued)

Vue script field name	Data type	Description
rx_ack_unsent_data	unsigned long long	The number of ACKs received for an unsent data.
rx_out_order_pkts	unsigned long long	The number of out-of-order packets received.
rx_out_order_bytes	unsigned long long	The number of out-of-order bytes received.
rx_after_close_pkts	unsigned long long	The number of packets received after the connection has been closed.
fast_lo_conns	unsigned long long	The number of fastpath loopback connections.
tx_fast_lo_pkts	unsigned long long	The number of packets transmitted through fast path loopback connections.
rx_fast_lo_pkts	unsigned long long	The number of packets received through fast path loopback connections.
tx_fast_lo_bytes	unsigned long long	The number of bytes transmitted through fast path loopback connections.
rx_fast_lo_bytes	unsigned long long	The number of bytes received through fast path loopback connections.
rx_bad_hw_cksum	unsigned long long	The number of packets received with bad hardware checksum.
rx_bad_cksum	unsigned long long	The number of packets discarded due to bad checksum errors.
rx_bad_off	unsigned long long	The number of packets discarded due to bad error offset fields.
rx_short_pkts	unsigned long long	The number of packets discarded because packets are short. The packet size is smaller than minimum TCP packet size.
rx_queue_overflow	unsigned long long	The number of packets discarded because the listener queue is full.
rx_after_win_pkts	unsigned long long	The number of packets received with data that exceeded the receivers window size.
rx_after_win_bytes	unsigned long long	The number of bytes received with data that exceeded the receivers window size.
initiated	unsigned long long	The number of requests for TCP connection.

Table 13. Network protocol based statistics (TCP) (continued)

Vue script field name	Data type	Description
accepted	unsigned long long	The number of TCP connections accepted.
established	unsigned long long	The number of TCP connections established.
closed	unsigned long long	The number of TCP connections closed including connections drop.
dropped	unsigned long long	The number of TCP connections dropped.
ecn_conns	unsigned long long	The number of connections with Explicit Congestion Notification (ECN) capability.
ecn_congestion	unsigned long long	The number of times responded to ECN.
conn_drops	unsigned long long	The number of embryonic connections dropped.
segs_timed	unsigned long long	The number of times the segments attempt to update round trip time (RTT).
rtt_updated	unsigned long long	The number of times the segments updated the RTT.
ecnce	unsigned long long	The number of segments with congestion experienced(CE) bit set.
ecnwr	unsigned long long	The number of segments with congestion window reduced (CWR) bit set.
pmtu_resends	unsigned long long	The number of resends due to path MTU discovery.
pmtu_halts	unsigned long long	The number of path MTU discovery terminations due to retransmits.
rexmt_timeout	unsigned long long	The number of retransmit timeouts.
timeout_drops	unsigned long long	The number of connections dropped due to re-transmit timeouts.
fast_rxmt	unsigned long long	The number of fast retransmits.
new_reno_rxmt	unsigned long long	The number of NewReno fast retransmits.
false_fast_rxmt	unsigned long long	The number of times false fast retransmits avoided.
persist_timeouts	unsigned long long	The number of persist timeouts.

Table 13. Network protocol based statistics (TCP) (continued)

Vue script field name	Data type	Description
persist_drops	unsigned long long	The number of connections dropped due to persist timeouts.
keep_alive_timeout	unsigned long long	The number of keep alive timeouts.
keep_alive_probe	unsigned long long	The number of keep alive probes transmitted.
keep_alive_drops	unsigned long long	The number of connections dropped by keep alive.
delay_ack_syn	unsigned long long	The number of delayed ACKs for SYN.
delay_ack_fin	unsigned long long	The number of delayed ACKs for FIN.
sack_blocks_upd	unsigned long long	The number of times Selective Acknowledgments (SACK) blocks array is extended.
sack_holes_upd	unsigned long long	The number of times SACK holes array is extended.
tx_drops	unsigned long long	The number of packets dropped due to memory allocation failures.
time_wait_reuse	unsigned long long	Number of times an existing connection in TIME_WAIT state was reused for a new outgoing connection.
send_and_disc	unsigned long long	The number of send and disconnects.
spliced_conns	unsigned long long	The number of TCP spliced connections.
splice_closed	unsigned long long	The number of TCP spliced connections closed.
splice_resets	unsigned long long	The number of TCP spliced connections reset.
splice_timeouts	unsigned long long	The number of TCP spliced connections timeout.
splice_persist_drops	unsigned long long	The number of TCP spliced connections persist timeout.
splice_keep_drops	unsigned long long	The number of TCP spliced connections keep alive timeout.
bad_ack_conn_drops	unsigned long long	The number of connections dropped due to bad ACKs.
dup_syn_conn_drops	unsigned long long	The number of connections dropped due to duplicate SYN packets .

Table 13. Network protocol based statistics (TCP) (continued)

Vue script field name	Data type	Description
auto_cksum_offload	unsigned long long	The number of connections where checksum offload was dynamically disabled.
bad_syn	unsigned long long	The number of invalid packets discarded by listeners.
limit_transmit	unsigned long long	The number of times fast retransmit assisted by limited transmit algorithm.
pred_acks	unsigned long long	The number of times ACK packet headers correctly predicted.
pred_dat	unsigned long long	The number of times data packet headers correctly predicted.
paws_drops	unsigned long long	The number of segments dropped due to PAWS.
persist_drops	unsigned long long	The number of connection drops in persist state.
fake_syn_drops	unsigned long long	The number of fake SYN segments dropped.
fake_rst_drops	unsigned long long	The number of fake RST segments dropped.
data_inject_drops	unsigned long long	The number of data injection segments dropped.
tr_max_conn_drops	unsigned long long	Maximum connections dropped for TCP traffic regulation.
tr_nomem_drops	unsigned long long	The number of connections dropped for traffic regulation due to no memory.
tr_max_per_host	unsigned long long	Maximum connections per host dropped for traffic regulation.

The User Datagram Protocol (UDP) based statistics can be accessed as "`__stat.net.protocol.udp->fieldname`"

Table 14. Network protocol based statistics (UDP)

Vue script field name	Data type	Description
opackets	unsigned long long	Total number of UDP datagrams transmitted.
ipackets	unsigned long long	Total number of UDP datagrams received.
hdr_drops	unsigned long long	The number of packets whose size is smaller than header size. IP and UDP headers does not fit into single memory buffer (MBUF).

Table 14. Network protocol based statistics (UDP) (continued)

Vue script field name	Data type	Description
bad_cksum	unsigned long long	The number of UDP packets received with bad checksum errors.
bad_len	unsigned long long	The number of bad length packets received. The UDP length specified in the packet is either bigger than total packet size specified in IP header or smaller than size of UDP header.
no_socket	unsigned long long	The number of packets dropped due to no socket on the port.
sock_buf_overflow	unsigned long long	The number of times the socket buffer overflows.
dgm_no_socket	unsigned long long	The number of broadcast or multicast datagrams dropped due to no socket.
pcb_cache_miss	unsigned long long	The number of times input packets missing PCB cache.

The Internet Control Message Protocol (ICMP) based statistics can be accessed as "`__stat.net.protocol.icmp->fieldname`"

Table 15. Network protocol based statistics (ICMP)

Vue script field name	Data type	Description
sent	unsigned long long	Total number of ICMP packets transmitted.
received	unsigned long long	Total number of ICMP packets received.
errors	unsigned long long	The number of ICMP errors.
bad_cksum	unsigned long long	The number of ICMP messages received with bad checksum errors.
bad_len	unsigned long long	The number of ICMP messages received with bad length.
bad_code	unsigned long long	The number of ICMP messages with bad code fields. These messages have out of range icmp_code .
old_msg	unsigned long long	The number of errors not generated because the old packet protocol was ICMP.
old_short_msg	unsigned long long	The number of errors not generated because old IP packet was too short.

Table 15. Network protocol based statistics (ICMP) (continued)

Vue script field name	Data type	Description
short_msg	unsigned long long	The ICMP message size is smaller than minimum length of ICMP message (packet size <ICMP_MINLEN).
reflect	unsigned long long	The number of ICMP message responses generated.

The ICMPV6 protocol based statistics can be accessed as "__stat.net.protocol.icmpv6->fieldname"

Table 16. Network protocol based statistics (ICMPV6)

Vue script field name	Data type	Description
tx_echo_reply	unsigned long long	Total number of ICMPv6 echo replies transmitted.
rx_echo_reply	unsigned long long	Total number of ICMPv6 echo replies received.
errors	unsigned long long	The number of ICMPv6 errors.
rx_bad_cksum	unsigned long long	The number of ICMPv6 messages received with bad checksum errors.
rx_bad_len	unsigned long long	The number of ICMPv6 messages received with bad length.
bad_code	unsigned long long	The number of ICMPv6 messages with bad code fields. These messages have out of range icmp6_code .
old_msg	unsigned long long	The number of errors not generated because the old packet protocol was ICMPv6.
short_msg	unsigned long long	The ICMPv6 message size is smaller than minimum length of ICMPv6 message (packet size <ICMP6_MINLEN).
reflect	unsigned long long	The number of ICMPv6 message responses generated.
err_rate_limit	unsigned long long	The number of ICMPv6 errors beyond error rate limit.
tx_unreach	unsigned long long	The number of unreachable messages transmitted.
rx_unreach	unsigned long long	The number of unreachable messages received.
tx_big_pkt	unsigned long long	The number of times big ICMPv6 packets transmitted.
rx_big_pkt	unsigned long long	The number of times big ICMPv6 packets received.

Table 16. Network protocol based statistics (ICMPV6) (continued)

Vue script field name	Data type	Description
tx_timxceed	unsigned long long	The number of times ICMPv6 message sent time exceeded.
rx_timxceed	unsigned long long	The number of times ICMPv6 message received time exceeded.
tx_param_prob	unsigned long long	The number of times ICMPv6 messages transmitted with parameter problems.
rx_param_prob	unsigned long long	The number of times ICMPv6 messages received with parameter problems.
tx_echo_req	unsigned long long	The number of times echo request messages transmitted.
rx_echo_req	unsigned long long	The number of times echo request messages received.
tx_mld_qry	unsigned long long	The number of times group query requests transmitted.
rx_mld_qry	unsigned long long	The number of times group query requests received.
tx_mld_report	unsigned long long	The number of times group reports transmitted.
rx_mld_report	unsigned long long	The number of times group reports received.
rx_bad_mld_qry	unsigned long long	The number of times bad group queries received.
rx_bad_mld_report	unsigned long long	The number of times bad group reports received.
rx_our_mld_report	unsigned long long	The number of times our group reports received.
tx_mld_term	unsigned long long	The number of times group terminations transmitted.
rx_mld_term	unsigned long long	The number of times group terminations received.
rx_bad_mld_term	unsigned long long	The number of times bad group terminations received.
tx_redirect	unsigned long long	The number of times redirects transmitted.
rx_redirect	unsigned long long	The number of times redirects received.
rx_bad_redirect	unsigned long long	The number of times bad redirects received.
tx_router_sol	unsigned long long	The number of times router solicitations transmitted.

Table 16. Network protocol based statistics (ICMPV6) (continued)

Vue script field name	Data type	Description
rx_router_sol	unsigned long long	The number of times router solicitations received.
rx_bad_router_sol	unsigned long long	The number of times bad router solicitations received.
tx_router_adv	unsigned long long	The number of times router advertisements transmitted.
rx_router_adv	unsigned long long	The number of times router advertisements received.
rx_bad_router_adv	unsigned long long	The number of times bad router advertisements received.
tx_nd_sol	unsigned long long	The number of times neighbor solicitations transmitted.
rx_nd_sol	unsigned long long	The number of times neighbor solicitations received.
rx_bad_nd_sol	unsigned long long	The number of times bad neighbor solicitations received.
tx_nd_adv	unsigned long long	The number of times neighbor advertisements transmitted.
rx_nd_adv	unsigned long long	The number of times neighbor advertisements received.
rx_bad_nd_adv	unsigned long long	The number of times bad neighbor advertisements received.
tx_router_renum	unsigned long long	The number of times router re-numberings transmitted.
rx_router_renum	unsigned long long	The number of times router re-numberings received.
tx_haad_req	unsigned long long	The number of times home agent address discovery (HAAD) requests transmitted.
rx_haad_req	unsigned long long	The number of times HAAD requests received.
rx_bad_haad_req	unsigned long long	The number of times bad HAAD requests received.
tx_haad_reply	unsigned long long	The number of times HAAD replies transmitted.
rx_haad_reply	unsigned long long	The number of times HAAD replies received.
rx_bad_haad_reply	unsigned long long	The number of times bad HAAD replies received.
tx_prefix_sol	unsigned long long	The number of times prefix solicitations transmitted.

Table 16. Network protocol based statistics (ICMPV6) (continued)

Vue script field name	Data type	Description
rx_prefix_sol	unsigned long long	The number of times prefix solicitations received.
rx_bad_prefix_sol	unsigned long long	The number of times bad prefix solicitations received.
tx_prefix_adv	unsigned long long	The number of times prefix advertisements transmitted.
rx_prefix_adv	unsigned long long	The number of times prefix advertisements received.
rx_bad_prefix_adv	unsigned long long	The number of times bad prefix advertisements received.
no_mobility	unsigned long long	The number of mobility calls when not started.
ndp_q_drops	unsigned long long	The number of held packets dropped while waiting for ndp to complete.

The Internet Group Management Protocol (IGMP) based statistics can be accessed as "`__stat.net.protocol.igmp->fieldname`"

Table 17. Network protocol based statistics (IGMP)

Vue script field name	Data type	Description
rx_total	unsigned int	Total number of IGMP messages received.
rx_queries	unsigned int	The number of IGMP membership queries received.
tx_reports	unsigned int	The number of IGMP membership reports transmitted.
rx_reports	unsigned int	The number of IGMP membership reports received.
rx_our_reports	unsigned int	The number of IGMP membership reports received for our groups.
rx_bad_cksum	unsigned int	The number of IGMP messages received with bad checksum errors.
rx_short_msg	unsigned int	The number of IGMP messages received with few bytes whose size is smaller than minimum size of IGMP message.
rx_bad_queries	unsigned int	The number of IGMP membership queries received with invalid fields.
rx_bad_reports	unsigned int	The number of IGMP membership reports received with invalid fields.

The Address Resolution Protocol (ARP) based statistics can be accessed as "`__stat.net.protocol.arp->fieldname`"

Vue script field name	Data type	Description
<code>purged</code>	unsigned int	The number of ARP packets purged. When there is no space in the bucket, remove the oldest ARP entries from the bucket.
<code>sent</code>	unsigned int	Total number of ARP packets transmitted.

Memory statistics

Unlike other statistics providers memory statistics provides statistics that are as accurate as statistics fetched in synchronous mode but the statistics collection does not involve penalties of fetching the statistics repetitively during script execution.

Some of the items in memory statistics, such as the number of page faults, change frequently. Therefore, collecting memory statistics asynchronously produces incorrect information. For example, if you want to collect memory statistics for a process that fails because of low memory, you can use the interval between the `fork` and the `exit` system calls to collect the statistics. If memory statistics are not accumulated during the same intervals, the memory statistics do not change.

Memory statistics can be accessed by using the following format:

```
__stat.mem-><field>
```

For example, the following command indicates the number of page faults that occurred during a specific interval:

```
__stat.mem->page_faults
```

where `__stat` represents a statistical item, and `mem` indicates that the statistics is memory-related information.

The following independent fields can be used to display different statistical items:

Vue script field name	Data type	Description
<code>page_faults</code>	unsigned long long	Number of page faults.
<code>page_reclaims</code>	unsigned long long	Number of page reclamation.
<code>lock_misses</code>	unsigned long long	Number of lock misses.
<code>back_tracks</code>	unsigned long long	Number of backtrack operations.
<code>pageins</code>	unsigned long long	Number of pages that are paged in during the collection of memory statistics.
<code>pageouts</code>	unsigned long long	Number of pages that are paged out during the collection of memory statistics.
<code>num_ios</code>	unsigned long long	Number of I/O start operations.
<code>num_iodone</code>	unsigned long long	Number of iodone operations.

Table 19. ProbeVue fields for memory statistics (continued)

Vue script field name	Data type	Description
zerofills	unsigned long long	Number of zero-filled pages.
exec_fills	unsigned long long	Number of pages that are filled with executable files.
page_scans	unsigned long long	Number of pages that are examined.
pager_cycles	unsigned long long	Number of pager clock hand cycles.
page_steals	unsigned long long	Number of least recently used pages that must be included in the free list.
free_frame_waits	unsigned long long	Number of page frames that must be added to the list of available page frames.
extnd_xpt_waits	unsigned long long	Number of extend XPT wait operations.
pending_io_waits	unsigned long long	Number of pending I/O wait operations.

CPU statistics

All the CPU-related statistics are grouped into a new level that is called as *cpu*. All the CPU statistics items start with `__stat.cpu`. This statistic item is followed by either CPU instance, for example, `cpu0`, `cpu1` and so on, or a sublevel `cpu_total`.

Per logical CPU statistics

The per logical CPU statistics provides information about each logical CPU. In Vue script, each statistics is specified in the following format:

```
__stat.cpu.cpuX->fieldname
```

Where `cpuX` represents a CPU instance, for example, `cpu0`, `cpu1` and *fieldname* represents the statistics identifier. These fields are similar to the fields of the `perfstat_cpu_total_t` structure, which is defined in the `/usr/include/libperfstat.h` header file. This header file is part of the AIX `perfstat` library.

Table 20. Supported ProbeVue fields for per logical CPU statistics

Vue script field name	Data type	Description
user	unsigned long long	Raw number of clock ticks that are spent in the user mode.
sys	unsigned long long	Raw number of clock ticks that are spent in the system mode.
idle	unsigned long long	Raw number of clock ticks that are spent in the idle mode.
wait	unsigned long long	Raw number of clock ticks that are spent waiting for the completion of I/O operations.

Table 20. Supported ProbeVue fields for per logical CPU statistics (continued)

Vue script field name	Data type	Description
pswitch	unsigned long long	Number of context switches.
syscall	unsigned long long	Number of system calls.
sysread	unsigned long long	Number of read system calls.
syswrite	unsigned long long	Number of write system calls.
sysfork	unsigned long long	Number of fork system calls.
sysexec	unsigned long long	Number of exec system calls.
readch	unsigned long long	Number of characters that are transferred by using read system calls.
writetech	unsigned long long	Number of characters that are transferred by using write system calls.
bread	unsigned long long	Number of blocks that are read.
bwrite	unsigned long long	Number of blocks that are written.
lread	unsigned long long	Number of logical read requests.
lwrite	unsigned long long	Number of logical write requests.
phread	unsigned long long	Number of physical read operations, that is, read operations on raw devices.
phwrite	unsigned long long	Number of physical write operations, that is, write operations on raw devices.
iget	unsigned long long	Number of inode lookup operations.
namei	unsigned long long	Number of vnode lookup operations from a path name.
dirblk	unsigned long long	Number of 512-byte blocks that are read by the directory search routine to locate an entry for a file.
msg	unsigned long long	Number of Inter Process Communication (IPC) message operations.
sema	unsigned long long	Number of IPC semaphore operations.
minfaults	unsigned long long	Number of page faults with no I/O.
majfaults	unsigned long long	Number of page faults with disk I/O.

Table 20. Supported ProbeVue fields for per logical CPU statistics (continued)

Vue script field name	Data type	Description
puser	unsigned long long	Raw number of physical processor ticks in the user mode.
psys	unsigned long long	Raw number of physical processor ticks in the system mode.
pidle	unsigned long long	Raw number of physical processor ticks spent when the processor is idle.
pwait	unsigned long long	Raw number of physical processor ticks spent when the processor is waiting for the I/O operations.
redisp_sd0	unsigned long long	Number of thread redispaches within the scheduler affinity domain 0.
redisp_sd1	unsigned long long	Number of thread redispaches within the scheduler affinity domain 1.
redisp_sd2	unsigned long long	Number of thread redispaches within the scheduler affinity domain 2.
redisp_sd3	unsigned long long	Number of thread redispaches within the scheduler affinity domain 3.
redisp_sd4	unsigned long long	Number of thread redispaches within the scheduler affinity domain 4.
redisp_sd5	unsigned long long	Number of thread redispaches within the scheduler affinity domain 5.
migration_push	unsigned long long	Number of thread migrations from the local run queue to another queue because of starvation load balancing.
migration_S3grq	unsigned long long	Number of thread migrations from the global run queue to the local run queue resulting in a move across scheduling domain 3.
migration_S3pul	unsigned long long	Number of thread migrations from another processor's run queue resulting in a move across scheduling domain 3.
invol_cswitch	unsigned long long	Number of involuntary thread context switches.
vol_cswitch	unsigned long long	Number of voluntary thread context switches.
runque	unsigned long long	Number of threads in the run queue.
bound	unsigned long long	Number of bound threads.
decrintrs	unsigned long long	Number of decremter interrupts.
mpcrintrs	unsigned long long	Number of <i>Multi-Processor Communication (MPC) receive</i> interrupts.
mpcsintrs	unsigned long long	Number of <i>MPC send</i> interrupts.

Table 20. Supported ProbeVue fields for per logical CPU statistics (continued)

Vue script field name	Data type	Description
devintrs	unsigned long long	Number of device interrupts.
softintrs	unsigned long long	Number of off-level handlers that are called.
phantintrs	unsigned long long	Number of phantom interrupts.
idle_donated_purr	unsigned long long	Number of idle cycles that are donated by a dedicated partition that is enabled for donation.
idle_donated_spurr	unsigned long long	Number of idle Scaled Processor Utilization Resource Register (SPURR) cycles that are donated by a dedicated partition that is enabled for donation.
busy_donated_purr	unsigned long long	Number of busy cycles that are donated by a dedicated partition that is enabled for donation.
busy_donated_spurr	unsigned long long	Number of busy SPURR cycles that are donated by a dedicated partition that is enabled for donation.
idle_stolen_purr	unsigned long long	Number of idle cycles that are stolen by the hypervisor from a dedicated partition.
idle_stolen_spurr	unsigned long long	Number of idle SPURR cycles that are stolen by the hypervisor from a dedicated partition.
busy_stolen_purr	unsigned long long	Number of busy cycles that are stolen by the hypervisor from a dedicated partition.
busy_stolen_spurr	unsigned long long	Number of busy SPURR cycles that are stolen by the hypervisor from a dedicated partition.
hpi	unsigned long long	Number of hypervisor page-in operations.
puser_spurr	unsigned long long	Number of SPURR cycles that are spent in the user mode.
psys_spurr	unsigned long long	Number of SPURR cycles that are spent in the kernel mode.
pidle_spurr	unsigned long long	Number of SPURR cycles that are spent in the idle mode.
pwait_spurr	unsigned long long	Number of SPURR cycles that are spent in the wait mode.
spurrflag	int	This field is set if the CPU is running in the SPURR mode.
localdispatch	unsigned long long	Number of local thread dispatches on the logical CPU.
neardispatch	unsigned long long	Number of near thread dispatches on the logical CPU.
fardispatch	unsigned long long	Number of far thread dispatches on the logical CPU.
cswitches	unsigned long long	Number of context switches.

Table 20. Supported ProbeVue fields for per logical CPU statistics (continued)

Vue script field name	Data type	Description
state	int	Indicates whether the CPU is offline or online. A value of zero for the state field indicates that the CPU is offline and a non-zero value indicates that the CPU is online.
tb_last	unsigned long long	Last timebase register value read.
vtb_last	unsigned long long	Last virtual timebase value read.
icount_last	unsigned long long	Last instruction count register value read.

Note: The following members require unpinned memory access: pswitch, syscall, sysread, syswrite, sysfork, sysexec, readch, writtech, bread, bwrite, lread, lwrite, pthread, phwrite, iget, namei, dirblk, msg, sema, decrintrs, mpcrintrs, mpcsintrs, devintrs, softintrs, phantintrs. Because the probes of ProbeVue scripts run in an environment, where all interrupts are disabled, if the required data is not available in the resident memory, the members that require unpinned memory access are assigned with 0 values.

Overall statistics for all CPUs

The overall statistics provides cumulative information about all the CPUs in the system. In Vue script, each statistics is specified by using the following format:

```
__stat.cpu.cpu_total->fieldname
```

Where `cpu_total` represents the predefined sublevel name and `fieldname` represents the actual statistics identifier. These fields are similar to the fields of the `perfstat_cpu_total_t` structure, which is defined in the `/usr/include/libperfstat.h` header file. This header file is part of the AIX `perfstat` library.

Table 21. Supported ProbeVue fields for overall CPU statistics

Vue script field name	Data type	Description
ncpus	int	Number of active logical processors.
ncpus_cfg	int	Number of configured processors.
description	String	Description of the processor in the format <code><type_official name></code> . For example, for POWER7 [®] processor-based systems, the description is <code>POWERPC_POWER7</code> .
processorHZ	unsigned long long	Speed of the processor that is specified in Hz.
user	unsigned long long	Total number of raw clock ticks that are spent in the user mode.
sys	unsigned long long	Total number of raw clock ticks that are spent in the system mode.
idle	unsigned long long	Total number of raw clock ticks that are spent in the idle mode.
wait	unsigned long long	Total number of raw clock ticks that are spent waiting for the I/O operations.

Table 21. Supported ProbeVue fields for overall CPU statistics (continued)

Vue script field name	Data type	Description
pswitch	unsigned long long	Number of process context switches.
syscall	unsigned long long	Number of system calls.
sysread	unsigned long long	Number of read system calls.
syswrite	unsigned long long	Number of write system calls.
sysfork	unsigned long long	Number of fork system calls.
sysexec	unsigned long long	Number of exec system calls.
readch	unsigned long long	Number of characters that are transferred by using read system calls.
writetech	unsigned long long	Number of characters that are transferred by using write system calls.
devintrs	unsigned long long	Number of device interrupts.
softintrs	unsigned long long	Number of software interrupts.
lbolt	unsigned long long	Number of clock ticks since last restart of the operating system.
loadavg1	unsigned long long	Specifies a value that is equal to the product of the expression $(1 \ll \text{SBITS})$ and the average number of runnable processes during the last 1 minute. The \ll symbol indicates that the left shift operator and SBITS value is defined in the <code>/usr/include/sys/proc.h</code> header file.
loadavg5	unsigned long long	Specifies a value that is equal to the product of the expression $(1 \ll \text{SBITS})$ and the average number of runnable processes during the last 5 minutes. The \ll symbol indicates that the left shift operator and SBITS value is defined in the <code>/usr/include/sys/proc.h</code> header file.
loadavg15	unsigned long long	Specifies a value that is equal to the product of the expression $(1 \ll \text{SBITS})$ and the average number of runnable processes during the last 15 minutes. The \ll symbol indicates that the left shift operator and SBITS value is defined in the <code>/usr/include/sys/proc.h</code> header file.
runque	unsigned long long	Length of the run queue, that is, processes that are ready to run.
swpque	unsigned long long	Length of the swap queue, that is, the number of processes that are waiting to be paged in.

Table 21. Supported ProbeVue fields for overall CPU statistics (continued)

Vue script field name	Data type	Description
bread	unsigned long long	Number of blocks read.
bwrite	unsigned long long	Number of block written.
lread	unsigned long long	Number of logical read requests.
lwrite	unsigned long long	Number of logical write requests.
phread	unsigned long long	Number of physical read operations, that is, read operations on raw devices.
phwrite	unsigned long long	Number of physical write operations, that is, write operations on raw devices.
runocc	unsigned long long	This field is updated whenever run queue is updated, that is, the run queue is occupied. This value can be used to compute the simple average of ready processes.
swpocc	unsigned long long	This field is updated whenever swap queue is updated. For example, when the swap queue is occupied. This value can be used to compute the simple average of processes that are waiting to be paged in.
iget	unsigned long long	Number of inode lookup operations.
namei	unsigned long long	Number of vnode lookup operations from a path name.
dirblk	unsigned long long	Number of 512-byte blocks that are read by the directory search routine to locate an entry for a file.
msg	unsigned long long	Number of IPC message operations.
sema	unsigned long long	Number of IPC semaphore operations.
rcvint	unsigned long long	Number of tty interrupts that are received by the process.
xmtint	unsigned long long	Number of tty transmit interrupts.
mdmint	unsigned long long	Number of modem interrupts.
tty_rawinch	unsigned long long	Number of raw input characters.
tty_caninch	unsigned long long	Number of canonical input characters.
tty_rawoutch	unsigned long long	Number of raw output characters.
ksched	unsigned long long	Number of kernel processes created.

Table 21. Supported ProbeVue fields for overall CPU statistics (continued)

Vue script field name	Data type	Description
koverf	unsigned long long	Number of kernel process creation attempts when the configuration limit of processes is reached or the user has forked to the maximum limit.
kexit	unsigned long long	Number of kernel processes that became zombie processes.
rbread	unsigned long long	Number of remote read operation requests.
rcread	unsigned long long	Number of cached remote read operations.
rbwrt	unsigned long long	Number of remote write operations.
rcwrt	unsigned long long	Number of cached remote write operations.
traps	unsigned long long	Number of traps.
ncpus_high	int	Index of highest processor online. The index for of the processor starts from 1 instead of 0.
puser	unsigned long long	Raw number of physical processor ticks in the user mode.
psys	unsigned long long	Raw number of physical processor ticks in the system mode.
pidle	unsigned long long	Raw number of physical processor ticks in the idle mode.
pwait	unsigned long long	Raw number of physical processor ticks when the processor is waiting for I/O operations.
decrintrs	unsigned long long	Number of decremter tick interrupts.
mpcrintrs	unsigned long long	Number of <i>Multi Processor Communication (MPC)</i> receive interrupts.
mpcsintrs	unsigned long long	Number of <i>MPC send</i> interrupts.
phantintrs	unsigned long long	Number of phantom interrupts that are received by the partition.
idle_donated_purr	unsigned long long	Number of idle cycles that are donated by a dedicated partition that is enabled for donation.
idle_donated_spurr	unsigned long long	Number of idle SPURR cycles that are donated by a dedicated partition that is enabled for donation.
busy_donated_purr	unsigned long long	Number of busy cycles that are donated by a dedicated partition that is enabled for donation.
busy_donated_spurr	unsigned long long	Number of busy SPURR cycles that are donated by a dedicated partition that is enabled for donation.

Table 21. Supported ProbeVue fields for overall CPU statistics (continued)

Vue script field name	Data type	Description
idle_stolen_purr	unsigned long long	Number of idle cycles that are stolen by the hypervisor from a dedicated partition.
idle_stolen_spurr	unsigned long long	Number of idle SPURR cycles that are stolen by the hypervisor from a dedicated partition.
busy_stolen_purr	unsigned long long	Number of busy cycles that are stolen by the hypervisor from a dedicated partition.
busy_stolen_spurr	unsigned long long	Number of busy SPURR cycles that are stolen by the hypervisor from a dedicated partition.
iowait	short	Number of processes that are asleep waiting for buffered I/O operations.
physio	short	Number of processes that are waiting for raw I/O operations.
twait	unsigned long long	Number of threads that are waiting for file system direct IO(DIO) or Concurrent IO (CIO).
hpi	unsigned long long	Number of hypervisor page-in operations.
puser_spurr	unsigned long long	Number of SPURR cycles that are spent in the user mode.
psys_spurr	unsigned long long	Number of SPURR cycles that are spent in the kernel mode.
pidle_spurr	unsigned long long	Number of SPURR cycles that are spent in the idle mode.
pwait_spurr	unsigned long long	Number of SPURR cycles that are spent in the wait mode.
spurrflag	int	This field is set if the CPU is running in the SPURR mode.
tb_last	unsigned long long	Last timebase register value read.

Note: The following members require unpinned memory access: `devintrs`, `softintrs`, `loadavg1`, `loadavg5`, `loadavg15`, `decrintrs`, `mpcrintrs`, `mpcsintrs`, `phantintrs`. Because the probes of Vue scripts execute in an environment, where all interrupts are disabled, if the required data is not available in the resident memory, the members that require unpinned memory access are assigned with 0 values.

System statistics behavior

ProbeVue is a long running command and the Vue script is run multiple times. The Vue script might be referring to devices or resources to print the statistics and those resources and devices can become temporarily unavailable or permanently unavailable. Thus, ProbeVue system statistics statements have the following behaviors for different conditions:

Table 22. ProbeVue system statistics conditions and behaviors

Serial Number	Conditions	Behavior
1	Resource indicated in the statement does not exist.	During Vue script validation if ProbeVue does not find resource specified in the Vue script, compilation fails.
2	Device or resource is not open/active.	<ul style="list-style-type: none"> ProbeVue will print all zeros for the statistics. Device might get opened to run ioctl() function and then might get closed.
3	Device or resource moved from active state to inactive state.	ProbeVue will print zeros or stale data.
4	Device gets closed during run time.	Same results when resources are moved to inactive state.
5	Asynchronous fetch API takes time more than fetch interval set by user.	ProbeVue will not be able to guarantee the fetch interval if fetch takes more time than interval.
6	Dynamic LPAR operation performed during run time and resource indicated by Vue script are getting removed.	ProbeVue will print all zeros or stale data.
7	Device statistic fetch call is blocked due to timeout in device drivers.	ProbeVue will not be able to provide the guarantee of fetch interval and user might see the stale data.

Value and type assignment

The classification in the preceding section is one way to view the variables in a Vue script. The variable classes can be examined from a different perspective, namely how their values are derived.

External variables

Kernel class variables, entry and exit class variables and built-in variables are all external variables.

They exist independent of the ProbeVue framework and derive their values outside the context of any Vue script. ProbeVue allows the current values of external variables to be made available inside a Vue script. These variables are always read-only within the context of the Vue script. Any program statements that attempt to modify the value of an external variable will be flagged by the compiler as an illegal statement.

Although external variables have a pre-defined type, ProbeVue requires explicit declarations of all external variables, except for the built-in ones, in the Vue script that accesses them. The following table describes how the types of external variables are determined:

Variable	Type
Kernel global class	From the __kernel declaration statement of the kernel variable.
Entry class	From the function prototype declaration in the Vue script. Must specify the data types of each argument being used in the Vue script.
Return value from kernel functions	From the function prototype declaration in the Vue script. Must provide the type of the return value.

Variable	Type
Built-ins	These are generally dependent upon the underlying kernel variable. Their defined types and the equivalent ProbeVue types are as follows:

Built-in	Defined type	ProbeVue type
__tid	tid_t	long long
__pid	pid_t	long long
__ppid	pid_t	long long
__pgid	pid_t	long long
__pname	char [32]	String [32]
__uid	uid_t	unsigned int
__euid	uid_t	unsigned int
__trcid	pid_t	long long
__errno	int	int
__kernelmode	int	int
__r3..__r10	32-bit for 32-bit process 64-bit for 64-bit process	unsigned long
__curthread	N/A	All the members is long long
__curproc	N/A	All the members except for cwd is long long. The cwd member is of type string.
__ublock	N/A	All the members is long long

Note: The maximum size of the returned data can be smaller than the size of the type. For example, process IDs in AIX can fit in a 32-bit integer, while the pid_t data type is a 64-bit integer for 64-bit processes and the kernel.

Script variables

A script variable is either an automatic, thread-local, or global class variable.

Script variables exist only inside the context of a Vue script and their values are assigned from the script. Further, they can only be accessed or modified inside the script that defines them.

In general, you must explicitly declare the data type of a script variable through a declaration statement. However, the compiler can implicitly determine the data type of a program variable in some limited cases if the first reference to the variable is an assignment operation with the variable on the left-hand side of the assignment operator.

Implicit type determination for integral types

To be assigned an integral type, the right-hand side of the assignment must be in one of the following situations:

- A constant number.
- Another variable of integral type including built-in variables. Assigning from a variable whose type is not known is an error.
- A Vue function that returns an integral type like the **diff_time** function.

- Casting the expression on the right-hand side to an integral type, although this might give a warning in some cases.
- An expression involving any of the preceding situations.

The variable takes its type in addition to its value based on the expression on the right-hand side. In addition, the class of the variable can be assigned to the variable by prefixing it to the variable. The following script demonstrates some examples:

```

/*
 * File:    implicit2.e
 * Usage:   Demonstrates implicit assignment for integer types
 */

int read(int fd, char *p, long size);

@@BEGIN
{
    count = 404;           /* count: int of global class */
    zcount = 2 * (count - 4); /* zcount: int of global class */
    llcount = 33459182089021LL; /* lcount: long long of global class */
    lxcount = 0xF00000000245B20LL; /* xcount: long long of global class */
}

@@syscall:$1:read:entry
{
    __auto probev_timestamp_t ts1, ts2;
    int gsize;
    ts1 = timestamp();
    auto:dcount = llcount - lxcount; /* dcount: long long of auto class */

    auto:mypid = __pid; /* mypid: pid_t (64-bit integer) of automatic class */
    fd = __arg1; /* fd: int of global class */

    /* The following cast will likely cause a compiler warning
     * but can be ignored here
     */
    global:bufaddr = (long)__arg2; /* bufaddr: long of global class */

    gsize = __arg3;
    thread:size = gsize + 400; /* size: int of thread-local class */

    printf("count = %d, zcount = %lld\n", count, zcount);
    printf("llcount = %lld, lxcount = 0x%016llx, diff = %lld\n",
           llcount, lxcount, dcount);
    printf("mypid = %ld, fd = %d, size = %d\n", mypid, fd, size);
    printf("bufaddr = 0x%08x\n", bufaddr);
    ts2 = timestamp();

    auto:diff = diff_time(ts1, ts2, MICROSECONDS); /* diff: int of automatic class */
    printf("Time to execute = %d microseconds\n", diff);

    exit();
}

```

Note: The presence in the preceding script of a shell positional-like parameter, namely the **\$1** symbol in the **@@syscall:\$1:read:entry** probe specification. The **syscall** probe manager allows a process ID for the second field to indicate that the system call probe point must be enabled for a specific process only. Rather than hard-code a specific process ID, the second field has been set to a shell positional parameter in this script to permit the actual process ID to be passed as an argument at the time the script is issued. The **probev** command replaces any shell positional parameters in the script with the respective arguments passed on the command line.

Assuming that you are probing a process that has process ID 250000, the following script shows an example of running the **implicit2.e** script.

```

# probev implicit2.e 250000
WRN-100: Line:29 Column:26 Incompatible cast
count = 404, zcount = 800
llcount = 33459182089021, lxcount = 0x0f00000000245b20, diff = -1080830451389212643
mypid = 250000, fd = 10, size = 4496
bufaddr = 0x20033c00
Time to execute = 11 microseconds

```

In the preceding example, the **\$1** symbol in the script is automatically replaced with "250000", thus restricting the read system call entry probe point to the process with process ID equal to 250000.

Implicit type determination for string type

To be assigned a string type, the right-hand side of the assignment must be in one of the following situations:

- A string literal that is a sequence of characters within double quotation mark.
- Another variable of type string including built-in variables.
- A `Vue` function that returns a string like the **`et_userstring`** function.
- An expression involving any of the preceding situations.

The following example demonstrates implicit string type assignment:

```
/*
 * File:    implicit3.e
 * Usage:   Demonstrates implicit assignment for string types
 */

int write(int fd, char *p, long size);

@@BEGIN
{
    s1 = "Write system call:\n";
}

@@syscall:$1:write:entry
{
    String s2[40];
    wbuf = get_userstring(__arg2, __arg3);
    s2 = s1;
    zbuf = s2;
    pstring = zbuf + wbuf;
    printf("%s\n", pstring);
}

@@syscall:$1:write:exit
{
    ename = __pname;
    printf("Exec name = %s\n", ename);
    exit();
}
```

A process ID must be passed as an argument to the script when issuing it to replace the **\$1** shell positional parameter variable.

Implicit type determination for list type

To be assigned a list type, the right-hand side of the assignment must be the **`list()`** function. **`list()`** function is supported from any clause.

Data models for 32-bit and 64-bit processes

AIX supports two development environments: the 32-bit and the 64-bit development environments.

Compilers on AIX offer the following two programming models:

ILP32

ILP32, acronym for integer, long, and pointer 32, is the 32-bit programming environment in AIX. The ILP32 data model provides a 32-bit address space with a theoretical memory limit of 4 GB.

LP64

LP64, acronym for long, and pointer 64, is the 64-bit programming environment on AIX. With the exception of data type size and alignments, LP64 supports the same programming features as the ILP32 model and is backward compatible with the most widely used int data type.

Accordingly, a program on AIX can be compiled to run as either a 32-bit program or a 64-bit program. The same Vue script can be issued for a process running in 32-bit or 64-bit mode. As per the data model specification, an external variable of type long accessed in a Vue script must be treated as 4 bytes long when the probed (or traced) process is a 32-bit process. The same variable must be treated as 8 bytes long when the probed process is a 64-bit process. The layout and size of a structure or union, which contains members that are pointers or long variables, will depend upon whether it is being viewed from the perspectives of a 32-bit process or a 64-bit process. To avoid confusion, Vue provides semantic rules for handling the two different data models in a logical and consistent manner based on the variable's class.

Size-invariant variable types

int, short, char and long long are size-invariant variable types.

They always have the same size whether in 32-bit mode or 64-bit mode irrespective of the class of the declared variable.

Type	Size
long long	8
int	4
short	2
char	1

Size-variant variable types

The size-variant variable types have both a 32-bit mode and a 64-bit mode.

Size variant data types are listed below:

Type	32-bit size	64-bit size
long	4	8
Pointer types	4	8

In the preceding table, a pointer type refers to types like **char ***, **int ***, **struct foo ***, **unsigned long ***, and so on.

The following semantic rules apply for variables that are defined with any of the preceding types, that is, for "longs" and "pointers". The rules apply whether the variables are members of a structure or union, or whether they are declared as individual variables:

Automatic class

The mode of the variable will depend upon the probed process's mode (32 or 64).

Thread-local class

The mode of the variable will depend upon the probed process's mode (32 or 64).

Global class

The variable is always treated as being in 64-bit mode irrespective of the probed process's mode. This allows the variable to be used safely by both 32-bit and 64-bit processes without losing any data.

Kernel global class

Kernel variables that are longs or pointers are always 64-bit mode as the only supported kernel for AIX 6.1 and beyond is the 64-bit kernel.

Entry class

If a long or pointer type is defined in the function prototype for any of the parameters to the function, the modes of the corresponding entry class variables (**__arg1** through **__arg32**) will depend upon the mode of the probed process (32 or 64).

Exit class

If a long or pointer type is defined in the function prototype as the type of the return value of the function, the mode of the exit class variable (**__rv**) will depend upon the mode of the probed process (32 or 64).

Built-in class

These variables generally have a size-invariant type with the exception of the **__r3** through **__r10** built-ins that are defined as having an unsigned long type and hence are 32-bit long for 32-processes and 64-bit long for 64-bit processes.

The **@@BEGIN** and **@@END** probes are always issued in 64-bit mode.

Data types in Vue

The Vue language accepts three special data types in addition to the traditional C-89 data types.

Data types derived from the C language

Data types derived from the C language

The Vue language supports most of the data types defined in the C-89 specification. They include the signed and unsigned versions of the integral data types: **char**, **short**, **int**, **long** and **long long**. A "plain" **char** is treated as unsigned while the other integral types, if unqualified, are treated as signed. This matches the implementation of C on PowerPC. The Vue language also supports the floating types: **float** and **double**. In addition to these basic types of the C language, Vue also supports derived types like the array, structure, union and pointer types, the enumeration type and some incomplete types like **void**.

Floating types

You can only use a floating point type in simple assignment expressions and as arguments for functions like **printf**. In particular, you cannot use floating-point variables as operands of any unary or binary operators other than the assignment operator.

Pointer types

You can use pointers to dereference kernel or application data. However, you cannot declare pointers to Vue script variables or take their addresses.

Character arrays

You cannot use a character array as a string as in C, but must use the string data type.

Incomplete types

You cannot use array types of unknown size.

Bit-field types

The Vue compiler ignores bit-field declarations and the layout of structure or union types that contain members, which are bit-fields, is undefined.

ILP32 and LP64 data models

Generally, a C program can be compiled in either 32-bit mode where it follows the ILP32 data model or in 64-bit mode where it follows the LP64 model. Because the same Vue clause can be issued by both 32-bit and 64-bit processes, Vue internally supports both models at the same time.

Range and bucket data type

Range and bucket data type

The range data type in Vue is designed to handle the distribution of data points for certain defined ranges. Each range of the variable defined of the range data type displays the count of the number of elements within the corresponding range value. The range data type will support Integral and String types of ranges. The distribution of range values for integral ranges can be the power of two distributions or a linear distribution. An example of linear and power of two distributions of range values follow:

Linear distribution:

Range	Count
0 - 5	2
5 - 10	4
10 - 15	1
Others	20

Power of 2 distribution:

Range	Count
1 - 2	2
2 - 4	9
4 - 8	2005
8 - 16	4
16 - 32	1999
32 - 64	7
Others	5

The previous distributions indicates the count of the number of elements greater than or equal to the lower bound of the range and less than the upper bound of the range value. For example, in the power of two distributions, the count of data starting from 4 and lesser than 8 is 2005. The count of the value that does not come under the defined ranges is displayed in the Others range.

Example of string ranges

Example:

Range	Count
Read, write, open	87
Close, foo1	3
foo2	1
Others	51

In the previous example, the distribution indicates the number of times a particular string occurs within the range values. In this example, read, write, and open has been called 87 times.

Declaration and initialization of the range data type:

The range data type can be declared by using the `range_t` keyword. For example, the following declaration in the Vue script defines two range data type variables:

```
range_t T1, T2; // T1 and T2 are of Range data type variables.
```

The `set_range` and the `add_range` routines are used initialize the integral and string ranges for any particular range data type variables.

Initializing integral range data type: The `set_range` routine will be used to initialize the integral ranges. The syntax of the `set_range` differs for linear and power of two distributions of range values. The syntax of the `set_range` routine for linear distribution follows:

```
void set_range(range_t range_data, LINEAR, int min, int max, int step);
```

Example:

```
set_range(T1, LINEAR, 0, 100, 10);
```

In the previous example, the `set_range` routine initializes the range data T1. The range data T1 has the linear distribution of values. The lower bound of T1 is 0 and the upper bound is 100. The size of each range is 10. The distribution for the previous example will look like the following:

Range	Count
0 - 10	...
10 - 20	...
20 - 30	...
...	...
...	...
90 - 100	...

The syntax for the initialization of the power of 2 distribution follows:

```
set_range(range_t range_data, POWER, 2);
```

Example:

```
set_range(T2, POWER, 2);
```

In this example the routine initializes the range data type T2 as the power of 2 distribution range type.

Initializing the string range data type: The `add_range` routine initializes the string range data type.

Syntax:

```
void add_range(range_t range_data , String S1, String S2, ..., String Sn);
```

Example:

```
add_range(T1, "read", "write", "open");
```

This routine add the strings `read`, `write` and `open` to a single slot of `range_t` data T1. Another `add_range` on the same `range_t` data T1 adds the strings to the next slot.

```
add_range(T1, "close", "func1", "func2");
```

This routine adds the strings `close`, `func1` and `func2` to the `range_t` data T1 in the next slot.

Note: The `range_t` range data type is a special data type for Vue, which only can be used to store as a value inside an associative array. For any other operations (such as arithmetic, logical, bitwise, and relational) on the `range_t` data type, they will fail and error out.

Notes: This information discusses various uses and initialization routines of the `range_t` data type.

1. Declaration of the `range_t` data type can be done only in the `@@BEGIN` clause.
2. Initialization of the `set_range` routine can be used only inside the `@@BEGIN` clause.
3. A range data type whose range values are integral can be initialized only once. The same variable cannot be initialized twice.

Example:

```
set_range(T1, LINEAR, 0, 50, 5); // Valid syntax
set_range(T1, LINERA, 10, 100, 10); // Error, cannot initialize an already
// initialized T1.
set_range(T1, POWER, 2); // Error, T1 has already initialized.
add_range(T1, "read", "write"); // Error, T1 has already initialized.
```

4. The parameters of `min`, `max` and `step` are integral constants for `set_range` routine.

Storing and printing the range data type:

The range data type can be stored in an associative array as a value by using the `qrangle` routine. The `qrangle` routine finds the slot number whose frequency and count need to be incremental.

Example:

For this example, T1 is a **range_t** data type whose range values are of the integral type.

```
qrangle(aso["read"], T1, time_spent);
```

In this example, the `qrangle` routine finds the slot number in which the `time_spent` fails and the count for that slot number are incremented for the associative `aso` array corresponding to the `read` key.

In the following example, the T2 is a **range_t** data type and range values are of the string type.

```
qrangle(aso["function usage"], T2, get_function());
```

In this example, `qrangle` routine finds the slot number in which the function passed as the third argument fails and increments the count for that slot for `aso` associative array corresponding to the **function usage** key.

Notes:

1. For any ASO, only one **range_t** type variable can be stored as a value. Using **qrangle** for two different types of **range_t** variable type for the same ASO will fail.

Example:

```
qrangle(aso["read"], T1, time_spent); // Correct syntax.
qrangle(aso["read"], T2, time_spent); // Error. Two different range_t types
// cannot be used for the same ASO.
```

The `quantize` and `lquantize` functions of the associative array whose value type is **range_t** shows the visual quantization of frequency and count of ranges.

2. While printing the string range a maximum of 40 characters (including the comma) can be printed for a particular slot. If the strings in a slot have more than 40 characters, the string range is truncated and is printed with the last 3 characters as dots(...).

Examples of range data type and the qrangle routine:

```
@@BEGIN
{
  __thread start ;
  range_t T1;
  set_range(T1, LINEAR, 0, 150, 10) ;
}
@@syscall :$__CPID :read :entry
{
  thread :tracing = 1 ;
  start = timestamp() ;
}
@@syscall :$__CPID :read :exit
  when(thread :tracing == 1)
{
  __auto long time_spent;
  currttime = timestamp() ;
  time_spent = diff_time(start, currttime, MICROSECONDS);
  qrangle(aso["read"], T1, time_spent);
}
@@END
{
  print(aso);
  quantize(aso);
}
```

Expected output for this example:

Key	Value
Read	count
0-11	4
10-20	6

	60-70	7	
	Others	32	
Key		Value	
Read	Range	count	
	0-10	4	===
	10-20	6	====
	60-70	7	=====
	Others	32	=====

Stack trace type

Stack trace type

A variable of type **stktrace_t** is used to hold the return value from the ProbeVue function **get_stktrace**, which returns the current stack trace. The stack trace returned is the stack trace of the current thread. This variable can also be stored in an associative array either as a key or as a value. The **stktrace_t** type is an abstract data type, and this variable cannot be used directly with the standard C unary or binary operators. Internally, this variable is an array of unsigned longs.

Vue supports the following characteristics and operations for the stack trace type variables:

Declaration of stack trace type variable

A variable can be declared to be of type stack trace by declaring it as follows in the script:

```
stktrace_t st;           // st is a stktrace_t variable.
st = get_stktrace(5);   // Get the stack trace up to five levels.
a_st[0] = get_stktrace(-1); // Get the stack trace up to the extent possible and
                        // store in the associative array a_st as value.
```

The qualifiers signed, unsigned, register, static, auto, thread, kernel, and const are not supported for the **stktrace_t** type variables.

Assignment operation

The assignment (=) operator allows a **stktrace_t** type variable to be assigned to another **stktrace_t** type variable. The original values in the target **stktrace_t** variables are destroyed. No type casting is allowed from or to the **stktrace_t** variable types. In the following example, the content of the stack trace *t1* is assigned to *t2*.

```
stktrace_t t1, t2;      // Declares two stack trace variables.
t1 = get_stktrace();    // Get the current stack trace in t1.
t2 = t1;                // Get the content of t1 into t2.
```

Comparison operation

Only equality (==) and inequality (!=) operators are allowed for the **stktrace_t** variables. The result of these operators is either True(1) or False(0) based on the entire entries of the **stktrace_t** variables. Comparisons of individual entries of **stktrace_t** variables are not allowed. No other comparison operators (>=, >, < or =<) are allowed for the **stktrace_t** type variables.

```
if( t1 == t2) // comparing two stktrace_t type variables.
    printf("Entries are similar");
else
    printf("Entries are not similar");
```

Printing stack trace type variable

The **stktrace_t** variable can be printed by using the **%t** or **%T** format specifiers in the **printf** function of Vue. The output is the symbolic stack trace of the thread that is saved in the variable. The symbol with the address (symbol plus address) are printed only when the thread that corresponds to the **stktrace_t** variable is in the running state and when the **%t** format specifier is used to print the stack trace; otherwise only the stack trace as the address is printed for the variable.

The **stktrace_t** type variable stored in associative array either as key or value can be printed with print function of associative array. Addresses with symbol (symbol name + offset) is printed if thread that

corresponds to the **stktrace_t** type stored in the associative array is running; otherwise only addresses is printed. When you set the **STKTRC_NO_SYM** flag in the **set_aso_print_options()** function, this variable prints raw addresses of the running thread.

```
stktrace_t t1;
t1 = get_stktrace (5);
printf ("%t", t1);      // Displays the stack trace stored in variable t1.
printf ("%T", t1);     // Displays the stack trace stored in variable t1 with raw
addresses.
a[__tid] = t1;         // Store t1 as value in an associative array a.
print(a) ;           // Print associative array a, whose value
// type is
stktrace_t variable.
```

Limitations for stack trace type variable

- The array of the **stktrace_t** variable cannot be declared.
- The **stktrace_t** variables cannot be used as a member of a struct or a union.
- Access of any individual entry of the stack trace is not allowed.
- The operations (assignment, comparison, and printing) of the **stktrace_t** type variables is not supported in the **sysrtrace** probe.

Special data types

In addition to the traditional C-89 data types, the Vue language also accepts seven special data types.

String type

The string data type is a representation of string literals. Unlike in C, the string is a basic data type in Vue

Having a string type avoids some of the confusion in C which does not support a string type but permits a string to be represented both by a pointer to a **char** type and by a character array. You can explicitly declare a string variable by using the string declaration statement. Explicitly declared string variables must also specify the maximum string length (similar to how character arrays are declared in C). Unlike C, strings in Vue are not explicitly terminated by a null character and you do not need to reserve space for it.

```
String s[40];          /* Defines a string 's' of length 40 */
s = "probevue";
```

Further, any string literal written in C-style with enclosing double-quotes is automatically assigned a string data type. Vue automatically converts an external variable that is declared as a C-style character data type (**char *** or **char[]**) to the string data type as needed.

You can use the following operators for the string data type:

- The concatenation operator: "+" .
- The assignment operator: "=" .
- Relative operators for comparing strings: "=", "!=", ">", ">=", "<" and "<=".

You can set a string variable to the empty string by assigning "" to it as the following example:

```
s = "";              /* Sets s to an empty string */
```

Unlike the C language, a pair of adjacent string literals is not concatenated automatically. The concatenation operator (+) must be explicitly applied as in the following example:

```
String s[12];
// s = "abc" "def";
/* ERROR: Commented out as this will result in a syntax error */
s = "abc" + "def";    /* Correct way to concatenate strings */
```

Vue supports several functions that accept a string data type as a parameter or return a value that has a string data type.

List type

A variable of type list collects a set of integral type values. The list type is an abstract data type and you cannot use a list variable directly with the standard C unary or binary operators.

You can use the following operations for the list type:

- A constructor function, **list()** to create a new list variable if its not defined before, if the variable is already defined - it should be cleared.
- A concatenation function, **append** to add an item to the list or to join two lists together.
- The "=" operator that allows a list to be assigned to another.
- A set of aggregation functions that operate on a list variable and return a scalar (integer) value like **sum**, **avg**, **min**, **max** and so on.

Although, you can use a list variable to collect any integral value, the values are always saved as 64-bit signed integers.

The **list()** function returns a new empty new list which must be assigned to a variable of list type. This will create a new variable of type list if the list variable on the left-hand side of the assignment operator had not been previously assigned to a list. It may also be assigned to an existing list variable in which case, any values collected in the target list are discarded. Further, a variable can be declared to be of type list by declaring it as follows anywhere in the Vue script:

```
__list l_opens;
```

The effect of this is as if the **list()** function was invoked in the `@@BEGIN` probe and the return value assigned to this list variable.

The following example creates a new list variable called **l_opens**:

```
l_opens = list();
```

The **list** function can be invoked from any clause. If you specify an existing list name when invoking the **list** function, the existing list is cleared.

You can use the **append()** function to add a value to a list variable. Each call to the **append** function adds a new value to the set of values already saved in the list variable. The following example shows how the size of the list variable grows with each call to the **append** function:

```
append(l_opens, n_opens1); /* l_opens = {n_opens1} */
append(l_opens, n_opens2); /* l_opens = {n_opens1, n_opens2} */
append(l_opens, n_opens3); /* l_opens = {n_opens1, n_opens2, n_opens3} */
append(l_opens, n_opens4); /* l_opens = {n_opens1, n_opens2, n_opens3, n_opens4} */
```

The second parameter to the **append()** function can also be a variable of type list which will append all its values to the target list specified by the first parameter. So **append** can also be used to join two lists.

In the following example, the contents of list b are added to list a:

```
a=list()
b=list()
append(a,b)
```

Note: The value added to the list must be a parameter of integral or list type and it is an error if any of the variables **n_opens1 -n_opens4** do not have an integral type. Any types that are smaller than a **long long** (like a short or an int) are automatically promoted to the **long long** type.

You can also use **append** to join two lists. The first argument is the target list and the second is the list source list. In the following example, the contents of list b are added to list a:

```
a=list()
b=list()
append(a,b)
```

The **append()** function has no return value.

A list can be assigned to another list using the assignment operator. The original values in the target list are destroyed. In the following example, the contents of the **l_opens2** list are lost (the items are removed) and the contents of the **l_opens** list are copied over to the **l_opens2** list.

```
l_opens2 = list();
append(l_opens2, n_opens5);

l_opens2 = l_opens;
/* l_opens and l_opens2 => {n_opens1, n_opens2, n_opens3, n_opens4} */
```

The aggregation functions can be applied on a list variable as shown in the following examples:

```
/* below we assume n_opens1=4, n_opens2=6, n_opens3=2 and n_opens4 = 4
 * at the time they were added to the l_opens list variable
 */
x = avg(l_opens); /* this will set x to 4 */
y = min(l_opens); /* this will set y to 2 */
z = sum(l_opens); /* this will set z to 16 */
a = count(l_opens) /* this will set a to 4 */
```

A list variable is useful when accurate aggregate values need to be recorded. List variables are updated atomically, so use them only when required, as they are less efficient than regular variables.

Associative array type

An associative array is a map or look-up table consisting of a collection of keys and their associated values. There is a one-to-one mapping between a set of keys and values. Associative arrays are supported by Perl, ksh93, and several other languages.

In Vue, each mapping is from one or more keys (dimensions) to a single value. Associative Array keys can have the following types

- integral
- floating point
- string
- timestamp
- stacktrace
- path
- MAC address
- IP address

Associative Array value can have the following types

- integral
- floating point
- string
- timestamp
- stacktrace
- path
- MAC address
- IP address
- list
- range

Associative arrays are abstract data types in Vue. Following actions can be performed on associative array data types.

- Binding keys to a value:

If the instance of the keys does not already exist, this action adds a key or set of keys along with the associated value to the associative array. Otherwise, this action replaces the associated old value with the new value. Keys that are not bound have a default value of 0 for the numerical type, empty string value for the string type, or a NULL value for other key types.

Following example explains binding keys to a value

```
/* single key dimension */
count["ksh"] = 1;
/* multiple key dimensions */
var[0]["a"][2.5] = 1;
var[1]["a"][3.5] = 2;
```

The very first use of an associative array variable sets the type of the keys, the key dimension size and the type of the value. It must remain the same throughout the other places in the Vue script.

For a key in ASO, you can bind a LIST value by completing the following different actions:

1. By assigning a LIST variable:

```
assoc_array["ksh"]=ll /* copies ll list into associative array */
assoc_array["ksh"]=assoc_array["abc"]; /* copies a list in ASO to another list in ASO.
Here the value type of assoc_array is LIST */
```

2. By assigning an empty list returned by **list()** constructor function:

```
assoc_array["ksh"]=list(); /* assigns an empty list */
```

3. By appending a list or integral value

```
append(assoc_array["ksh"], 5); /* integral value 5 is appended to the list in ASO */
append(assoc_array["ksh"], ll); /* appends the contents of LIST variable ll to the list in ASO*/
append(assoc_array["ksh"], assoc_array["abc"]); /* appends the contents of list in ASO to another list
in ASO */
```

- Unbinding a key or a set of keys and deleting the associated value: The **delete()** function is used to remove keys and binding values from the associative array. An unbound key is assigned a value of 0 or an empty string.

Following example explains how to use delete function for unbinding a key

```
delete(count, "ksh");
delete(var, 0, "a", 2.5);
```

The first argument is the name of the associative array variable. The array variable name must be followed by N comma separated keys, where N is the key dimension size. If you want to only delete an associated value on the basis of key dimensions other than N, you can specify ANY for that key dimension. For example, to delete all the elements that have the string "a" as the second dimension, enter the following command.

```
delete(var, ANY, "a", ANY);
```

All of the keys in **delete()** function can be specified as ANY, in which case, all elements of the associative array will be deleted. This function returns 0 if matching elements are found and are deleted. Otherwise, the **delete()** function returns a value of 1.

- Finding the value for a set of keys: This operation looks up values that are bound to a single key or multiple keys.

```
total = count["ksh"] + count["csh"];
prod = var[0]["a"][2.5] * var[1]["a"][3.5];
```

A LIST value for a key can be retrieved by indexing the associative array with the key. All LIST functions, **sum()**, **min()**, **max()**, **count()**, and **avg()**, can be used on a List in an Associative array. You can also assign a list in associative array to a LIST variable.

Example:

```
/* copies associative array list into list variable "ll" */
ll=assoc_array["ksh"];
/* prints the sum of all elements of list in associative array indexed with ksh */
printf("sum of assoc_array %d\n",sum(assoc_array["ksh"]));
/* prints the minimum value */
printf("min of assoc_array %d\n",min(assoc_array["ksh"]));
/* prints the maximum value */
printf("max of assoc_array %d\n",max(assoc_array["ksh"]));
/* prints the number of values in list */
printf("count of assoc_array %d\n",count(assoc_array["ksh"]));
/* prints average value of the list */
printf("avg of assoc_array %d\n",avg(assoc_array["ksh"]));
```

- Checking whether a key or set of keys exists: The `exists()` function checks whether the associative array has any element that corresponds to the given keys. The `exists()` function returns a value 1 if an element is found and the function returns 0 if an element is not found.

Following code block checks whether a key or set of keys exist.

```
if (exists(count, "ksh"))
    printf("Number of ksh calls = %d\n", count["ksh"]);
if (exists(var, 0, "a", 2.5))
    printf("Found value = %d\n", var[0]["a"][2.5]);
```

If you specify ANY keyword for a particular key dimension, the dimension becomes insignificant for search operation. All keys in the `exists()` function can be specified as ANY, in which case, the `exists()` function will check if the Associative Array has any element or not.

```
my_key = "a";
if (exists(var, ANY, my_key, ANY))
    printf("Found element with second key as %s \n", my_key);
```

- Increment and decrement operation: This operation can be used to increment or decrement the associative array values. To use this operation, you must specify an integer as the value type for the key. The following examples show the usage of the increment and decrement operation:

1. `printf("Incremented value = %d\n", ++count["ksh"]);`
2. `printf("Incremented value = %d\n", count["ksh"]++);`
3. `printf("Decrement value = %d\n", --count["ksh"]);`
4. `printf("Decrement value = %d\n", count["ksh"]--);`

In example 1, the value that corresponds to **ksh** key, incremented and the incremented value is printed.

In example 2, the value that corresponds to **ksh** is first printed and then the value is incremented.

The decrement operation works the same way. However, the increment or decrement operation can be performed only on associative arrays whose value type is integer. The increment or decrement operation can also be used as an aggregator, where the value type of the associative array by default is set as integer. For example, on encountering statement, `a[100]++` the first time, associative array `a` is created with the integer key type and the integer value type. The value stored for key 100 is 1. However, for `a[100]--`, -1 would be stored as the value for key 100. On encountering subsequent increment or decrement operations for the same associative array `a` increment and decrement operations is performed on the value for the specified key.

The increment and decrement behavior is exactly the same for associative arrays with more than one key dimension:

```
++var[0]["a"][2.5];
var[0]["a"][2.5]++;
--var[1]["a"][3.5];
var[1]["a"][3.5]--;
```

- Printing contents of an associative array: This operation prints the key and associated value of the elements for an associative array. You can specify the following print options:

Associative array print option	Description	Possible Values	Default value
<i>num-of-entries</i>	Specifies to print the first number of key-value pairs.	$n \geq 0$. (If 0, all the entries are displayed.)	0
<i>sort-type</i>	Specifies the sorting order.	SORT_TYPE_ASCEND, SORT_TYPE_DESCEND	SORT_TYPE_ASCEND
<i>sort-by</i>	Specifies whether to sort based on key or value.	SORT_BY_KEY, SORT_BY_VALUE	SORT_BY_KEY
<i>list-value</i>	Specifies which LIST attribute to sort or quantize when the value of the associative array is the list type	USE_LIST_SUM, USE_LIST_MIN, USE_LIST_MAX, USE_LIST_COUNT, USE_LIST_AVG	USE_LIST_AVG
<i>sort-key-index</i>	Specifies the key index (dimension), based on which the output will be sorted.	-1 or k , where $0 \leq k < \text{number_of_keys}$	0
<i>stack-raw</i>	Specifies to print the stack trace in RAW address format.	STKTRC_NO_SYM	0

When the sort-by flag is SORT_BY_KEY, SORT_BY_VALUE and the key (given by sort-key-index) and value pair is of a type where sorting cannot be done, then the *num-of-entries* option and other print options are applied to the printing of the individual key and value pair, if applicable. For example, if the sorting is by range type, the *num-of-entries* option, and the other print options are reserved for slots of each range.

The default associative array print options can be changed by using the **set_aso_print_options()** function in the BEGIN probe.

Example:

```
set_aso_print_options (10, SORT_TYPE_DESCEND|SORT_BY_VALUE);
```

As shown in the example, multiple flags can be provided by inserting a vertical bar symbol between them.

Note: The sort-key-index option cannot be set by the set_aso_print_options() function because it cannot be generalized for associative arrays of different key dimension sizes.

The print() function prints the keys and associated value for all elements or a subset of elements of the associative array by using the default print options. To override the default print options, you must use additional arguments in the print() function. For more information about the print() function, see the [Vue functions](#) topic.

The **print()** function prints the key and value pairs of the associative array by using the default print options. If you want to view the associative array contents in a different format, he will provide the *num-of-entries* option and the print option flags as additional parameters to the **print()** function.

Example:

```
/* uses default print options to display the contents of associative array 'count' */
print(count);
/* prints the first 10 entries of sorted associative array 'count'.
Default sort-by and sort-type options are used */
print(count, 10);
/* sorts the associative array 'count' in descending order of values and
displays the first 10 entries of 'count' */
print(count, 10, SORT_BY_VALUE|SORT_TYPE_DESCEND);

/* print elements which have first key as 0 */
print(var[0][ANY][ANY]);
```

- The clear() routine is used to clear the keys and associated values for elements of the associative array. The clear() routine is also used to reset the value of the associated array key without clearing the keys. If the clear() subroutine successfully clear one or more elements, it returns a 0, and the subroutine returns 1 when no elements are cleared.

```
clear(count); // count is an associative array.
```

The previous routine with only one argument of the associative array type clears all the key pairs present in the associative array **count**. After the previous clear operation, the associative array **count** is empty.

```
clear(count, RESET_VALUE); // count is an associative array.
clear(var); // var is an associative array with three key dimensions
```

The previous clear routine resets the value of all the key pairs in the associative array without clearing the key. The following default value is reset based on the type of value of the associative array:

To clear elements with specific keys, you must specify the keys in the first argument. Also, to ignore any particular key dimension (all values of that particular key dimension match), you can specify ANY. If the keys are specified, all the key dimensions of the associative array must be specified as either genuine values of the matching key type or ANY.

```
clear(var[ANY][“a”][ANY]); // clear all elements with second key as “a”
```

You can specify a second parameter in the `clear()` routine as `RESET_VALUE`. If you specify `RESET_VALUE`, the keys of the associative array are retained and only the values are reset.

```
clear(count, RESET_VALUE);
clear(var[0][ANY][ANY], RESET_VALUE);
```

The `RESET_VALUE` is dependent on type of the value. The following table displays the data types and default values to which the data type is reset:

Type	Default value
Integral types (int, long, short, long long)	0
LIST	Empty
Foat and double	0.0000000
String	Empty
stktrace_t	Empty
probev_timestamp_t	0
path_t	Empty
mac_addr_t	0
ip_addr_t	0

- The `Quantize` operation prints the keys and values of the given associative array in a graphical format based on the linear scaling of the values.

```
quantize(count);
```

`count` is an associative array, and it prints out the following contents:

```
key          value          =====
1            1              =====
2            2              =====
3            3              =====
4            4              =====
5            5              =====
6            6              =====
```

Similar to `print()` function, you can provide `quantize()` function print options to override the default print options.

Example:

```
/* sorts the associative array 'count' in descending order of values and displays
the first 10 entries of 'count' in graphical format*/
quantize(count, 10, _BY_VALUE|SORT_TYPE_DESCEND);
```

For associative arrays with multi-dimensional keys, keys can be specified in the first argument to restrict specific elements:

```
quantize(var[0][ANY][ANY]); //quantize elements with first key as 0
```

- Lquantize on associative array: This operation prints the keys and values of the given associative array in graphical format based on the logarithmic scaling of the values.

```
lquantize (count);
```

Where count is an associative array would print out the following contents :

key	value	
500	500	====
1000	1000	====
2000	2000	=====
4000	4000	=====
8000	8000	=====
16000	16000	=====
32000	32000	=====
64000	64000	=====

Similar to **print()** function, you can provide with **lquantize()** function print options to override the default print options.

Example:

```
/* sorts the associative array 'count' in descending order of values, and displays
the first 10 entries of 'count' in graphical
format based on the logarithmic value*/
lquantize(count, 10, _BY_VALUE|SORT_TYPE_DESCEND);
```

For associative arrays with multi-dimensional keys, keys can be provided in the first argument to restrict to specific elements:

```
lquantize(var[0][ANY][ANY]); //lquantize elements with first key as 0
```

The following example show how to use an associative array:

Example:

```
# Trace all the alloc- related calls and store the entry
# Time in 'entry_time' associative array
#
@@uft:$__CPID:*:"/alloc/":entry
{
    entry_time[get_function()]=timestamp();
}
#
# At exit, first check if entry for this function was traced
# If so, delete the entry time from 'entry_time' associative array
# To ensure that next time no action is taken on exit if entry was not traced.

@@uft:$__CPID:*:"/alloc/":exit
{
    func =get_function();
    if(exists(entry_time, func) )
    {
        append(time_taken[func],
            diff_time(timestamp(),entry_time[func],MICROSECONDS));
        delete(entry_time, func);
    }
}
#
# Print the list attributes sum, min, max, count, and avg time taken in every
# Alloc function.
#
@@syscall:$__CPID:exit:entry
{
    print(time_taken);
    exit();
}
```


Note: With this way, you need not define multiple list variables explicitly and still get the complete functions of the list with help of associative arrays.

Timestamp data type

A variable of type **probev_timestamp_t** holds the return value from the **timestamp** ProbeVue function, which returns a timestamp in internal AIX format.

The **probev_timestamp_t** type variable can be later passed as a parameter into the **diff_time** function that returns the difference between two timestamps. This data type can also be stored in an associative array either as a key or as a value.

Although the ProbeVue compiler does not type check when you use a **long** data type instead of the **probev_timestamp_t** data type for storing timestamps, if possible, avoid this usage.

The following operations are acceptable for a variable of the **probev_timestamp_t** type:

- Can be explicitly initialized to zero.

Note: A timestamp variable of global class or thread-local class is initialized to zero on start of the ProbeVue session.

- Can be compared against zero. Timestamp values returned by the **timestamp** function are always greater than zero.
- Can be compared against another timestamp variable. A later timestamp is guaranteed to be larger than an earlier timestamp.
- Can be passed as a parameter to the **diff_time** function.
- Can be printed using the **printf** or **trace** function.

File path data type

A variable of type **path_t** can be used to hold the value of **__file->path** (refer to **__file** built in for I/O probe manager) or function **fd_path()**. Only local or global variables of type **path_t** are supported. A variable of this type can also be the key or value in an associative Array.

Declaration of file path variable

```
path_t pth; // global variable of type path_t
auto:pth2 = fd_path(fd); // store in a local path_t variable
my_aso[__file->fname] = __file->path; // store in associative array
```

The signed, unsigned, register, static, thread, and kernel qualifiers are not supported for the **path_t** type variables.

Assignment operation

The assignment (=) operator allows one **path_t** variable to be assigned to another **path_t** variable. The original value of the variable is overwritten.

In the following example, after the assignment, **p1** and **p2** refer to the same file path:

```
path_t p1, p2;
p1 = fd_path(fd); // fd is a valid file descriptor value
p2 = p1;
```

Comparison operation

Only equality (==) and inequality (!=) operators are allowed for the **path_t** variables. The result of the equality operator is true (1) if both represent the same absolute file path or false (0) otherwise.

The inequality operator is the exact compliment of this behavior. No other comparison operators (>=, >, <, or <=) are allowed for the **path_t** type variables.

Printing file path type variables

A **path_t** type variable can be printed with “%p” format specifier in the **printf()** function.

This printing of file path involves a time-consuming file search operation in the corresponding file system. Hence, it must be judiciously used in Vue scripts.

Note: For transient files, which might no longer exist when the printf() message, can get printed; a null string is printed as file path.

An associative array that has the path_t type variables as key or value (or both) can be printed by using the print() function.

```
printf("file path=[%p]\n", __file->path);
my_aso[0] = fd_path(fd); // fd is valid file descriptor value
print(my_aso);
```

Limitations file path type variables

- Array of the path_t variables cannot be declared.
- The path_t variables cannot be used as members of struct or union.
- A pointer to the path_t variable is not allowed.
- Typcasting of the path_t variable to any other type or typcasting any other type to the path_t type is not allowed.
- No arithmetic operator (+, -, *, /, ++, --, and so on) can be used with the path_t type variable.

MAC address data type

A variable of type mac_addr_t is used to hold the value of MAC address. The MAC address data type is an abstract data type and you cannot use it directly with standard C unary or binary operators.

Only local or global variables of type mac_addr_t are supported.

A variable of this type can also be stored in an associate array either as a key or as a value.

Vue supports the following characteristics and operations for the MAC address type variables:

Declaration of MAC address variable

```
mac_addr_t
m1; //
    global variable of type
__auto mac_addr_t m2; // auto variable of type
m2 = __etherhdr->src_addr; // store source MAC address in a local
variable
mac_aso["src_mac_addr"] = __etherhdr->src_addr ; // store in an associative array.
```

The signed, unsigned, register, static, thread, and kernel qualifiers are not supported for the mac_addr_t type variables.

Assignment operation

The assignment (=) operator allows a mac_addr_t type variable to be assigned to another mac_addr_t type variable. The original values of the variable is overwritten. No typcasting is allowed from or to the mac_addr_t variable types.

In the following example, the content of the mac_addr_t m1 is assigned to m2.

```
mac_addr_t m1, m2; // Declares two MAC address variables.
m1 = __etherhdr->src_addr; // Get the source MAC address of the packet in m1.
m2 = m1 ; // Get the content of m1 into m2.
```

Comparison operation

Only equality (==) and inequality (!=) operators are allowed for the mac_addr_t variables. The result of the equality operator is True (1) if both contains the same MAC address values or False (0) otherwise.

The inequality operator is the exact complement of this behavior. No other comparison operators (\geq , $>$, $<$, or \leq) are allowed for the `mac_addr_t` type variables.

```
if ( m1 == m2) // comparing two mac_addr_t type variables.
printf("Mac addresses are equal"); else printf("Mac addresses are not equal");
```

Printing MAC address type variables

A `mac_addr_t` type variable can be printed with “%M” format specifier in the `printf()` function of `Vue`. An associative array that has `mac_addr_t` type variables as key or value (or both) can be printed by using the `print()` function.

```
printf(" Source MAC address=[%M]\n", __etherhdr->src_addr);
mac_aso["src_mac_address"] = __etherhdr->src_addr ; // Store source MAC address as value in an associative array mac_aso.
print(mac_aso);
```

Limitations for MAC address type variable

- The array of the `mac_addr_t` variable cannot be declared.
- The `mac_addr_t` variables cannot be used as a member of a struct or a union.
- Pointer to `mac_addr_t` variable is not allowed.
- Typecasting of `mac_addr_t` variable to any other type or typecasting any other type to `mac_addr_t` type is not allowed.
- No arithmetic operator (+, -, *, /, ++, -- etc) can be used with `mac_addr_t` type variable.

IP address data type

This is an abstract data type and cannot be used directly with standard C unary or binary operators. Only local or global variables of type `ip_addr_t` are supported. A variable of this type can also be stored in an associate array either as a key or as a value.

`Vue` supports the following characteristics and operations for the IP address type variables:

Declaration of IP address variable

```
ip_addr_t i1; // global variable of type ip_addr_t
__auto ip_addr_t i2; // auto variable of type
ip_addr_t i2 = __ip4hdr->src_addr; // store source IP address in a local ip_addr_t
variable.
ip_aso["src_ip_addr"] = __ip4hdr->src_addr; // store in an associative array.
```

The qualifiers `signed`, `unsigned`, `register`, `static`, `thread` and `kernel` are not supported for the `ip_addr_t` type variables.

Assignment operation

The assignment (`=`) operator allows a `ip_addr_t` type variable to be assigned to another `ip_addr_t` type variable and also it allows constant IP address or hostname to be assigned to `ip_addr_t` type variable. The original values of the variable is overwritten. No type casting is allowed from or to the `ip_addr_t` variable types.

In the following example, the content of the `ip_addr_t` `i1` is assigned to `i2`.

```
ip_addr_t i1, i2; // Declares two IP address variables.
ip_addr_t i3, i4, i5; // Declares three IP address variables.
i1 = __ip4hdr->src_addr; // Get the source IP address of the packet in i1.
i2 = i1 ; // Get the content of i1 into i2.
i3 = "10.10.10.1"; // Assign the constant IPv4 address to i3 variable.
i4 = "fe80::2c0c:33ff:fe48:f903"; // Assign the Ipv6 address to i4 variable.
i5 = "example.com"; // Assign the hostname to i5 variable.
// Get the content of i1 into i2.
```

Comparison operation

Only equality (`==`) and inequality (`!=`) operators are allowed for `ip_addr_t` types variables. The comparison allowed only between two `ip_addr_t` type variables and with constant string type (IP address or hostnames are provided in double quotes “192.168.1.1” or “example.com”).

The result of the equality operator is True (1) if both contains the same IP address type (IPV4 or IPV6) and values. or False (0) otherwise. The inequality operator is the exact compliment of that. No other comparison operators (>=, >, < or =<) are allowed for the ip_addr_t type variables.

```

if( i1 == i2)                // comparing two ip_addr_t type variables.
                                //IP address
string
printf("IP addresses are equal");
else printf("IP addresses are not equal");
or
if( i1 == "192.168.1.1")      // comparing ip_addr_t type variable and constant string.
    printf("IP addresses are equal");
else printf("IP addresses are not equal");
or
if( i1 = "example.com")      // comparing ip_addr_t type variable and constant
                                //IP address string
printf("IP addresses are equal");
else printf("IP addresses are not equal");

```

Printing IP address type variables

A ip_addr_t type variable can be printed with “%I” format specifier to print IP address in dotted decimal or hex format and “%H” format specifier to print hostname in the printf() function of Vue. This printing hostname involves a time consuming dns lookup operation. Hence it should be judiciously used in VUE scripts.

Note: When user uses the format specifier “%H” to print host name for IP address which may not exists in dns, for those IP addresses, it prints the IP addresses in dotted decimal/hex format instead of hostname.

An associative array that has ip_addr_t type variables as key or value (or both) can be printed using the print() function.

```

printf(" Source IP address=[%I]\n", __ip4hdr->src_addr);
ip_asof["src_ip_address"] = __ip4hdr->src_addr ; // Store source IP address as value in an associative array
print(ip_aso);

```

Limitations for IP address type variable

- The array of the ip_addr_t variable cannot be declared.
- Pointer to ip_addr_t variable is not allowed.
- Typecasting of ip_addr_t variable to any other type or typecasting any other type to ip_addr_t type is not allowed.
- No arithmetic operator (+, -, *, /, ++, -- etc) can be used with ip_addr_t type variable.

net_info_t data type

The net_info_t variable is structure or composite variable that is used to hold the network four tuples (local and remote IP addresses and port numbers) information from the specific socket descriptor through the sockfd_netinfo Vue function.

The members of net_info_t structure are accessed like any other user-defined structure in Vue script. The net_info_t type is an abstract data type and this variable cannot be used directly with standard C unary or binary operators. This variable is a structure containing 4 tuple information. This variable elements can be accessed by using the “.” operator like C structure elements.

Elements of the net_info_t data type are as follow:

```

net_info_t
{
    int local_port;
    int remote_port;
    ip_addr_t local_addr;
    ip_addr_t remote_addr;
};

```

Vue supports the following characteristics and operations for the net_info_t type variables:

Declaration of net_info_t type variable

```

net_info_t n1,n2
// n1 is variable of type net_info_t
sockfd_netinfo(fd, n1);
// fd is socket descriptor and n1 contains network
// four tuple information from sockfd_netinfo Vue function.

```

```
n2.local_addr = __ip4hdr->src_addr;
n2.remote_addr = __ip4hdr->dst_addr;
n1.local_port = __tcp_hdr->src_port;
n1.remote_port = __tcp_hdr->dst_port;
```

The signed, unsigned, register, static, thread, local, global, and kernel qualifiers are not supported for the `net_info_t` type variables.

Limitations for `net_info_t` type variable

- Structure and union member variable cannot be supported.
- Pointer to the `net_info_t` variable cannot be declared.
- This variable is not supported in associative array.
- The array of the `net_info_t` variable cannot be declared.
- Typecasting of the `net_info_t` variable to any other type or typecasting any other type to `net_info_t` type is not allowed.
- Arithmetic operator (+, -, *, /, ++, --, and so on) cannot be used with `net_info_t` type variable.

Vue library functions

Unlike programs written in C or in FORTRAN, or in a native language, scripts written in Vue do not have access to the subroutines provided by the AIX system libraries or any user libraries. However, Vue supports its own special internal library of functions useful for dynamic tracing programs.

Tracing-specific functions

get_function

Returns the name of the function that encloses the current probe. When the **get_function** function is called from interval, systrace, BEGIN, and END clause, the function returns an empty string.

timestamp

Returns the current timestamp.

diff_time

Finds the difference between two time stamps in microseconds or milliseconds.

Trace capture functions

printf

Formats and prints values of variables and expressions.

trace

Prints data without formatting.

stktrace

Formats and prints the stack trace.

List functions

list

Instantiates a list variable.

append

Appends a new item to list.

sum, max, min, avg, count

Aggregation functions that can be applied to a list variable.

C-library functions

atoi, strstr

Standard string functions.

Functions to support tentative tracing

start_tentative, end_tentative

Indicators for start and end of tentative tracing.

commit_tentative, discard_tentative

Commits or discards tentative trace data.

Miscellaneous functions

exit

Terminates the Vue script.

get_userstring

Reads string (or data) from user memory.

ptree

Prints the process tree of the probed process.

You can apply the Vue string functions only on variables of string type and not on a pointer variable. Standard string functions like **strcpy**, **strcat**, and so on are not necessary in Vue, because they are supported through the language syntax itself.

The ProbeVue compiler validates the data types of the parameters passed to Vue functions.

For **printf** function, validation is done to check whether there is an argument supplied in the **printf** function for each format specifier given in the format string. The total number of format specifiers and the total number of arguments passed to the **printf** function should be equal. In addition to this, validation is also done to match whether the type of the argument passed is compatible with the actual type mentioned as format specifier in the format string. If these checks fail, the ProbeVue throws an error message.

For example,

```
printf("hello world %s, %d\n", str);
```

would throw up an error message from the compiler as no argument is passed for **%d**. Similarly,

```
Printf("The total count of elements is %d\n", str);
```

also throws an error message as the format specified is **%d** whereas the argument passed, **str** variable is a string.

Other Features Functions

However, when given as

```
printf ("The total count of elements is %lld\n", i);
```

Where *i* is a variable of type int, no error message is thrown because the variable *i* is a compatible type for the format specifier requested. Hence, no exact type checking is done; however, compatible type checking is.

You cannot put functions in the predicate section of a Vue clause.

Predicates

You cannot use predicates when execution of clauses at probe points must be performed conditionally. The predicate section is identified by the presence of the **when** keyword immediately after the probe specification section. The predicate itself consists of regular C-style conditional expressions with the enclosing parentheses.

There are some restrictions on the expressions inside the predicate section:

- Kernel class variables are not permitted in the predicate.
- Automatic class variables are not permitted in the predicate.
- Floating-point type variables are not permitted in the predicate.
- Vue functions are not permitted inside a predicate.
- Side effects are not permitted inside a predicate and so the = assignment operator and its derivatives like +=, |=, and so on are not permitted.
- The ninth and higher parameters passed to a function (the entry class variables **__arg9**, **__arg10**, and so on) are not permitted in the predicate.

Conditional execution of specific actions within a clause is possible by using the **if ... else** statement which works like the analogous statement in C. However, if the entire clause is to be issued conditionally, it is preferable to use predicates instead because ProbeVue is designed to optimize execution of predicates.

Note: When a probe point can fire for more than one process, using thread-local variables inside the predicate is an excellent way to reduce the overall performance impact of enabling the probe. Putting conditional checks inside a predicate is preferable to using the **if** statement.

The following script uses thread-local variables inside predicates to efficiently detect when a particular character string is written to a specific file. It also shows an example of using the **if** statement within the action block of a clause with a predicate. Both the file name and character string are passed as parameters to the script using shell positional parameters.

```
/*
 * Filename : chkfilewrite.e
 *
 * Capture when someone writes a particular word to a specific file
 * takes 2 arguments: filename and word
 *
 * assumes file name is < 128
 *
 * Usage: probevue chkfilewrite.e \"<filename>\" \"<string>\"
 *
 * The backslashes above are necessary to prevent shell
 * from stripping the double quotation mark.
 */

int open(char *fname, int m, int p);
int write(int fd, char *s, int size);

@@syscall:*:open:entry
{
    __auto String fname[128];
    fname = get_userstring(__arg1, -1);
    if (fname == $1)
        thread:opening = 1;
}

@@syscall:*:open:exit
when (thread:opening == 1)
{
    thread:fd = __rv;
    thread:opening = 0;
}

@@syscall:*:write:entry
when (thread:fd == __arg1)
{
    __auto String buf[128];
    if (__arg3 < 128)
        buf = get_userstring(__arg2, __arg3);
    else
        buf = get_userstring(__arg2, 128);

    if (strstr(buf, $2)) {
        printf("%d wrote word to file.\n", __pid);
        exit();
    }
}
}
```

To run this program to check when someone writes the string "Error" to the **foo.log** file, you can issue the following command:

```
probevue chkfilewrite.e \"foo.log\" \"Error\"
```

Note: You can enhance the preceding script by adding a **close** probe to detect when the file is closed to prevent the script from catching the word after the original file is closed and a new one is opened and the same file descriptor number is reused.

Symbolic constants

Vue supports some pre-defined symbolic constants, which are commonly used in AIX programming. These constants are treated as keywords in Vue. During compilation, the constants are replaced by their definitions in system header files. Probe manager-specific symbolic constants are explained in their respective sections. Following are the generic symbolic constants.

AF_INET

This specifies the address family of type IPv4. This ensures that the data is of IPV4 type.

AF_INET6

This specifies the address family of type IPv6. This ensures that the data is of IPV6 type.

NULL

To set pointer types to a NULL or zero value. You cannot use NULL to set a String variable to the empty string.

Error numbers or "errno" names

These are the standard error names like **EPERM**, **EAGAIN**, **ESRCH**, **ENOENT**, and so on, specified by the POSIX and ANSI standards and defined in the `/usr/include/sys/errno.h` header file.

The following script traces when the **bind** system call fails with **errno** set to **EADDRINUSE** (address already in use).

```
/*
 * File: bind.e
 */

/*
 * Okay to use void for parameters since we are not planning to
 * access them in this script.
 */
int bind(void);

@@syscall:*:bind:exit
  when (__rv == -1)
{
  /*
   * The following check could also be moved to the predicate,
   * although it may not buy a lot because we are already in an
   * error path that should be executed only rarely
   */
  if (__errno == EADDRINUSE)
  /* This check could also be moved to the predicate */
    printf("%d failed with EADDRINUSE for bind() call.\n", __pid);
}
```

Signal names

These are the standard signal names like **SIGSEGV**, **SIGHUP**, **SIGILL**, **SIGABRT**, and so on, specified by the ANSI standards and defined in the `/usr/include/sys/signal.h` header file.

The following script shows how to debug "who" killed a particular process by sending it a specific signal.

```
/*
 * File: signal.e
 *
 * Who sent SIGKILL to my process ?
 */

/* Process IDs are < 2^32, so using an 'int' here instead of pid_t is
 * good enough
 */
int kill(int pid, int signo);

@@syscall:*:kill:entry
  when (__arg1 == $1 && __arg2 == SIGKILL)
{
  /* Trace sender of SIGKILL */
  printf("Stack trace of %s: (PID = %d)\n", __pname, __pid);
  stktrace(PRINT_SYMBOLS|GET_USER_TRACE, -1);
  exit();
}
```


FUNCTION_ENTRY

Identifies if a probe point is a function entry point. Used with the **get_location_point** function.

FUNCTION_EXIT

Identifies if a probe point is a function exit point. Used with the **get_location_point** function.

Header files

You can specify multiple header files on the command line either by separating the header files with a comma (with no spaces between the comma and the file names) or by specifying each one separately with the **-I** flag. The following two examples are equivalent:

```
probevue -I myheader.i,myheader2.i myscript.e
probevue -I myheader.i -I myheader2.i myscript.e
```

C++ header file can be included for struct/class definitions and allows a **probevue** script to access struct/class data fields through a pointer. All C++ header files can be listed using **#include** directives between **##C++** and **##Vue** directive in the ProbeVue script. For using this option IBM C++ compiler must be installed on the system. Another option to include C++ header file is to first preprocess the C++ header file with **-P** option of **probevue** and then include the preprocessed file with **-I** option of **probevue**. With **-P** option **probevue** will generate the out file with same name as input C++ header file with a **.Vue** suffix.

The advantage of using **-I** option for preprocessed C++ header file is that IBM C++ compiler need not be installed on the system.

You can run the following command to preprocess C++ header file.

```
probevue -P myheader.h
```

Note: To run the above command IBM C++ compiler is a pre-requisite.

The above command will generate a file called **myheader.Vue**. This file can be further shipped to another system and can be used to probe C++ application by including with **-I** option of **probevue**. While using the shipped preprocessed C++ header file the systems environment should be same for system being used for generating preprocessed C++ header file and system being used to include preprocessed header file with **-I** option of **probevue** for probing C++ application.

The C++ header file which is either being used for precompilation with **-P** option or being included between **##C++** and **##Vue** should have **.h** extension for including standard Input/Output C++ header file. For including IOstream header use **#include<iostream.h>** instead of **#include<iostream>**.

You can run the following command for a C++ executable named **cpp_executable** and the script named **myscript.e** to probe the C++ application.

```
probevue -I myheader.Vue -X cpp_executable myscript.e
```

Note: To run the above command IBM C++ compiler is a not a pre-requisite.

Supported shell elements

The Vue language syntax includes support for shell variables identified by the **\$** prefix like exported shell variables and positional parameters (arguments to the script).

Vue shell variables can appear anywhere in the Vue script. They can be part of the probe specification, be used in predicates or within statements in action blocks. However, unlike in a shell script, they are not expanded if used within double quoted strings.

The arguments passed from the command line to the script are referenced within the script as **\$1**, **\$2**, **\$3**, and so on. Consider the following Vue script:

```
/* Program name: myscript.e */
@@syscall:*:read:entry
  when (__pid == $1)
{
```

```

    int count;
    count++;
}

```

In the following example, the process ID of the process running the **myprog** program replaces **\$1** in the preceding script. It assumes that the **prgrep** shell program, which prints the process ID given the process name, is used to invoke the Vue script.

```

probevue myscript.e `prgrep myprog`

```

Environment variables exported from the shell can also be referenced in the script using the **\$** operator. Consider the following Vue script:

```

/* Program name: myscript2.e */
@@syscall:*:read:entry
  when (__pid == $PID)

{
  int count;
  count++;
}

/* program to be traced has a function called 'foo' */

@@uft:$PID:*:foo:entry
{
  printf("Read system call was invoked %d times\n", count);
}

```

In the following example, 3243 replaces \$PID in the preceding script:

```

PID=3243 probevue myscript2.e

```

If an environment variable needs to be recognized as a string inside the ProbeVue script, the value of the environment variable must include the enclosing double quotation mark that identify it as a string. For example, the following script captures trace output when a specific file is opened in the system:

```

/* Program name: stringshell.e */
int open(char *path, int oflag);
@@syscall:*:open:entry
{
  String s[40];
  s = get_userstring(__arg1, -1);
  if (s == $FILE_NAME) {
    printf("pid %d (uid %d) opened %s\n", __pid, __uid, s);
    exit();
  }
}

```

The script expects that \$FILE_NAME is the name of an exported shell environment variable, which includes the double quotation mark in its value. The following script is an example:

```

export FILE_NAME="/etc/passwd"
probevue stringshell.e

```

If the value of an existing environment variable that does not have double quotation mark is required in a script, then a new environment variable will need to be constructed using double quotation mark around the existing environment variable. The following script is an example:

```

export FILE_NAME="$HOME"
probevue stringshell.e

```

Vue supports two special environment variables which are useful when the process to be probed is started by the **probevue** command itself using the **-X** flag. The **\$_CPID** environment variable indicates the process ID of the child process created by the **probevue** command, and the **\$_CTID** environment variable indicates its thread ID. The **-X** flag is useful to probe short-lived processes especially for debugging purposes.

A Vue script can be executed directly (like a shell script) by setting the first line to the following script:

```
#!/usr/bin/probevue
```

The **probevue** command can also read the Vue script from standard input like the shell does. This can be accomplished by omitting the script filename from the command line. This is useful to test short scripts.

Vue does not support special shell parameters such as **\$\$** and **\$@** which are created internally by the shell.

Trace capture facilities

ProbeVue supports comprehensive trace capture facilities. The basic trace capture action is provided through the **printf** function that can be invoked from any probe as part of the action block. The Vue version of **printf** function is equipped with most of the power of the C library version. A second trace capture function is the **trace** function. The **trace** function accepts a single variable as a parameter and copies its value in printable hexadecimal format to the trace buffer. This function is particularly useful for dumping the contents of strings and structures. The **stktrace** function is another trace capture function which captures the stack trace of the traced thread at the current probe point.

In addition to the values of internal script variables, external variables like kernel global variables, context-specific data like parameters to the function being probed, return values from a function, and so on can also be captured and displayed through these trace capture functions.

The trace reporter always displays trace data in order of time of occurrence and thus the data captured from different CPUs are internally sorted before they are output.

Tentative tracing

Tentative tracing permits intelligent filtering of data and will reduce the actual amount of trace data that is presented to you and which you need to analyze. This has the excellent side effect of preventing buffer overflow problems if you can ensure that the tentatively collected data be discarded or committed early.

The following script is an example of using tentative tracing functions to capture trace data only when necessary:

```
/*
 * File: tentative.e
 *
 * Print details when write system call takes longer than a
 * specified number of microseconds
 *
 * Usage: probevue tentative.e <processID> <microseconds>
 */
int write(int fd, char *buf, int size);

@@BEGIN
{
    probev_timestamp_t ts1, ts2;
}

@@syscall:$1:write:entry
{
    __auto String buf[256];

    if (__arg3 < 256)
        buf = get_userstring(buf, __arg3);
    else
        buf = get_userstring(buf, 256);

    start_tentative("write");

    /* print out all the data associated with the write */
    stktrace(PRINT_SYMBOLS|GET_USER_TRACE, -1);

    printf("fd = %d, size = %d\n", __arg1, __arg3);

    /* Prints 256 bytes of buf, even though size may be < 256 */
    trace(buf);

    end_tentative("write");
}
```

```

    /* Get timestamp for when we entered write: do this at the end of
     * the probe to reduce probe effect
     */
    ts1 = timestamp();
}

/* If we started probing in the middle of write, ts1 will be zero,
 * ignore that case with a predicate
 */
@@syscall:$1:write:exit
when (ts1 != 0)
{
    /* diff_time() may return up to a 64-bit value, but we
     * use an int here since we don't expect the difference to
     * larger than a few hundred microseconds at the most.
     */
    int micros;

    /* Get timestamp for when we exited write: do this at the beginning of
     * the probe to reduce probe effect
     */
    ts2 = timestamp();

    micros = diff_time(ts1, ts2, MICROSECONDS);

    start_tentative("write");
    printf("Return value from write = %d\n", __rv);
    end_tentative("write");

    if (micros > $2) {
        /* Can mix normal trace with tentative also */
        printf("Time to write = %d, limit = %d micro seconds\n",
            micros, $2);
        commit_tentative("write");
        exit();
    }
    else
        discard_tentative("write");
}

```

Running ProbeVue

Dynamic tracing is only allowed for users with privileges or for the superuser.

Authorizations and privileges

This is unlike the static tracing facilities in AIX, which enforce relatively limited privilege checking. There is a reason for requiring privileges to run the **probevue** command. A Vue script can potentially produce more severe impacts on system performance than a static tracing facility like AIX system trace. This is because probe points for system trace are pre-defined and restricted. ProbeVue can potentially support many more probe points and the probe locations can potentially be defined almost anywhere. Further, ProbeVue trace actions at a probe point can take much longer to issue than the system trace actions at a probe point since those are limited to explicit data capture.

In addition, ProbeVue allows you to trace processes and read kernel global variables, both of which need to be controlled to prevent security exposures. A ProbeVue session can also consume a lot of pinned memory and restricting usage of ProbeVue to users with privilege reduces the risk of denial of service attacks. ProbeVue also allows administrators to control the memory usage of ProbeVue sessions through the SMIT interface.

Privileges for dynamic tracing are obtained differently depending upon whether role-based access control (RBAC) is enabled or not. Please refer to the AIX man pages for more information about enabling and disabling RBAC.

Note that in legacy or RBAC-disabled mode, there are no authorizations. Regular users cannot acquire privileges to run the **probevue** command to start a dynamic tracing session or run the **probevctrl** command to administer ProbeVue. Only the superuser can have privileges for both these functions. Do not disable RBAC when using ProbeVue unless you prefer to restrict this facility to root users only.

RBAC-enabled mode

Privileges in an RBAC system are obtained through authorizations. An authorization is a text string associated with security-related functions or commands. Authorizations provide the mechanism to grant rights to you to perform privileged actions. Only a user with sufficient authorization can issue the **probevue** command and start a dynamic tracing session.

aix.ras.probevue.trace.user.self

This authorization allows you to trace their applications in user space. The user ID of the process to be traced must be equal to the real user ID of the user invoking the **probevue** command.

This authorization allows you to enable probe points provided by the uft probe manager for your processes. However, the effective, real and saved user IDs of the process to be traced must be equal. Thus, you cannot trace `setuid` programs with just this authorization.

aix.ras.probevue.trace.user

This authorization allows you to trace any application in user space including `setuid` programs and applications started by the superuser. Be careful when handing out this authorization. This authorization allows you to issue the **probevue** command and enable probe points provided by the uft probe manager for any process on the system.

aix.ras.probevue.trace.syscall.self

This authorization allows you to trace system calls made by their applications. The effective, real and saved user IDs of the process making the system call must be the same and equal to the real user ID of the user invoking the **probevue** command. This authorization allows you to enable probe points provided by the syscall probe manager for your processes. The second field of the probe specification must indicate the process ID for a process started by you.

aix.ras.probevue.trace.syscall

This authorization allows you to trace system calls made by any application on the system including `setuid` programs and applications started by the superuser. Be careful when handing out this authorization. This authorization allows you to issue the **probevue** command and enable probe points provided by the syscall probe manager for any process. The second field of the probe specification can either be set to a process ID to probe a specific process or to `*` to probe all processes.

aix.ras.probevue.trace

This authorization allows you to trace the entire system and includes all the authorizations defined in the preceding sections. You can also access and read kernel variables when running the **probevue** command, trace system trace events by using the `systrace` probe manager and trace the CPU bound probes by using the interval probe manager. Be careful while using this authorization.

aix.ras.probevue.manage

This authorization allows you to administer ProbeVue. This includes changing the values of the different ProbeVue parameters, starting or stopping ProbeVue and viewing details of dynamic tracing sessions of all users when running the **probevctrl** command. Without this authorization, you can use the **probevctrl** command to view session data for dynamic tracing sessions started by you or view the current values for ProbeVue parameters.

aix.ras.probevue.rase

This authorization allows you to access to a highly privileged set of "RAS events" Vue functions which can produce system and LMT trace records, create live dumps, and even lead to the system `abend`. This privilege must be very carefully controlled.

aix.ras.probevue

This authorization grants all dynamic tracing privileges and is equivalent to all the preceding authorizations combined.

The superuser (or root) has all these authorizations assigned by default. Other users will need to have authorizations assigned to them by first creating a role with a set of authorizations and assigning the role to the user. The user will also need to switch roles to a role that has the required authorizations defined for dynamic tracing before invoking the **probevue** command. The following script is an example of how to provide user "joe" authorization to enable user space and system call probes for processes started by "joe".

```
mkrole authorizations=  
"aix.ras.probevue.trace.user.self,aix.ras.probevue.trace.syscall.self"  
appttrace
```

```
chuser roles=apptrace joe
setkst -t roleTR
```

ng command:

```
swrole apptrace
```

Note: The interval probe manager does not have a specific authorization associated with it. You can enable interval probe points if you have any of the **aix.ras.probevue.trace*** authorizations.

ProbeVue privileges

The privileges that are available for ProbeVue are listed in the following table. A description of each privilege and the authorizations that map to that privilege is provided. Privileges form a hierarchy where the parent privilege contains all of the rights that are associated with the privileges of its children, but it can include additional privileges also.

Privilege	Description	Authorizations	Associated command
PV_PROBEVUE_TRC_USER_SELF	Allows a process to enable dynamic user space probe points on another process with the same real user ID.	aix.ras.probevue.trace.user.self aix.ras.probevue.trace.user aix.ras.probevue.trace aix.ras.probevue	probevue
PV_PROBEVUE_TRC_USER	Allows a process to enable dynamic user space probe points. Includes the PV_PROBEVUE_TRC_USER_SELF privilege.	aix.ras.probevue.trace.user aix.ras.probevue.trace aix.ras.probevue	probevue
PV_PROBEVUE_TRC_SYSCALL_SELF	Allows a process to enable dynamic system call probe points on another process with the same real user ID.	aix.ras.probevue.trace.syscall.self aix.ras.probevue.trace.syscall aix.ras.probevue.trace aix.ras.probevue	probevue
PV_PROBEVUE_TRC_SYSCALL	Allows a process to enable dynamic system call space probe points. Includes the PV_PROBEVUE_TRC_SYSCALL_SELF privilege.	aix.ras.probevue.trace.syscall aix.ras.probevue.trace aix.ras.probevue	probevue
PV_PROBEVUE_TRC_KERNEL	Allows a process to access kernel data when dynamic tracing.	aix.ras.probevue.trace aix.ras.probevue	probevue

Table 23. ProbeVue privileges (continued)

Privilege	Description	Authorizations	Associated command
PV_PROBEVUE_MANAGE	Allows a process to administer ProbeVue.	aix.ras.probevue.manage aix.ras.probevue	probevctrl
PV_PROBEVUE_RASE	Authorizes the use of the restricted "RAS events" functions.	aix.ras.probevue.rase aix.ras.probevue	probevue
PV_PROBEVUE_*	Equivalent to all the preceding privileges (PV_PROBEVUE_*) combined.	aix.ras.probevue	probevue probevctrl

ProbeVue parameters

AIX provides a set of parameters that you can use to tune ProbeVue or the ProbeVue framework. The parameters allow you to specify both global limits on resource usage by the ProbeVue framework and to specify resource usage for individual users.

Note: Probe managers are not contained within the ProbeVue framework and hence these limits do not apply them.

All ProbeVue parameters can be modified through the SMIT interface (use the "smit probevue" fast path) or directly through the **probevctrl** command. ProbeVue can be stopped if there are no active dynamic tracing sessions and it can be restarted after stopping it without requiring a reboot. ProbeVue can fail to stop if any sessions that used thread-local variables had been previously active.

The following table summarizes the parameters defined for dynamic tracing sessions. In the description, a privileged user refers to the superuser or a user with the **aix.ras.probevue.trace** authorization and a non-privileged user is one who does not have this authorization.

Table 24. Parameters for dynamic tracing session

Description as in SMIT	Maximum value	Initial high configuration value	Initial low configuration value	Minimum value	Associated command
MAX pinned memory for ProbeVue framework	64 GB	10% of available memory or the maximum value, whichever is smaller.	16 MB	3 MB	<p>Maximum pinned memory in MB that is allocated for ProbeVue data structures, including per-CPU stacks and per-CPU local table regions and by all dynamic tracing sessions. It does not include any memory allocated by Probe Managers.</p> <p>Note: Although, this parameter can be modified at any time, the value takes effect only the next time ProbeVue is started.</p>
Default per-CPU trace buffer size	256 MB	128 KB	8 KB	4 KB	<p>Default size in KB of per-CPU trace buffer. Two trace buffers are allocated per CPU for each dynamic tracing session by ProbeVue, one active and used by the writer or the Vue program when it captures trace data and one inactive and used by the reader or the trace consumer.</p>

Table 24. Parameters for dynamic tracing session (continued)

Description as in SMIT	Maximum value	Initial high configuration value	Initial low configuration value	Minimum value	Associated command
					<p>For example, on an 8-way with per-CPU trace buffer size set to 16 KB, the total memory consumed by the trace buffers for a ProbeVue session is 256 KB. You can specify a different buffer size (larger or smaller) when you start the probevue command until it is within the session memory limits.</p>
MAX pinned memory for regular user sessions	64 GB	2 MB	2 MB	0 MB	<p>Maximum pinned memory allocated for a non-privileged user ProbeVue session including memory for the per-CPU trace buffers. A value of 0 effectively disables all non-privileged users. Privileged users have no limits on the memory used by their ProbeVue sessions. However, they are still limited by the maximum pinned memory allowed for the ProbeVue framework.</p>

Table 24. Parameters for dynamic tracing session (continued)

Description as in SMIT	Maximum value	Initial high configuration value	Initial low configuration value	Minimum value	Associated command
MIN trace buffer read rate for regular user	5000 ms	100 ms	100 ms	10 ms	The minimum period, in milliseconds, that a non-privileged user can request the trace consumer to check for trace data. This value is internally rounded to the next highest multiple of 10 milliseconds. Privileged users are not limited by this parameter, but the fastest read rate that they can specify is 10 milliseconds.
Default trace buffer read rate	5000 ms	100 ms	100 ms	10 ms	The default period in milliseconds that the in-memory trace buffers are checked for trace data by the trace consumer. You can specify a different read rate (larger or smaller) when starting the probevue command until it is larger than the minimum buffer read rate.

Table 24. Parameters for dynamic tracing session (continued)

Description as in SMIT	Maximum value	Initial high configuration value	Initial low configuration value	Minimum value	Associated command
MAX concurrent sessions for regular user	8	1	1	0	Number of concurrent ProbeVue sessions allowed for a non-privileged user. A value of zero effectively disables all non-privileged users.
Size of per-CPU computation stack	256 KB	20 KB	12 KB	8 KB	The size of the per-CPU computation stack used by ProbeVue when issuing the Vue script. The value is rounded to the next highest multiple of 8 KB. ProbeVue allocates a single stack per-CPU for all ProbeVue sessions. The memory consumed for the stacks is not included in the per-session limits.

Table 24. Parameters for dynamic tracing session (continued)

Description as in SMIT	Maximum value	Initial high configuration value	Initial low configuration value	Minimum value	Associated command
					<p>Note: Although, this parameter can be modified at any time, the value takes effect only after AIX kernel boot image is rebuilt and rebooted. You have to configure ProbeVue stack to use 96K virtual memory to get the current directory listing.</p>
Size of per-CPU local table size	256 KB	32 KB	4 KB	4 KB	<p>The size of the per-CPU local table used by ProbeVue for saving variables of automatic class and for saving temporary variables. ProbeVue uses half of this area for automatic variables and the remaining half for saving temporary variables.</p>

Table 24. Parameters for dynamic tracing session (continued)

Description as in SMIT	Maximum value	Initial high configuration value	Initial low configuration value	Minimum value	Associated command
					<p>The value is always rounded to the next highest multiple of 4 KB. ProbeVue allocates a single local table and a single temporary table per-CPU used by all ProbeVue sessions. The memory consumed for the local tables is not included in the per-session limits.</p> <p>Note: Although, this parameter can be modified at any time, the value takes effect only the next time ProbeVue is started.</p>
MIN interval allowed in an interval probe	N/A	1		1	Minimum timer interval, in milliseconds, allowed for global root user in interval probes.

Table 24. Parameters for dynamic tracing session (continued)

Description as in SMIT	Maximum value	Initial high configuration value	Initial low configuration value	Minimum value	Associated command
Number of threads to be traced	N/A	32	32	1	Maximum number of threads that a ProbeVue session can support when it has thread-local variables. The ProbeVue framework allocates the thread-local variables to the maximum number of threads that are specified with this attribute, at the start of the session. If more than the specified number of threads hit the probe that has a thread-local variable, the ProbeVue session is abruptly stopped.

Table 24. Parameters for dynamic tracing session (continued)

Description as in SMIT	Maximum value	Initial high configuration value	Initial low configuration value	Minimum value	Associated command
Number of page faults to be handled	1024	0	0	0	Number of page fault contexts for handling page faults for the entire framework. A page fault context includes stack and local table for saving automatic class variables and temporary variables. A page fault context is required to access the paged-out data. If there are no page fault context that is free at the time of a page fault, ProbeVue does not fetch the paged-out data.
Maximum probe execution time for systrace probes when fired in interrupt context	N/A	0	0	0	This number limits the maximum time, in milliseconds, a systrace probe executing in interrupt context can take. By default, the value is zero, which means the systrace probe can any time.

Table 24. Parameters for dynamic tracing session (continued)

Description as in SMIT	Maximum value	Initial high configuration value	Initial low configuration value	Minimum value	Associated command
Maximum probe execution time for io probes when fired in interrupt context	N/A	0	0	0	This number limits the maximum time, in milliseconds, an io probe executing in interrupt context can take. By default, the value is zero, which means it can any time
Maximum probe execution time for sysproc probes when fired in interrupt context	N/A	0	0	0	This number limits the maximum time, in milliseconds, a sysproc probe executing in interrupt context can take. By default, the value is zero, which means it can any time.
Maximum probe execution time for network probes when fired in interrupt context	N/A	0	0	0	This number limits the maximum time, in milliseconds, a network probe executing in interrupt context can take. By default, the value is zero, which means it can any time.

Table 24. Parameters for dynamic tracing session (continued)

Description as in SMIT	Maximum value	Initial high configuration value	Initial low configuration value	Minimum value	Associated command
Max network buffer size	64 KB	64 bytes	96 bytes	96 bytes	This value is a pre-allocated buffer size (in bytes) used by network probe manager for bpf probe points. This value is allocated when the first bpf probe is enabled and exists in the system till the last bpf probe is disabled. When the last bpf probe type is disabled, this buffer is released. This buffer is used to copy the data when packet data is spanned across multiple packet buffers.
Asynchronous statistics fetch interval, in milliseconds	NA	1000 milliseconds (1 second)	1000 milliseconds (1 second)	100 milliseconds	The interval, in milliseconds, to fetch the asynchronous statistics. This value is global and is applicable to all ProbeVue sessions.
Fetch statistics in asynchronous mode only	NA	No	No	NA	Specifies ProbeVue that statistics must be fetched in an asynchronous mode even if synchronous mode is available.

Table 24. Parameters for dynamic tracing session (continued)

Description as in SMIT	Maximum value	Initial high configuration value	Initial low configuration value	Minimum value	Associated command
Maximum probe execution time for CPU bound interval probes when fired in interrupt context.	60 seconds	60 seconds	100 milliseconds	100 milliseconds	This number limits the maximum time, in milliseconds, a CPU bound interval probe executing in interrupt context can take. By default, the value is 60secs.

Profiling ProbeVue Session

The ProbeVue framework provides a profiling facility that can be turned on or off to estimate the impact of enabled probes on the application. This facility accumulates the time taken by probe actions when they are started and reports when requested or when the session ends.

The profiling report displays the probe string and the time taken by the action corresponding to that probe string. The time that is consumed by the probe action is maintained as a list where the data collected is total, minimum, maximum, and average time taken by probe action. Profiling data also displays number of times that the probe action was timed. When you are looking up the profile for multiple functions through one probe string (by using regular expression or * in place of function name), profiling data provides an accumulated data of probes started for all such functions. It does not provide timing details for functions that are probed separately but only per-probe action.

The BEGIN and END probe actions are not profiled with this facility. These profiling details are session-specific details. You can enable probevue session profiling along with session start by using the **probevue** command or **probevctrl** command.

For more information, see the [probevue](#) and [probevctrl](#) commands.

Sample programs

Example 1

The following canonical "Hello World" program prints "Hello World" into the trace buffer and exits:

```
#!/usr/bin/probevue

/* Hello World in probevue */
/* Program name: hello.e */

@@BEGIN
{
    printf("Hello World\n");
    exit();
}
```

Example 2

The following "Hello World" program prints "Hello World" when you types Ctrl-C on the keyboard:

```
#!/usr/bin/probevue

/* Hello World 2 in probevue */
/* Program name: hello2.e */

@@END
{
```

```

    printf("Hello World\n");
}

```

Example 3

The following program shows how to use thread-local variables. This Vue script counts the number of bytes written to a particular file. It assumes that the processes are single-threaded or those threads that open files are the same ones that write to them. It also assumes that all write operations are successful. The script can be terminated at any time and you can obtain the current count of bytes written by typing Ctrl-C on the terminal.

```

#!/usr/bin/probevue

/* Program name: countbytes.e */
int open( char * Path, int OFlag, int mode );
int write( int fd, char * buf, int sz);
int done;

@@syscall:*.open:entry
  when (done != 0 )
{
  if (get_userstring(__arg1, -1) == "/tmp/foo") {
    thread:trace = 1;
    done = 1;
  }
}

@@syscall:*.open:exit
  when (thread:trace)
{
  thread:fd = __rv;
}

@@syscall:*.write:entry
  when (thread:trace && __arg1 == thread:fd)
{
  bytes += __arg3; /* number of bytes is third arg */
}

@@END
{
  printf("Bytes written = %d\n", bytes);
}

```

Example 4

The following tentative tracing program shows how to trace the arguments passed to the read system call only if it returns zero bytes when reading the **foo.data** file:

```

#!/usr/bin/probevue
/* File: ttrace.e */
/* Example of tentative tracing */
/* Capture parameters to read system call only if read fails */
int open ( char* Path, int OFlag , int mode );
int read ( int fd, char * buf, int sz);

@@syscall:*.open:entry
{
  filename = get_userstring(__arg1, -1);
  if (filename == "foo.data") {
    thread:open = 1;
    start_tentative("read");
    printf("File foo.data opened\n");
  }
}

@@syscall:*.open:exit
  when (thread:open == 1)
{
  thread:fd = __rv;
  start_tentative("read");
  printf("fd = %d\n", thread:fd);
  thread:open = 0;
}

@@syscall:*.read:entry
  when (__arg1 == thread:fd)

```

```

{
    start_tentative("read");
    printf("Read fd = %d, input buffer = 0x%08x, bytes = %d,",
        __arg1, __arg2, __arg3);
    end_tentative("read");
    thread:read = 1;
}

@@syscall:*.read:exit
    when (thread:read == 1)
{
    if (__rv < 0) {
        /* The printf below, even though non-tentative, is only
        * executed in error cases and merges with the
        * previously printed tentative data
        */
        printf(" errno = %d\n", __errno);
        commit_tentative("read");
    }
    else
        discard_tentative("read");
    thread:read = 0;
}

```

A possible output if the read failed because a bad address (say 0x1000) was passed as input buffer pointer could look like the following output:

```

#probevue ttrace.e
File foo.data opened
fd = 4
Read fd = 4, input buffer = 0x00001000, bytes = 256, errno = 14

```

Example 5

The following Vue script prints the values of some kernel variables and exits immediately. Pay attention to the **exit** function in the **@@BEGIN** probe:

```

/* File: kernel.e */
/* Example of accessing kernel variables */
/* System configuration structure from /usr/include/sys/systemcfg.h */
struct system_configuration {
    int architecture; /* processor architecture */
    int implementation; /* processor implementation */
    int version; /* processor version */
    int width; /* width (32 || 64) */
    int ncpu; /* 1 = UP, n = n-way MP */
    int cache_attrib; /* L1 cache attributes (bit flags) */
        /* bit 0/1 meaning */
        /* ----- */
        /* 31 no cache / cache present */
        /* 30 separate I and D / combined */
    int icache_size; /* size of L1 instruction cache */
    int dcache_size; /* size of L1 data cache */
    int icache_asc; /* L1 instruction cache associativity */
    int dcache_asc; /* L1 data cache associativity */
    int icache_block; /* L1 instruction cache block size */
    int dcache_block; /* L1 data cache block size */
    int icache_line; /* L1 instruction cache line size */
    int dcache_line; /* L1 data cache line size */
    int L2_cache_size; /* size of L2 cache, 0 = No L2 cache */
    int L2_cache_asc; /* L2 cache associativity */
    int tlb_attrib; /* TLB attributes (bit flags) */
        /* bit 0/1 meaning */
        /* ----- */
        /* 31 no TLB / TLB present */
        /* 30 separate I and D / combined */
    int itlb_size; /* entries in instruction TLB */
    int dtlb_size; /* entries in data TLB */
    int itlb_asc; /* instruction tlb associativity */
    int dtlb_asc; /* data tlb associativity */
    int resv_size; /* size of reservation */
    int priv_lck_cnt; /* spin lock count in supevisor mode */
    int prob_lck_cnt; /* spin lock count in problem state */
    int rtc_type; /* RTC type */
    int virt_alias; /* 1 if hardware aliasing is supported */
    int cach_cong; /* number of page bits for cache synonym */
    int model_arch; /* used by system for model determination */
    int model_impl; /* used by system for model determination */
}

```

```

int Xint;          /* used by system for time base conversion */
int Xfrac;        /* used by system for time base conversion */
int kernel;       /* kernel attributes */
/* bit          0/1 meaning */
/* -----*/
/* 31    32-bit kernel / 64-bit kernel */
/*          /* 30    non-LPAR / LPAR */
/*          /* 29    old 64bit ABI / 64bit Large ABI */
/*          /* 28    non-NUMA / NUMA */
/*          /* 27    UP / MP */
/*          /* 26    no DR CPU add / DR CPU add support */
/*          /* 25    no DR CPU rm / DR CPU rm support */
/*          /* 24    no DR MEM add / DR MEM add support */
/*          /* 23    no DR MEM rm / DR MEM rm support */
/*          /* 22    kernel keys disabled / enabled */
/*          /* 21    no recovery / recovery enabled */
/*          /* 20    non-MLS / MLS enabled */
long long physmem; /* bytes of OS available memory */
int slb_attr;      /* SLB attributes */
/* bit          0/1 meaning */
/* -----*/
/* 31    Software Managed */
int slb_size;     /* size of slb (0 = no slb) */
int original_ncpus; /* original number of CPUs */
int max_ncpus;    /* max cpus supported by this AIX image */
long long maxrealaddr; /* max supported real memory address +1 */
long long original_entitled_capacity;
/* configured entitled processor capacity */
/* at boot required by cross-partition LPAR */
/* tools. */
long long entitled_capacity; /* entitled processor capacity */
long long dispatch_wheel; /* Dispatch wheel time period (TB units) */
int capacity_increment; /* delta by which capacity can change */
int variable_capacity_weight; /* priority weight for idle capacity*/
/* distribution */
int splpar_status; /* State of SPLPAR enablement */
/* 0x1 => 1=SPLPAR capable; 0=not */
/* 0x2 => SPLPAR enabled 0=dedicated; */
/* 1=shared */
int smt_status; /* State of SMT enablement */
/* 0x1 = SMT Capable 0=no/1=yes */
/* 0x2 = SMT Enabled 0=no/1=yes */
/* 0x4 = SMT threads bound true 0=no/1=yes */
int smt_threads; /* Number of SMT Threads per Physical CPU */
int vmx_version; /* RPA defined VMX version, 0=none/disabled */
long long sys_lmbsize; /* Size of an LMB on this system. */
int num_xcpus; /* Number of exclusive cpus on line */
signed char errchecklevel; /* Kernel error checking level */
char pad[3]; /* pad to word boundary */
int dfp_version; /* RPA defined DFP version, 0=none/disabled */
/* if MSbit is set, DFP is emulated */
};

__kernel struct system_configuration _system_configuration;

@@BEGIN
{
String s[40];
int j;
__kernel int max_sdl; /* Atomic RAD system decomposition level */
__kernel long lbolt; /* Ticks since boot */

printf("No. of online CPUs\t\t= %d\n", _system_configuration.ncpus);

/* Print SMT status */
printf("SMT status\t\t=");
if (_system_configuration.smt_status == 0)
printf(" None");
else {
if (_system_configuration.smt_status & 0x01)
printf(" Capable");
if (_system_configuration.smt_status & 0x02)
printf(" Enabled");
if (_system_configuration.smt_status & 0x04)
printf(" BoundThreads");
}
printf("\n");

/* Print error checking level */
if (_system_configuration.errchecklevel == 1)
s = "Minimal";
else if (_system_configuration.errchecklevel == 3)

```

```

    s = "Normal";
else if (_system_configuration.errchecklevel == 7)
    s = "Detail";
else if (_system_configuration.errchecklevel == 9)
    s = "Maximal";
printf("Error checking level\t\t= %s\n",s);

printf("Atomic RAD system detail level\t= %d\n", max_sdl);

/* Long in the kernel is 64-bit, so we use %lld below */
printf("Number of ticks since boot\t= %lld\n", lbolt);

    exit();
}

```

The following output is a possible output when you run the preceding script on a Power 5 dedicated partition with default kernel attributes:

```

# probevue kernel.e
No. of online CPUs           = 4
SMT status                   = Capable Enabled BoundThreads
Error checking level         = Normal
Atomic RAD system detail level = 2
Number of ticks since boot   = 34855934

```

Probe managers

Probe managers are not part of the basic ProbeVue framework. Probe managers enable probe points that can be used by ProbeVue for dynamic tracing.

Probe managers generally support a set of probe points that belong to some common domain and share some common feature or attribute that distinguishes them from other probe points. Probe points are useful at points where control flow changes significantly, at points of state change or at other points of significant interest. Probe managers are careful to select probe points only in locations that are safe to instrument.

Probe managers can choose to define their own distinct rules for the probe specifications within the common style that must be followed for all probe specifications.

ProbeVue supports the following probe managers:

System call probe manager

The syscall probe manager supports probes at the entry and exit of well-defined and documented base AIX system calls. These are the system calls that have the same interface at the **libc.a** (or C library) entry point and in the kernel entry point. Either the system call is a pass-through (the C library simply imports the symbol from the kernel and the exports it with no code in the library) or there is trivial code for the interface inside the library.

The syscall probe manager accepts a 4-tuple probe specification in one of the following formats:

- `syscall:*:<system_call_name>:entry`
- `syscall:*:<system_call_name>:exit`

where the `system_call_name` field is to be substituted by the actual system call name. These indicate that a probe be placed at the entry and exit of system calls. Assigning the `*` to the second field indicates that the probe will be fired for all processes.

Note: Different privileges are required for enabling system call probes. Probing every process in the system requires higher privileges than probing your own processes.

Additionally, the syscall probe manager also accepts a 4-tuple probe specification in one of the following formats:

- `syscall:<process_ID>:<system_call_name>:entry`

- `syscall:<process_ID>:<system_call_name>:exit`

where a process ID can be specified as the second field of the probe specification to support probing of specific processes.

The system call names accepted by the syscall probe manager are the names of the *libc.a* interfaces and not the kernel's internal system call names. For example, the **read** subroutine is exported by *libc.a*, but the actual system call name or kernel entry point is **kread**. The syscall probe manager will internally translate a *libc* interface to its kernel entry point and enable the probe at entry into the **kread** kernel routine. Because of this, if multiple C library interfaces invoke the **kread** routine, the probe pointfires for those interfaces also. Generally, this is not a problem because for most of the system calls supported by the syscall probe manager, there is a 1-to-1 mapping between the *libc* interface and the kernel routine.

For each syscall probe, there is an equivalent probe point in the library code provided by the uft probe manager. The uft probe manager does support all library interfaces (unless it is a passthrough interface and there is no code for the call or references to it in the library at all) including those not supported by the syscall probe manager. However, the syscall probe manager has two advantages:

- The syscall probe manager can probe every process in the system by specifying asterisk as the second field.
- The syscall probe manager is more efficient than the uft probe manager because it does not need to switch from user mode to kernel mode and back to run the probe actions.

For more information about the full list of system calls supported by the syscall probe manager see [ProbeVue](#).

UFT probe manager

The uft or the user function tracing probe manager supports probing user space functions that are visible in the XCOFF symbol table of a process. The uft probe manager supports probe points that are at entry and exit points of functions whose source is a C or FORTRAN language text file even though the symbol table can contain symbols whose sources are from a language other than C or FORTRAN.

The tracing of Java™ applications in a way identical to the existing tracing mechanism from the users point of view and the JVM is one that performs most of the real tasks on behalf of Probevue.

The uft probe manager accepts a 5-tuple probe specification in the following format:

```
uft:<processID>:*:<function_name>:<entry|exit>
```

Note: The uft probe manager requires the process ID for the process to be traced and the complete function name of the function at whose entry or exit point the probe is to be placed.

PID or program name

Specifies process ID of the process to be traced or the name of the program (name of the executable file) to be probed. You can provide the process ID or the name of the executable file in double quotes ("`/usr/bin/test`") or in single quotes ('`/usr/bin/test`') or as a file path (`/usr/bin/test`). You can provide the absolute path or the relative path of the executable file. You can also provide the hard link or the symbolic link to the executable file as a value to the program name parameter. The following examples show different ways to specify the second tuple probe specification.

Examples

```
@@uft:123450 :* :foo :entry
@@uft : /usr/bin/test:* :foo:entry
@@uft:./test:*:foo:entry
@@uftxlc++:/usr/bin/test:*:foo:exit
@@uftxlc++:"/usr/bin/testxxx":*:foo:exit
@@uft : "/usr/bin/xxx":* :foo:entry
@@uft: '/home/xxx/test':*:func1:entry
```

If `/usr/bin/test` is the executable file and `/testln` is the soft link to the `/usr/bin/test` file (`/testln ->/usr/bin/test`), run the following command to probe all processes that were started by running the `/usr/bin/test` executable file or the `/testln` soft link for that user.

```
@@uft:/testln:*:foo:entry
```

Note: At a time a user can trace an executable by using one session only.

When the third field is set to `*`, the UFT probe manager searches the function in all of the modules loaded into the process address space including the main executable and shared modules. This implies that if a program contains more than one C function with this name (for example, functions with static class that are contained in different object modules), then probes will be applied to the entry point of every one of these functions.

If a function name in a specific module needs to be probed, the module name needs to be specified in the third field. The probe specification syntax to provide the library module name is illustrated below:

```
# Function foo in any module
@@uft:<pid>:*:foo:entry
# Function foo in any module in any archive named libc.a
@@uft:<pid>:libc.a:foo:entry
# Function foo in the shr.o module in any archive named libc.a
@@uft:<pid>:libc.a(shr.o):foo:entry
```

The function name in the fourth tuple can be specified as an Extended Regular Expression (ERE). The ERE should be enclosed between `/` and `/` like `"/<ERE>/"`.

When the function name is specified as an ERE, all the functions matching the specified regular expression in the specified module is probed.

```
/* Probe entry of all libc.a functions starting with "malloc" word */
@@uft:$__CPID:libc.a:"/^malloc.*"/:entry
/* Probe exit of all functions in the executable a.out */
@@uft:$__CPID:a.out:"/.*"/:exit
```

In the entry probes, where a function name is specified as a regular expression, individual arguments cannot be accessed. However, probevue function **print_args** can be used to print the function name and its arguments. The argument values is printed based on the argument type information available in the traceback table of the function.

In the exit probes, where a function name is specified as a regular expression, return value cannot be accessed.

Probevue supports enabling probes in more than one process at the same time. However, you will need privileges even for probing processes that belong to you.

Probevue enforces a restriction that prevents processes with user-space probes from being debugged using the **ptrace** or **procfs** based APIs.

As indicated above, the uft probe manager supports probes in shared modules like shared library modules. The following script shows an example that traces mutex activity by enabling probes in the thread library's mutex lock and unlock subroutines.

```
/* pthreadlocks.e */
/* Trace pthread mutex activity for a given multithreaded process */
/* The following defines are from /usr/include/sys/types.h */

typedef long long pid_t;
typedef long long thread_t;

typedef struct {
    int    __pt_mutexattr_status;
    int    __pt_mutexattr_pshared;
    int    __pt_mutexattr_type;
} pthread_mutexattr_t;

typedef struct __thrq_elt thrq_elt_t;
```



```

struct __thrq_elt {
    thrq_elt_t    *__thrq_next;
    thrq_elt_t    *__thrq_prev;
};

typedef volatile unsigned char _simplelock_t;

typedef struct __lwp_mutex {
    char          __wanted;
    _simplelock_t  __lock;
} lwp_mutex_t;

typedef struct {
    lwp_mutex_t    __m_lmutex;
    lwp_mutex_t    __m_sync_lock;
    int            __m_type;
    thrq_elt_t     __m_sleepq;
    int            __filler[2];
} mutex_t;

typedef struct {
    mutex_t        __pt_mutex_mutex;
    pid_t          __pt_mutex_pid;
    thread_t       __pt_mutex_owner;
    int            __pt_mutex_depth;
    pthread_mutexattr_t  __pt_mutex_attr;
} pthread_mutex_t;

int pthread_mutex_lock(pthread_mutex_t *mutex);
int pthread_mutex_unlock(pthread_mutex_t *mutex);

@@uft:$ _CPID:*:pthread_mutex_lock:entry
{
    printf("thread %d: mutex 0x%08x locked\n", __tid, __arg1);
}

@@uft:$ _CPID:*:pthread_mutex_unlock:entry
{
    printf("thread %d: mutex 0x%08x unlocked\n", __tid, __arg1);
}

```

- User has to map the Fortran data types to ProbeVue data types and use the same in the script. The mapping of Fortran basic data types to ProbeVue data types is listed in the below table.

Fortran data-type	ProbeVue data-type
INTEGER * 2	short
INTEGER * 4	int/long
INTEGER * 8	long long
REAL	float
DOUBLE PRECISION	double
COMPLEX	No equivalent basic data type. This needs to be mapped to a structure as shown below: <pre>typedef struct complex { float a; float b; } COMPLEX;</pre>
LOGICAL	int (The Fortran standard requires logical variables to be the same size as INTEGER/ REAL variables)
CHARACTER	char
BYTE	signed char

- Fortran passes IN scalar arguments of internal procedures by value, and other arguments by reference. Arguments passed by reference should be accessed with `copy_userdata()`. More information on argument association in fortran can be found in the [Argument association](#) topic.
- Routine names in a Fortran program is case in-sensitive. But, while specifying them in a ProbeVue script, they should be in lower-case .

The following sample script illustrates how to map Fortran data types to ProbeVue data types:

```

/* cmp_calc.e */
/* Trace fortran routines
cmp_calc(COMPLEX, INTEGER) and
cmplx_d(void) */

typedef struct complex{
    float a;
    float b;
} COMPLEX;

typedef int INTEGER;

/* arguments are indicated to be of pointer type as they are passed by reference */
void cmp_calc(COMPLEX *, INTEGER *);
void cmplx_d();

@@uft:$__CPID:*:cmplx_d:entry
{
printf("In cmplx_d entry \n");
}

@@uft:$__CPID:*:cmp_calc:entry
{
COMPLEX c;
int i;
copy_userdata(__arg1, c);
copy_userdata(__arg2, i);
printf("%10.7f+j%9.7f  %d \n", c.a,c.b,i);
}

```

- Fortran stores arrays in column-major form, whereas ProbeVue stores in row-major form and the below script shows how users can retrieve the array elements.

```

/* array.e*/
/* ProbeVue script to probe fortran program array.f */

void displayarray(int **, int, int);
@@uft:$__CPID:*:displayarray:entry
{
int a[5][4];          /* row and column sizes are interchanged */
copy_userdata(__arg1, a);
/* to print the first row */
printf("%d %d %d \n", a[0][0], a[1][0], a[2][0]);
/* to print the second row */
printf("%d %d %d \n", a[0][1], a[1][1], a[2][1]);
}

/* Fortran program array.f */

PROGRAM ARRAY_PGM
IMPLICIT NONE
INTEGER, DIMENSION(1:4,1:5) :: Array
INTEGER :: RowSize, ColumnSize
CALL ReadArray(Array, RowSize, ColumnSize)
CALL DisplayArray(Array, RowSize, ColumnSize)
CONTAINS
SUBROUTINE ReadArray(Array, Rows, Columns)
IMPLICIT NONE
INTEGER, DIMENSION(1:,1:), INTENT(OUT) :: Array
INTEGER, INTENT(OUT) :: Rows, Columns
INTEGER :: i, j
READ(*,*) Rows, Columns
DO i = 1, Rows
READ(*,*) (Array(i,j), j=1, Columns)
END DO
END SUBROUTINE ReadArray
SUBROUTINE DisplayArray(Array, Rows, Columns)
IMPLICIT NONE
INTEGER, DIMENSION(1:,1:), INTENT(IN) :: Array
INTEGER, INTENT(IN) :: Rows, Columns
INTEGER :: i, j
DO i = 1, Rows
WRITE(*,*) (Array(i,j), j=1, Columns )
END DO
END SUBROUTINE DisplayArray
END PROGRAM ARRAY_PGM

```

- Intrinsic or built-in functions cannot be probed with ProbeVue . All FORTRAN routines as listed in the XCOFF symbol table of the executable/linked libraries can be probed. ProbeVue uses the XCOFF symbol table to identify the location of these routines. However, the prototype for the routine has to be provided by the user and ProbeVue tries to access the arguments according to the prototype provided. For routines where the compiler mangles the routine names, the mangled name should be provided. Since Vue is a C-style language, user should ensure that the FORTRAN function/subroutine prototype is appropriately mapped to C language style function prototype. Please refer to the linkage conventions for argument passing and function return values in the [Passing data from one language to another](#) topic. The below example illustrates this:

```

/* Fortran program ext_op.f */
/* Operator "*" is extended for rational multiplication */
MODULE rational_arithmetic
IMPLICIT NONE
  TYPE RATNUM
    INTEGER :: num, den
  END TYPE RATNUM
  INTERFACE OPERATOR (*)
    MODULE PROCEDURE rat_rat, int_rat, rat_int
  END INTERFACE
CONTAINS
  FUNCTION rat_rat(l,r)      ! rat * rat
    TYPE(RATNUM), INTENT(IN) :: l,r
    TYPE(RATNUM) :: val,rat_rat
    val.num=l.num*r.num
    val.den=l.den*r.den
    rat_rat=val
  END FUNCTION rat_rat
  FUNCTION int_rat(l,r)     ! int * rat
    INTEGER, INTENT(IN)    :: l
    TYPE(RATNUM), INTENT(IN) :: r
    TYPE(RATNUM) :: val,int_rat
    val.num=l*r.num
    val.den=r.den
    int_rat=val
  END FUNCTION int_rat
  FUNCTION rat_int(l,r)    ! rat * int
    TYPE(RATNUM), INTENT(IN) :: l
    INTEGER, INTENT(IN)    :: r
    TYPE(RATNUM) :: val,rat_int
    val.num=l.num*r
    val.den=l.den
    rat_int=val
  END FUNCTION rat_int
END MODULE rational_arithmetic
PROGRAM Main1
Use rational_arithmetic
IMPLICIT NONE
  TYPE(RATNUM) :: l,r,l1
  l.num=10
  l.den=11
  r.num=3
  r.den=4
  l1=l*r
END PROGRAM Main1

/* ext_op.e */
/* ProbeVue script to probe routine that gets called when "*"
   is used to multiply rational numbers in ext_op.f */

struct rat
{
    int num;
    int den;
};
void __rational_arithmetic_NMOD_rat_rat(struct rat*,
    struct rat*,struct rat*);
/* Note that the mangled function name is provided. */
/* Also, the structure to be returned is sent in the buffer whose address is provided as the first
argument. */
/* The first explicit parameter is in the second argument. */
@@BEGIN
{
    struct rat* rat3;
}
@@uft:$__CPID:*:__rational_arithmetic_NMOD_rat_rat:entry
{
    struct rat rat1,rat2;
    copy_userdata((struct rat *)__arg2,rat1);
    copy_userdata((struct rat *)__arg3,rat2);
    rat3=__arg1;
    /* The address of the buffer where the returned structure will be stored is saved at the function
entry */
    printf("Argument Passed rat_rat = %d:%d,%d:%d\n",rat1.num,rat1.den,rat2.num,rat2.den);
}
@@uft:$__CPID:*:__rational_arithmetic_NMOD_rat_rat:exit

```

```

{
    struct rat rrat;
    copy_userdata((struct rat *)rat3,rrat);
    /* The saved buffer address is used to fetch the returned structure */
    printf("Return from rat_rat = %d:%d\n",rrat.num,rrat.den);
    exit();
}

```

- ProbeVue won't support direct inclusion of Fortran header files in the script. However, a mapping of Fortran data types to ProbeVue data types can be provided in a ProbeVue header file and included with the "-I" option.

Related concepts

[Java applications probe manager](#)

C++ applications probe manager

C++ Probe Manager supports probing of C++ applications in a way identical to C probe managers. Support for "uft" style entry/exit probes on any C++ function, including member, overloaded, operator, and template functions in the core executable. A function entry/exit probe in C++ must use the **@@uftxlc++** probe manager.

All tuples in the **@@uftxlc++** style probe specifications have the same usage and format as for the **@@uft** style probe strings, with the exception of the function name. Because C++ allows a single function name to be overloaded, the function name specified in the probe string may have to include the function's argument types to uniquely identify the function being probed.

For example:

```
@@uftxlc++:12345:*:"foobar(int, char *)":entry
```

```
@@uftxlc++:/usr/test:*:"foobar(int, char *)":entry
```

Note: The return type is missing from the above probe string because it does not take part in the name mangling algorithm for regular functions. In case of a template function, the user must specify an explicit template instantiation to probe on and must also specify the return type of the template instantiation:

```
@@uftxlc++:12345:*:void foobar<int>(int, char *):entry
```

Note: The probe strings must use quotes around the function name as specified in above two examples and the probevue command will signal an error if the quotes are missing. The quotes are not only because of the colon ":" but also because of the comma ",". The comma operator is used for separating multiple probes on the same line and without the quotes it takes precedence. This results in very strange error messages for the user.

When probing a class member function or a function defined in a namespace, the fully qualified function name must be used in the probe string. To avoid any ambiguity between the single colon (:) tuple separator in probe strings and the double colon (::) scope resolution operator in a fully qualified C++ name, the entire function name tuple in the probe string must be quoted.

```
@@uftxlc++:12345:*:"Foo::bar(int)":entry
```

Limitations:

1. Access to data fields that are inherited from a virtual base class is not supported.
2. Template classes are not supported and must not be included in the C++ header.
3. Pointers to members are not supported.
4. To probe a class with the class definition, an object of the class is instantiated in the header file either as a global object or in a dummy function.

Example:

Below is c++ application

```
#include "header.cc"
main()
{
int      i = 10;
incr_num(i);
float    a = 3.14;
incr_num(a);
char     ch = 'A';
incr_num(ch);
double   d = 1.11;
incr_num(d);
}
```

Content of the "header.cc"

```
# cat header.cc
#include <iostream.h>
template <class T>
T incr_num( T a)
{
return (++a);
}
int dummy()
{
int i=10,j=20;
incr_num(i);
float a=3.14;
incr_num(a);
char ch='A',dh='Z';
incr_num(ch);
double d=1.1,e=1.11;
incr_num(d);
return 0;
}
```

Content of the Vue script vue_cpp.e

```
##C++
#include "header.cc"
##Vue
@@uftxlc++:$__CPID:*:"incr_num<int>(int)":entry
{
printf("Hello1_%d\n",__arg1 );
}
@@uftxlc++:$__CPID:*:"incr_num < float > (float)" :entry
{
printf("Hello2_%f\n",__arg1 );
}
@@uftxlc++:$__CPID:*:"incr_num < char > ( char)":entry
{
printf("Hello3_%c\n",__arg1 );
}
@@uftxlc++:$__CPID:*:"incr_num < double > ( double)":entry
{
printf("Hello4_%lf\n",__arg1 );
exit();
}
```

Execution :

```
/usr/vacpp/bin/xlc app.c++
# probevue -X ./a.out vue_cpp.e
Hello1_10
Hello2_3.140000
Hello3_A
Hello4_1.110000
```

The function prototype in the fourth tuple can be specified as an Extended Regular Expression (ERE). The ERE should be enclosed between “/” and “/” like “/<ERE>/”. When function prototype is specified as an ERE, all the functions matching the specified regular expression in the specified module will be probed.

```
/* Probe entry of all the C++ functions in the executable a.out */
@@uftxlc++:$__CPID:a.out:"/.*"/:entry
```

```
/* Probe exit of all the C++ functions with 'foo' word in it */  
@@uftxlc++:$__CPID:*:"/foo/":exit
```

In the entry probes, where a function name is specified as a regular expression, individual arguments cannot be accessed. However, probevue function **print_args()** can be used to print the function name and its arguments. The argument values is printed based on the argument type information available in the traceback table of the function.

In the exit probes, where a function name is specified as a regular expression, return value cannot be accessed.

Java applications probe manager

Java Probe Manager (JPM) supports probing of Java applications in a way identical to C and C++ probe managers. A single Vue script should be able to trace multiple java applications at the same time by using different process IDs of the JVMs. The same script can be used to probe syscalls or C/C++ applications along with Java applications and can use other probe managers.

Like uft (user function tracing) probe manager java probe manager also accepts 5-tuple probe specification in the following format:

```
uftjava :< process_ID> :*< _qualified_function_name >: entry
```

Where the second tuple is the process ID of JVM process corresponding to the Java application that is being traced.

Third field: reserved for future use.

Fourth field: where the java method needs to be specified.

This name is a completely qualified name as used in java applications like Mypackage.Myclass.Mymethod.

Some of the restrictions that may apply are

- Only pure java methods can be probed, Native (shared library calls) or encrypted codes are not traceable.
- Only entry probes are supported.
- Can support only JVM v 1.5 and above that supports JVMTI interface.
- At any given point of time, no two Probevue sessions can probe the same Java application with @@uftjava.
- Polymorphic/Overloaded methods are not supported.
- Tracing/accessing external variables with same name as any of the Probevue keywords or built-in names are not supported. This may need those external symbols (Java application variable names) to be renamed.
- Accessing arrays of java applications is not supported in this release.
- Accessing arrays of java applications is not supported in this release.
- get_function () built-in for java language is not supported in this release.

Note: In case of tracing non static methods, argument number starts with __arg2 like non static methods of C++. The __arg1 is used for self reference (this pointer).

Data Access: The action blocks of java probes can access the following data similar to existing behavior.

- Action block can access global, local and kernel script variables.
- Action block can access method arguments (Entry class variables) of primitive types.
- Action block can access the built-in variables.
- Action block can access Java application variables through fully qualified names, only static (class members).

```
x = some_package.app.class.var_x; //Access static/class member.
```

- Accessing java application primitive types variables is supported; they must be converted/promoted/ casted implicitly without losing value to equivalent types in Vue language. But the actual memory usage (size) may differ from that of Java language.

The functions supported in the context of Java probe manager are listed in the following table:

Table 26. Supported functions by Java probe manager

Function	Description
stktrace()	Provides the Stack trace of the Java application (running thread) that is being traced.
copy_userdata()	Copy data from java application into script variables.
get_probe()	Returns the probe string.
get_stktrace	Returns the runtime stack trace.
get_location_point()	Returns the current probe location.
get_userstring()	Copy string data from java application.
exit()	exits from the probevue trace session.

Changes to Probevue command:

Table 27. probevue command change

Command	Description
-X option	This option can be used (along with -A option) to launch Java application, in the current release the user has to manually pass an additional optional string agentlib:probevuejava along with all the other options that are needed to run the java application.

For Example:

```
probevue -X /usr/java5/bin/java -A -agentlib:probevuejava myjavaapp myscript.e
```

When running the 64 bit JVM, we have to use "agentlib:probevuejava64" as in:

```
probevue -X /usr/java5_64/bin/java -A -agentlib:probevuejava64 myjavaapp myscript.e
where myjavaapp is the java class of myjavaapp.java application
```

Example ExtendedClass.java Source:

```
class BaseClass
{
    static int i=10;

    public static void test(int x)
    {
        i += x;
    }
}

public class ExtendedClass extends BaseClass
{
    public static void test(int x, String msg)
    {
        i += x;
        System.out.print("Java: " + msg + "\n\n");
        BaseClass.test(x);
    }

    public static void main(String[] args)
    {
        BaseClass.test(5);
        ExtendedClass.test(10, "hello");
    }
}
```

```
}  
}
```

Example test.e script for above Java application:

```
@@uftjava:$__CPID:*:"BaseClass.test":entry  
{  
    printf("BaseClass.i: %d\n", BaseClass.i);  
    printf("BaseClass.test: %d\n", __arg1);  
    stktrace(0, -1);  
    printf("\n");  
}  
  
@@uftjava:$__CPID:*:"ExtendedClass.test":entry  
{  
    printf("BaseClass.i: %d\n", BaseClass.i);  
    printf("ExtendedClass.test: %d, %s\n", __arg1, __arg2);  
    stktrace(0, -1);  
    printf("\n");  
}
```

Example ProbeVue session with above script:

```
# probevue -X /usr/java5/jre/bin/java \  
-A "-agentlib:probevuejava ExtendedClass" test.e  
Java: hello  
  
BaseClass.i: 10  
BaseClass.test: 5  
BaseClass.test()+0  
ExtendedClass.main()+1  
  
BaseClass.i: 15  
ExtendedClass.test: 10, hello  
ExtendedClass.test()+0  
ExtendedClass.main()+8  
  
BaseClass.i: 25  
BaseClass.test: 10  
BaseClass.test()+0  
ExtendedClass.test()+39  
ExtendedClass.main()+8
```

Related concepts

[UFT probe manager](#)

Interval probe manager

The interval probe manager provides probe points that fire at a user-defined time-interval. The probe points are not located in kernel or application code, but instead are based on wall clock time interval based probe events.

The interval probe manager is useful for summarizing statistics collected over an interval of time. It accepts a 4-tuple or a 5-tuple probe specification in the following format:

```
@@interval:*:clock:<# milliseconds>:[*|cpu_ids]
```

The interval probe manager will filter probe events by process ID if it is provided in the second field. Assigning the * to the second field indicates that the probe will be fired for all processes. Further, the only value supported by the interval probe manager for the third field is the clock keyword that identifies the probe specification as being for a wall clock probe. The fourth field, that is the <# milliseconds> field, identifies the number of milliseconds between firings of the probe. The interval probe manager requires that the value for this field consist only of digits 0-9.

The fifth field, <cpu_ids> identifies CPU IDs of the CPUs on which the probe occurs. The value of this field must be specified in the range of CPU IDs, or * or individual CPU IDs. The process ID specified in the second tuple and the CPU ID specified in the fifth tuple are mutually exclusive. Hence, the interval probe manager filters the probe events either by the specified process ID in tuple 2 running on any CPU or filters events for all the processes running on the specified CPUs. The minimum supported clock value is 100 milliseconds. Only one CPU bound interval probe per logical CPU in the system. This is an optional field.

For interval probes without process Id, intervals should be exactly divisible by 100. Thus, probe events that are apart by 100ms, 200ms, 300ms, and so on, are allowed in non-profiling interval probes. For interval probes with process Id specified, intervals should be greater or equal to minimum interval allowed for global root user or exactly divisible by 10 for other users. Thus, probe events that are apart by 10ms, 20ms, 30ms, and so on, are allowed for normal users in profiling interval probes. Only one profiling interval probe can be active for a process.

Examples for Interval probe specification for CPU bound probes are listed below:

1. To get the context information, run the user defined ProbeVue clause on all CPUs (Specifying “*” in the cpu_ids field denotes all the CPUs).

```
@@interval:*:clock:100:*
```

2. To get the context information for CPUs in the range 10 to 20 (Specify the format as x-y, where x=10 and y=20 for the range of CPUs).

```
@@interval:*:clock:100:10-20
```

3. To get the context information for CPUs 10 and 12 (Specify the format with the ‘|’ (pipe) symbol which separates each CPU IDs. This format is used to specify more than one CPU ID.

```
@@interval:*:clock:100:10|12
```

4. To get the context information for set of CPUs 10 and 12-20.

```
@@interval:*:clock:100:10|12-20
```

Note: User can provide either “*” or set of CPUIDs.

Note: The interval probe manager does not guarantee that a probe will be fired exactly the number of milliseconds apart as indicated by value of the fourth field. Higher-priority interrupts and code that runs after disabling all interrupts can cause the probe to fire later than the specification.

The interval probe manager requires only basic dynamic tracing privileges. The interval probe manager enforces the following limits on the number of probes it supports to prevent malicious users from running the kernel out of memory by creating huge numbers of interval probes.

Table 28. Limits specified by the interval probe manager

Interval	Count
Maximum number of interval probes per user	32
Maximum number of interval probes in system	1024

The interval probe manager does not support the following functions. If used inside an interval manager probe point, these functions will generate an empty string or zero as output.

- **get_function**
- **get_probe**
- **get_location_point**

When process ID is not specified, an interval probe can trigger in the context of any process depending upon when the probe fires since the probe event is based on wall clock time. Because of this, the ProbeVue framework does not allow the use of any of the following functions inside the interval probe manager's action block to prevent unauthorized access to a process's internal data. This security violation is caught only in the kernel. The Vue script will successfully compile but the session will fail to initialize.

- **stktrace**
- **get_userstring**

These functions provide no value when used from the probe manager. Even if you are the root user, you cannot call these functions inside the interval probe manager.

When the process ID is specified, the interval probe is triggered for all the threads within the process at the specified time interval. As the probe is fired in the context of the process, **stktrace()** function and **__pname** built-in is allowed inside the interval probe manager's action block, unlike when process ID is not specified.

System trace probe manager

The system trace probe manager provides probe points wherever existing system trace hooks to trace channel zero (system event channel) occur, both within the kernel and within applications. To use this probe manager, you must have the kernel access privilege, and not be running in a WPAR.

The system trace probe manager accepts a 3-tuple probe specification in the following format:

```
@@systrace:*:<hookid>
```

where the *hookid* argument specifies the ID for the specific system trace hook of interest. The *hookid* argument consists of 4 hex digits typically of the form hhh0. For example, to specify the *hookid* argument for the **fork** system call, specify 1390. See the `/usr/include/sys/trchkid.h` file for examples, such as `HKWD_SYSC_FORK`. The entries in this file are hook words, where the *hookid* value is in the upper halfword. Because hook words can be arbitrary, no validation of the *hookid* argument beyond checking that it is a valid hex string of up to 4 hex digits is performed. It is not an error to specify a *hookid* value that never occurs.

As a convenience, you can specify the *hookid* argument with fewer than 4 hex digits. In this case, first a trailing zero is assumed, and then additional leading zeroes as necessary to implicitly define the required 4 digits. For example, you can use 139 as an abbreviation of 1390. Similarly, 0100, 010, and 10 all specify the same *hookid* value, taken from `HKWD_USER1`.

You can specify the *hookid* argument with the * wildcard character. This will probe all system tracing, with likely unacceptable performance implications. Hence, such a specification must be used only when absolutely necessary.

The second tuple is reserved, and must be specified as an asterisk, as shown.

Only system trace events that actually occur and record system trace data trigger probes. In particular, a system trace probe can only occur when system trace is active. The systrace probe manager is an event-based probe manager. Hence, probe name, function name, and location point are not available. As the hookword is passed to the script, this is not a significant restriction.

A non-root user is limited to at most 64 systrace probes simultaneously enabled. No more than 128 explicit systrace probes can be enabled system-wide.

ProbeVue built-in register variables allow access to the data traced. You cannot use the `__arg*` variables for this purpose. There are two general styles for system tracing.

The following style is for the `trchhook(64)/utrchhook(64)` (or the equivalent `TRCHKLx` macros in C) hooks:

- `__r3` contains the 16 bit hookid.
- `__r4` contains the subhookid.
- `__r5` contains traced data word D1.
- `__r6` contains traced data word D2.
- `__r7` contains traced data word D3.
- `__r8` contains traced data word D4.
- `__r9` contains traced data word D5.

Not all trace hooks contain all 5 data words. Undefined data words from a given trace hook will appear as zero. The Vue clause for a given hook ID must know exactly what and how much data its hook ID traces.

If the trace record was produced by one of the functions in the `trcgen` or `trcgent` family, use the following style:

- `__r3` contains the 16 bit hookid.

- __r4 contains the subhookid.
- __r5 contains traced data word D1.
- __r6 contains the length of the traced data.
- __r7 contains the address of the traced data.

The following script shows a simple example of the systrace probe manager:

```

@@systrace:*:1390
{
    if (__r4 == 0) { /* normal fork is traced with subhookid zero */
        printf("HKWD_SYSC_FORK: %d forks child %d\n", __pid, __r5);
        exit();
    }
}

```

Systrace probe manager is independent of the system trace function and can trace hook locations even if the system trace function is not in an active state. You can turn on the system trace function when a ProbeVue session is active.

Systrace probe manager uses AIX operating systems kernel to facilitate some of its constructs. Tracking hooks that are currently used by ProbeVue will adversely affect AIX reliability. Therefore, some ProbeVue constructs cannot be used when you trace such hook IDs. The following table displays the exception of the constructs:

Note: The following constructs are ignored if you want to trace all hooks by using the @@systrace:*:* specification. You might not be able to display stack trace if the AIX kernel prohibits generating an exception in the environment where a hook is located. The ability of ProbeVue to display stack traces is determined at run time.

Number	Trace-Hook	Construct
1	HKWD_KERN_HCALL	ALL
2	HKWD_KERN_SLIH	Associative array, range, stktrace variable, __stat
3	HKWD_KERN_LOCK	Associative array, range, stktrace variable, __stat
4	HKWD_KERN_UNLOCK	Associative array, range, stktrace variable, __stat
5	HKWD_KERN_DISABLEMENT	ALL
6	HKWD_KERN_DISPATCH	__ublock, stktrace, get_stktrace, __pname, __execname, __errno
7	HKWD_KERN_DISPATCH_SRAD	__ublock, stktrace, get_stktrace, __pname, __execname, __errno
8	HKWD_KERN_DISP_AFFIN	__ublock, stktrace, get_stktrace, __pname, __execname, __errno
9	HKWD_KERN_UNDISP	__ublock, stktrace, get_stktrace, __pname, __execname, __errno
10	HKWD_KERN_IDLE	__ublock, stktrace, get_stktrace, __pname, __execname, __errno
11	HKWD_KERN_FLIH	Associative array, range, stktrace variable, __stat

Number	Trace-Hook	Construct
12	HKWD_KERN_RESUME	Associative array, range, stktrace variable, __stat
13	HKWD_KERN_VPM	Associative array, range, stktrace variable, __stat
14	HKWD_PM_NOTIFY	Associative array, range, stktrace variable, __stat

With the appropriate privilege, a Vue script can generate system trace records by using the reliability, availability, and serviceability (RAS) events Vue functions. However, the Systrace probe manager does not detect trace records generated from a Vue script.

Extended system call probe manager (syscallx)

The syscallx probe manager, on the other hand, allows all base system calls to be traced. Base system calls is the set of system calls exported by the kernel and base kernel extensions, which are available immediately after boot-up. System calls that are exported from kernel extensions that may loaded later are not supported. Either a specific system call or all system calls can be specified through the probe point tuple. However, unlike the syscall probe manager, the third field of the probe point tuple for the syscallx must identify the actual kernel entry point function. The syscallx probe manager also limit probes to fire in a specific process if the process ID is specified as the second field of the probe point tuple.

The following are some examples:

```
/* Probe point tuple to probe the read system call entry for all processes */
@@syscallx:*:kread:entry
```

```
/* Probe point tuple to probe the fork system call exit for process with ID 434 */
@@syscallx:434:kfork:exit
```

```
/* Probe point tuple to probe entry for all base system calls */
@@syscallx:*:*:entry
```

```
/* Probe point tuple to probe exit for all base system calls for process 744 */
@@syscallx:744:*:exit
```

System calls supported by the syscall probe manager

The following table lists the system calls supported by the syscall probe manager along with the actual entry name in the kernel.

Note: The kernel entry name is provided here only for documentation purposes. The kernel entry names can change between releases or even after a service update.

System call name	Kernel entry name
absinterval	absinterval
accept	accept1
bind	bind
close	close
creat	creat
execve	execve

Table 30. System calls supported by the syscall probe manager (continued)

System call name	Kernel entry name
exit	_exit
fork	kfork
getgidx	getgidx
getgroups	getgroups
getinterval	getinterval
getpeername	getpeername
getpid	_getpid
getppid	_getppid
getpri	_getpri
getpriority	_getpriority
getsockname	getsockname
getsockopt	getsockopt
getuidx	getuidx
incinterval	incinterval
kill	kill
listen	listen
lseek	klseek
mknod	mknod
mmap	mmap
mq_close	mq_close
mq_getattr	mq_getattr
mq_notify	mq_notify
mq_open	mq_open
mq_receive	mq_receive
mq_send	mq_send
mq_setattr	mq_setattr
mq_unlink	mq_unlink
msgctl	msgctl
msgget	msgget
msgrcv	__msgrcv
msgsnd	__msgsnd
nsleep	_nsleep
open	kopen
pause	_pause

Table 30. System calls supported by the syscall probe manager (continued)

System call name	Kernel entry name
pipe	pipe
plock	plock
poll	_poll
read	kread
reboot	reboot
recv	_erecv
recvfrom	_enrecvfrom
recvmsg	_erecvmsg
select	_select
sem_close	_sem_close
sem_destroy	sem_destroy
sem_getvalue	sem_getvalue
sem_init	sem_init
sem_open	_sem_open
sem_post	sem_post
sem_unlink	sem_unlink
sem_wait	_sem_wait
semctl	semctl
semget	semget
semop	__semop
semtimedop	__semtimedop
send	_esend
sendmsg	_esendmsg
sendto	_esendto
setpri	_setpri
setpriority	_setpriority
setsockopt	setsockopt
setuidx	setuidx
shmat	shmat
shmctl	shmctl
shmdt	shmdt
shmget	shmget
shutdown	shutdown
sigaction	_sigaction

Table 30. System calls supported by the syscall probe manager (continued)

System call name	Kernel entry name
sigpending	_sigpending
sigprocmask	sigprocmask
sigsuspend	_sigsuspend
socket	socket
socketpair	socketpair
stat	statx
waitpid	kwaitpid
write	kwrite

Running in a WPAR

Workload partitions or WPARs are virtualized operating system environments within a single instance of the AIX operating system. The WPAR environment is somewhat different from the standard AIX operating system environment.

Dynamic tracing is supported in the WPAR environment. By default, when creating a WPAR, only the **PV_PROBEVUE_TRC_USER_SELF** and the **PV_PROBEVUE_TRC_USER** privileges are assigned to the WPAR and the superuser (root) on a WPAR system will be granted these privileges. An admin user from the global partition can change the value of the default WPAR privilege set or can explicitly assign additional privileges when creating the WPAR.

Privileges on WPAR have generally the same meanings as on a global partition. Be careful when assigning **PV_PROBEVUE_TRC_KERNEL** or the **PV_PROBEVUE_TRC_MANAGE** to a WPAR. Any user with **PV_PROBEVUE_TRC_KERNEL** privilege can access global kernel variables while a user with **PV_PROBEVUE_TRC_MANAGE** privilege can change the values of ProbeVue parameters or shutdown ProbeVue. These changes affect all users even those in other partitions.

When you issue the **probevue** command in a WPAR, processes running in other WPARs or in the global partition are not visible to it. Because of this, you can only probe processes in your same WPAR. The **probevue** command will fail if the probe specification contains a process ID that is outside its partition. The **PV_PROBEVUE_TRC_USER** and **PV_PROBEVUE_TRC_SYSCALL** privileges in a WPAR only allow you to probe user space functions or system calls of processes that are in your WPAR. When probing system calls, the second field of the syscall probe specification must be set to a valid WPAR-visible process ID. Assigning the value * to the second field is not supported.

When a ProbeVue session is initiated in a mobile WPAR, it temporarily switches the WPAR to a non-checkpointable state. After the ProbeVue session terminates, the WPAR is checkpointable again.

I/O probe manager

I/O probe manager provides capabilities to trace I/O operation events in various layers of AIX I/O stack. Use the syscall probe manager to trace application I/O request that is triggered by a read/write system call. Use I/O probe manager to probe further into the syscall layer.

Use I/O probe manager to analyze response time of I/O operations of a block device that segregates the service time and queuing delay.

The following layers are supported:

- Logical File System (LFS)
- Virtual File System (VFS)
- Enhanced Journaled File Systems (JFS2)
- Logical Volume Manager (LVM)

- Small Computer System Interface (SCSI) disk driver
- Generic block devices

The primary use cases for I/O probe manager are as follows:

- Identify the following patterns of I/O usage of a device. Valid devices can be a disk, logical volume, or volume group, or file system (type or mount path) in a specified time period:
 - I/O operation count
 - Size of I/O operations
 - Type of I/O operation (read/write)
 - Sequential or random nature of I/O
- Get process or thread-wise usage information of a file system (type or mount path), logical volume, volume group, or disk.
- Get an end-to-end mapping of I/O flow among various layers (wherever possible).
- Monitor a specific I/O resource usage. For example:
 - Trace any write operations of the `/etc/passwd` file.
 - Trace read operation on block 0 of the `hdisk0` device.
 - Trace when a new logical volume is opened in root volume group (`rootvg`).
- For Multipath I/O (MPIO) disks, get path-specific information by the following actions:
 - Get path-wise usage and response time information.
 - Identify path switching or path failure.
- For I/O errors, get more details about the error in disk driver layer.

Probe specification

I/O probes must be specified in the following format in Vue script:

```
@@io:sub_type:io_event:operation_type:filter[|filter ...]
```

This specification consists of five tuples that are separated by colon (:). The first tuple is always `@@io`.

Probe sub type

The second tuple signifies the sub type of the probe that indicates the layer of AIX I/O stack that contains the probe. This tuple can have one of the following values:

<i>Table 31. Second tuple for probes</i>	
Second tuple (sub type)	Description
disk	This probe starts for disk driver events. Currently, the I/O probe manager supports only the <code>scsidisk</code> driver.
lvm	This probe starts for Logical Volume Manager (LVM) events.
bdev	This probe starts for any block I/O device. Disk, CD-ROM, diskette are examples of block devices. This sub type is used only when no other sub type is applicable. For example, if a block device is not a disk, volume group, or logical volume, this sub type is applicable.
jfs2	This probe starts for JFS2 file system events.
vfs	This probe starts for any read/write operation on a file.

Note: The second tuple cannot have a value of asterisk (*).

For a disk type of second tuple, the third tuple can have the following values:

Sub type (Second tuple)	I/O event (Third Tuple)	Description
disk	entry	This probe starts whenever disk driver receives an I/O request to process.
	iostart	This probe starts when the disk driver picks up an I/O request from its ready queue and sends it down to lower layer (for example, adapter driver). A single original I/O request to disk driver can send multiple command requests (some might be driver-related task management command requests) to lower layer. However, sometimes the driver can combine multiple original requests and send a single request to lower layer.
	iodone	This probe starts when the lower layer (for example, adapter driver) returns an I/O request (successful or failed) to disk driver.
	exit	This probe starts when disk driver returns an I/O request (successful or failed) to its upper layer.

Note: The members of the following built-in values are available in the probes that are mentioned for the probe sub type: `__iobuf`, `__diskinfo`, `__diskcmd` (only in `disk:iostart` and `disk:iodone`), and `__iopath` (only in `disk:iostart` and `disk:iodone`).

For every entry, a corresponding exit probe is defined that has the same `__iobuf->bufid` value available at both the probe points. The entry event can be followed by multiple `iostart` events, but at least one of them must have the same `__iobuf->bufid` value. Every `iostart` event has a matching `iodone` event that has the same `__iobuf->child_bufid` value.

For an LVM type of second tuple, the third tuple can have the following values:

Sub type (second tuple)	I/O event (third tuple)	Description
lvm	entry	This probe starts whenever the LVM layer receives an I/O request to process.
	iostart	This probe starts when LVM picks an I/O request from its ready queue and sends down to the lower layer (usually the disk driver).
	iodone	This probe starts when the lower layer (for example, disk driver) returns an I/O request (successful or failed) to LVM.
	exit	This probe starts when LVM returns an I/O request (successful or failed) to its upper layer.

Note: The members of the following built-ins values are available in the probes that are mentioned for LVM: `__iobuf`, `__lvol`, and `__volgrp`. Every entry has a corresponding exit probe, which has the same `__iobuf->bufid` value available at both the probe points.

The entry event can be followed by multiple `iostart` events, but at least one of them has the same `__iobuf->bufid` value. Every `iostart` event has a matching `iodone` event that has the same `__iobuf->child_bufid` value.

For generic block device probes, the third tuple can have the following values:

<i>Table 34. Generic block device second tuple: Third tuple values</i>		
Sub type (second tuple)	I/O event (third tuple)	Description
bdev	iostart	This probe gets fired when any block I/O (for example, disk, logical volume, CD-ROM) device is initiated. It happens when the AIX devstrat kernel service is called by any code.
	iodone	This probe gets fired when a block I/O request completion happens, when the AIX iodone kernel service is called by any code.

Note: The members of the following built-in values are available in the probes that are mentioned in bdev: __iobuf. Every iostart event has a matching iodone event that has the same __iobuf->bufid value.

For JFS2 file system probes, the third tuple can have the following values:

<i>Table 35. JFS2 second tuple: Third tuple values</i>		
Sub type (second tuple)	I/O event (third tuple)	Description
jfs2	buf_map	This probe starts when a logical file extent gets mapped to an I/O buffer and is sent to the underlying logical volume.

Note: The members of the following built-in values are available in the probe that is mentioned for JFS2 file system probes: __j2info.

For Virtual file system (VFS) probes, the third tuple can have the following values:

<i>Table 36. VFS second tuple: Third tuple values</i>		
Sub type (second tuple)	I/O event (third tuple)	Description
vfs	entry	This probe starts when any read/write operation on a file is initiated.
	exit	This probe starts when any read/write operation on a file is completed (whether success or failure).

Note: The members of the following built-in are available in the probe that is mentioned in VFS probes: __file.

For the same thread, every entry is followed by an exit event that has the same __file->inode_id value.

Probe operation type

The fourth tuple indicates the type of I/O operation that is specified by the probe. The fourth tuple can have one of the following values:

<i>Table 37. Fourth tuple for I/O operation</i>	
Fourth tuple	Description
read	The probe starts for only the read operation.
write	The probe starts for only the write operation.
*	The probe starts for both read and write operations.

Probe filter

The fifth tuple is the filter tuple that helps in filtering more specific probes according to the requirement. The possible values are subtype dependent. Multiple values can be specified separated by | character, and the probe starts if it matches any of those filters. If the value of the fifth tuple is *, no filtering occurs and the probe starts if other tuples match. If multiple selectors are specified, and one of them is *, it is equivalent to the whole tuple value of *.

For disk probes, the fifth tuple can have the following values:

<i>Table 38. Disk filter tuple</i>	
Filter (fifth tuple)	Description
Disk name. For example, hdisk0	The probe action is run only for the particular disk.
Disk type. Allowed symbols: FC, ISCSI, VSCSI, SAS	The probe action is run only for disks with matching type. The meanings of the symbols are as follows: <ul style="list-style-type: none"> • FC: Fibre Channel disk • ISCSI: iSCSI disk • VSCSI: Virtual SCSI disk (on VIOS client) • SAS: Serial Attached SCSI disks

Note: The disk name and disk type can be combined as filters. For example, the following probe starts for either hdisk0 or any other FC disk (at disk entry event, for both read/write operation type)

```
@@io:disk:entry:*:hdisk0|FC
```

For Logical Volume Manager (LVM) probes, the fifth tuple can have the following values:

<i>Table 39. LVM filter tuple</i>	
Filter (fifth Tuple)	Description
Logical volume name, for example hd5, lg_dumplv	The probe action is run only for the particular logical volume.
Volume group name, for example rootvg	The probe action is run only for those logical volumes that belong to a particular volume group.

The following probe starts for any logical volume that belongs to either root volume group (rootvg), or test volume group (testvg) (at iostart event, for write operation only):

```
@@io:lvm:iostart:write:rootvg|testvg
```

For generic block device probes, fifth tuple can have following values:

<i>Table 40. Generic block device filter tuple</i>	
Filter (fifth tuple)	Description
Block device name, for example: hdisk0, hd5, cd0	The probe action is run only for the particular block device.

Consider the following examples for generic block device probes:

```
@@io:bdev:iostart:*:cd0
```

```
@@io:bdev:iodone:read:hdisk3|hdisk5
```

For JFS2 file system probes, the fifth tuple can have following values:

<i>Table 41. JFS2 filter tuple</i>	
Filter (fifth tuple)	Description
File system mount path, for example: /usr	The probe action is run only for the file system with the particular mount path. It must be a JFS2 file system, otherwise ProbeVue rejects that probe specification.

Consider following examples for the JFS2 file system probes:

```
@@io:jfs2:buf_map:*:/usr|/tmp
```

For Virtual file system (VFS) probes, the fifth tuple can have following values:

<i>Table 42. VFS filter tuple</i>	
Filter (fifth Tuple)	Description
File system mount path. For example, /tmp	The probe action is run for files that belong to the file system.
File system type. The allowed symbols are JFS2, NAMEFS, NFS, JFS, CDROM, PROCFS, SFS, CACHEFS, NFS3, AUTOFS, POOLFS, VXFS, VXODM, UDF, NFS4, RFS4, CIFS, PMEMFS, AHAFS, STNFS, ASMFS	The probe action is run for files of the particular file system. The symbols correspond to the AIX file systems defined in the exported header file sys/vmount.h.

Consider the following examples for the Virtual file system (VFS) probes:

```
@@io:vfs:entry:read:JFS2
```

```
@@io:vfs:exit:*:/usr|JFS
```

I/O probe related built-in variables for Vue scripts

`__iobuf` built-in variable

You can use the special `__iobuf` built-in variable to access various information about the I/O buffer that is employed in the current I/O operation. It is accessible in probes of sub types: `disk`, `lvm`, and `bdev`. Its member elements can be accessed by using the `__iobuf->member` syntax.

Note: Whenever the actual value cannot be obtained, the value that is marked as `Invalid Value` is returned. This value is returned because of one of the following reasons:

- Page fault context is required, but the current `probevctrl` tunable value, `num_pagefaults`, is either 0 or not sufficient.
- The memory location that is containing the value is paged out.
- Any other severe system error such as invalid pointer or corrupted memory.

`__iobuf` built-in variable has the following members:

Table 43. The `__iobuf` built-in variable members

Member name	Type	Description	Invalid Value
blknum	unsigned long long	Starting block number of the I/O request.	0xFFFFFFFFFFFFFFFF
bcount	unsigned long long	Requested number of bytes in the I/O operation.	0xFFFFFFFFFFFFFFFF
bflags	unsigned long long	The flags that are associated with the I/O operation. The following symbols are available: <code>B_READ</code> , <code>B_ASYNC</code> , <code>B_ERROR</code> . The symbols can be used along with the <code>bflags</code> value to see whether it is set. For example, if <code>(__iobuf->bflags & B_READ)</code> is true, then it is a read operation. Note: There is no <code>B_WRITE</code> flag. If the <code>B_READ</code> flag is not set, it is considered to be write operation.	0
devnum	unsigned long long	The device number of the target device that is associated with the I/O operation. It has the device major number and minor number that is embedded in it.	0
major_num	int	The major number of the target device of the I/O operation.	-1
minor_num	int	The minor number of the target device of the I/O operation.	-1
error	int	In case of any error in the I/O operation, this value is the error number. This value is defined in the exported <code>errno.h</code> header file.	-1

Table 43. The `__iobuf` built-in variable members (continued)

Member name	Type	Description	Invalid Value
residue	unsigned long long	The remaining number of bytes from the original request that might not be read or written. On the I/O completion events, this value is ideally zero. But for read operation, a nonzero value might mean that you are trying to read more than what is available, which is acceptable. This value is considered only when error value is nonzero.	0xFFFFFFFFFFFFFFFF
bufid	unsigned long long	A unique number that is associated with the I/O request. While the I/O is in progress, the bufid value uniquely identifies the I/O request in all the events of a particular sub type. For example, in <code>disk: entry</code> , <code>disk: iostart</code> , <code>disk: iodone</code> , and <code>disk: exit</code> . If the <code>__iobuf->bufid</code> matches, it is the same I/O request at various stages).	0
parent_bufid	unsigned long long	If the value is not 0, this value provides the bufid of the upper layer buffer that is associated with this I/O request. You can now link the current I/O operation with the upper layer I/O request. For example, in a disk I/O request, the corresponding LVM I/O can be determined. Note: The <code>parent_bufid</code> field is not set in all code paths, and hence it is not always useful. Use the <code>child_bufid</code> field to link I/O requests between two adjacent layers.	0

Table 43. The `__iobuf` built-in variable members (continued)

Member name	Type	Description	Invalid Value
<code>child_bufid</code>	unsigned long long	If the value is not 0, this value provides the <code>bufid</code> of the new I/O request that is sent to the lower layer. The best events to record are <code>disk:iostart</code> , <code>lvm:iostart</code> , and <code>bdev:iostart</code> . You can identify the I/O in the lower adjacent layer by matching the <code>__iobuf->bufid</code> value to this <code>child_bufid</code> value. For example, in <code>lvm:iostart</code> , you can record the <code>__iobuf->child_bufid</code> value. Then, in <code>disk:entry</code> , you can match it with <code>__iobuf->bufid</code> to identify the corresponding I/O request.	0

__file built-in variable

You can use the `__file` special built-in variable to get various information about file operation. It is available in probes of sub type VFS. Its member elements can be accessed by using the `__file->member` syntax.

Note: Whenever the actual value cannot be obtained, the value that is marked as invalid is returned. The invalid value is returned because of one of the following reasons:

- Page fault context is required, but the current `probevctrl` tunable value `num_pagefaults` is either 0 or not sufficient.
- The memory location, which contains the value, is paged out.
- Any other severe system error such as invalid pointer, or corrupted memory.

The `__file` built-in variable has the following members:

Table 44. The `__file` built-in variable members

Member name	Type	Description	Invalid Value
f_type	int	<p>Specifies the type of the file. It can match one of the following built-in constant values:</p> <ul style="list-style-type: none"> • F_REG (regular file) • F_DIR (directory) • F_BLK (block device file) • F_CHR (character device file) • F_LNK (file link) • F SOCK (socket) <p>Note: The value might not match any of the built-in constants because the list does not include every possible file type, but only the most useful ones.</p>	-1

Table 44. The `__file` built-in variable members (continued)

Member name	Type	Description	Invalid Value
<code>fs_type</code>	<code>int</code>	<p>Specifies the type of the file system to which this file belongs. It can match one of the following built-in constant values:</p> <ul style="list-style-type: none"> • <code>FS_JFS2</code> • <code>FS_NAMEFS</code> • <code>FS_NFS</code> • <code>FS_JFS</code> • <code>FS_CDROM</code> • <code>FS_PROCFS</code> • <code>FS_SFS</code> • <code>FS_CACHEFS</code> • <code>FS_NFS3</code> • <code>FS_AUTOFS</code> • <code>FS_POOLFS</code> • <code>FS_VXFS</code> • <code>FS_VXODM</code> • <code>FS_UDF</code> • <code>FS_NFS4</code> • <code>FS_RFS4</code> • <code>FS_CIFS</code> • <code>FS_PMEMFS</code> • <code>FS_AHAFS</code> • <code>FS_STNFS</code> • <code>FS_ASMFS</code> <p>The built-in constants corresponds to the AIX file system types defined in the exported <code>sys/vmount.h</code> header file.</p>	-1
<code>mount_path</code>	<code>char *</code>	Specifies the path where the associated file system is mounted.	null string
<code>devnum</code>	<code>unsigned long long</code>	Specifies the device number of the associated block device of the file. Both the major and minor numbers are embedded in it. If there is no associated block device, then it is 0.	0

Table 44. The `__file` built-in variable members (continued)

Member name	Type	Description	Invalid Value
major_num	int	Specifies the major number of the associated block device of the file.	-1
minor_num	int	Specifies the minor number of the associated block device of the file.	-1
offset	unsigned long long	Specifies the current read/write byte offset of the file.	0xFFFFFFFFFFFFFFFF
rw_mode	int	Specifies the read/write mode of the file. It matches one of the built-in constant values: <code>F_READ</code> or <code>F_WRITE</code> .	-1
byte_count	unsigned long long	At <code>vfs: entry</code> event, <code>byte_count</code> provides the byte count of the read or write request. At <code>vfs: exit</code> event, it provides the number of bytes that remained unfulfilled. For example, the difference of this value between these two events determines how many bytes were processed in the operation.	0xFFFFFFFFFFFFFFFF
fname	char *	Specifies the name of the file (only base name, not path).	null string
inode_id	unsigned long long	Specifies a system-wide unique number that is associated with the file. Note: It is different from file inode number.	0
path	path_t (new data type in VUE)	Specifies the complete file path. It can be printed by using <code>printf()</code> and the format specifier <code>%p</code> .	null string as file path

Table 44. The `__file` built-in variable members (continued)

Member name	Type	Description	Invalid Value
error	int	If the read/write operation failed, the error number as defined in the exported <code>errno.h</code> header file. If there is no error, it is 0.	-1

`__lvol` built-in variable

You can use the `__lvol` special built-in variable to get various information about the logical volume in an LVM operation. It is available in probes of sub type `lvm`. Its member elements can be accessed by using the `__lvol->member` syntax.

Note: Whenever the actual value cannot be obtained, the value, which is marked as Invalid Value, is returned. There might be following reasons for getting this invalid value:

- Page fault context is required, but the current `probevctrl` tunable value `num_pagefaults` is either 0 or not sufficient.
- The memory location that contains the value is paged out.
- Any other severe system error such as invalid pointer or corrupted memory.

`__lvol` built variable in has following members:

Table 45. The `__lvol` built-in variable members

Member name	Type	Description	Invalid Value
name	char *	The name of the logical volume.	null string
devnum	unsigned long long	The device number of the logical volume. It has both major number and minor number that is embedded in it.	0
major_num	int	The major number of the logical volume.	-1
minor_num	int	The minor number of the logical volume.	-1

Table 45. The `__lvol` built-in variable members (continued)

Member name	Type	Description	Invalid Value
<code>lv_options</code>	unsigned int	<p>The options that are related to the logical volume. The following values are defined as built-in constants:</p> <ul style="list-style-type: none"> • <code>LV_RDONLY</code> (read-only logical volume) • <code>LV_NOMWC</code> (no mirror write consistency checking) • <code>LV_ACTIVE_MWC</code> (active mirror write consistency) • <code>LV_PASSIVE_MWC</code> (passive mirror write consistency) • <code>LV_SERIALIZE_IO</code> (I/O is serialized) • <code>LV_DUMPDEV</code> (This LV is a dump device) <p>You can check whether one of these values is set by having condition such as <code>__lvol->lv_options & LV_RDONLY</code>.</p> <p>Note: All possible values are not defined, and hence other options might be available in the value.</p>	0xFFFFFFFF

__volgrp built-in variable

You can use `__volgrp` special built-in variable to get various information about the volume group in an LVM operation. It is available in probes of sub type `lvm`. Its member elements can be accessed by using the `__volgrp->member` syntax.

Note: Whenever the actual value cannot be obtained, the value that is marked as `Invalid Value` is returned. The value could be invalid because of the following reasons:

- Page fault context is required, but the current `probevctrl` tunable value `num_pagefaults` is either 0 or not sufficient.
- The memory location that contains the value is paged out.
- Any other severe system error such as invalid pointer or corrupted memory.

`__volgrp` built-in variable has following members:

Table 46. The __volgrp built-in variable members

Member name	Type	Description	Invalid Value
name	char *	The name of the volume group.	null string
devnum	unsigned long long	The device number of the volume group. It has major number and minor number that is embedded in it.	0
major_num	int	The major number of the volume group.	-1
minor_num	int	The minor number of the volume group. Note: For volume group, AIX always assigns 0 as the minor number.	-1
num_open_lvs	int	The number of open logical volumes that belong to this volume group.	-1

__diskinfo built-in variable

You can use the __diskinfo special built-in variable to get various information about the disk in a disk I/O operation. It is available in probes of sub type disk. Its member elements can be accessed by using the __diskinfo->member syntax.

Note: Whenever the actual value cannot be obtained, the value that is marked as “Invalid Value” is returned. There might be following reasons for getting this value:

- Page fault context is required, but the current probevctrl tunable value num_pagefaults is either 0 or not sufficient.
- The memory location that contains the value is paged out.
- Any other severe system error such as invalid pointer or corrupted memory.

__diskinfo built-in variable has following members:

Table 47. The __diskinfo built-in variable members

Member name	Type	Description	Invalid Value
name	char *	The name of the disk.	null string.
devnum	unsigned long long	The device number of the disk. It has major number and minor number that are embedded in it.	0
major_num	int	The major number of the disk.	-1
minor_num	int	The minor number of the disk.	-1
lun_id	unsigned long long	The Logical Unit Number (LUN) for the disk.	0xFFFFFFFFFFFFFFFF

Table 47. The `__diskinfo` built-in variable members (continued)

Member name	Type	Description	Invalid Value
transport_type	int	The transport type of the disk. It can match one of the following built-in constant values: <ul style="list-style-type: none"> • T_FC (Fibre Channel) • T_ISCSI (iSCSI) • T_VSCSI (Virtual SCSI) • T_SAS (Serial Attached SCSI) 	-1
queue_depth	int	The queue depth of the disk. It indicates how many maximum simultaneous I/O requests that the disk driver can pass on to the lower layer (for example, adapter). If the number of incoming I/O requests is more than queue_depth, the request is handled differently. The extra request is handled by the disk driver in its wait queue until lower layer responds to at least one of the outstanding I/O requests.	-1
cmds_out	int	Number of outstanding I/O command requests to the lower layer (for example, adapter).	-1
path_count	int	Number of MPIIO paths of the disk (Only if the disk is MPIIO capable, else it is 0).	-1

Table 47. The __diskinfo built-in variable members (continued)

Member name	Type	Description	Invalid Value
reserve_policy	int	<p>The SCSI reservation policy of the disk. It matches one of the following built-in constant values:</p> <ul style="list-style-type: none"> • DK_NO_RESERVE (no_reserve) • DK_SINGLE_PATH (single_path) • DK_PR_EXCLUSIVE (PR_exclusive) • DK_PR_SHARED (PR_shared) <p>Refer to AIX MPIO documentation to know more about the reservation policies.</p>	-1
scsi_flags	int	<p>The SCSI flags of the disk. The following built-in flag values are defined:</p> <ul style="list-style-type: none"> • SC_AUTOSENSE_ENABLED (On error, target sends sense data in the response. Initiator needs not send request sense command.) • SC_NACA_1_ENABLED (Normal ACA is enabled and the target goes to ACA state if it is returning check condition.) • SC_64BIT_IDS (64-bit SCSI ID and logical unit number(LUN)) • SC_LUN_RESET_ENABLED (LUN reset command can be sent.) • SC_PRIORITY_SUP (Device supports I/O priority.) 	0

Table 47. The `__diskinfo` built-in variable members (continued)

Member name	Type	Description	Invalid Value
		<ul style="list-style-type: none"> • <code>SC_CACHE_HINT_SUP</code> (Device supports cache hints.) • <code>SC_QUEUE_UNTAGGED</code> (Device supports queuing of untagged commands.) <p>Note: All flag values are not defined, hence other flags present might be available in the value.</p>	0

`__diskcmd` built-in variable

You can use the `__diskcmd` special built-in variable to get various information about the SCSI I/O command for the current operation. It is available in probes of sub type disk (but only `iostart` and `iodone` events). Its member elements can be accessed by using syntax `__diskcmd->member`.

Note: Whenever the actual value cannot be obtained, the value that is marked as “Invalid Value” is returned. There might be following reasons for getting value:

- Page fault context is required, but the current `probevctrl` tunable value `num_pagefaults` is either 0 or not sufficient.
- The memory location that contains the value is paged out.
- Any other severe system error such as invalid pointer or corrupted memory.

`__diskcmd` built-in variable has following members:

Table 48. The `__diskcmd` built-in variable members

Member name	Type	Description
cmd_type	int	<p>The type of the SCSI command (both type and subtype are merged together). The following built-in constant values are available as command type:</p> <ul style="list-style-type: none"> • DK_BUF (normal I/O read/write) • DK_IOCTL (ioctl) • DK_REQSNS (Request sense) • DK_TGT_LUN_RST (target or LUN reset) • DK_TUR (Test unit ready) • DK_INQUIRY (Inquiry) • DK_RESERVE (SCSI-2 RESERVE, 6-byte version) • DK_RELEASE (SCSI-2 RELEASE, 6-byte version) • DK_RESERVE_10 (SCSI-2 RESERVE, 10-byte version) • DK_RELEASE_10 (SCSI-2 RELEASE, 10-byte version) • DK_PR_RESERVE (SCSI-3 Persistent Reserve, RESERVE) • DK_PR_RELEASE (SCSI-3 Persistent Reserve, RELEASE) • DK_PR_CLEAR (SCSI-3 Persistent Reserve, CLEAR)
		<ul style="list-style-type: none"> • DK_PR_PREEMPT (SCSI-3 Persistent Reserve, PREEMPT) • DK_PR_PREEMPT_ABORT (SCSI-3 Persistent Reserve, PREEMPT AND ABORT) • DK_READCAP (READ CAPACITY, 10-byte version) • DK_READCAP16 (READ CAPACITY, 16-byte version) <p>Note: The built-in constants are bit position values and hence their presence must be checked by using '&' operator (the '==' operator must not be used). For example: <code>__diskcmd->cmd_type & DK_IOCTL.</code></p>

Table 48. The `__diskcmd` built-in variable members (continued)

Member name	Type	Description
retry_count	int	It indicates whether the I/O command is retried after any failure. Note: The value of 1 means that it is the first attempt. Any larger value indicates actual retrials.
path_switch_count	int	It indicates how many times the path was changed for this particular I/O operation (usually indicates some I/O path failure, either transient or permanent).
status_validity	int	In case of any error, this value indicates whether it is a SCSI error or adapter error. It can match one of the following built-in constant values: <code>SC SCSI_ERROR</code> or <code>SC_ADAPTER_ERROR</code> . If there is no error, then it is 0.

Table 48. The `__diskcmd` built-in variable members (continued)

Member name	Type	Description
scsi_status	int	<p>If the <code>status_validity</code> field is set to <code>SC_SCSI_ERROR</code>, this field gives more details about the error. It can match one of the built-in constant values:</p> <ul style="list-style-type: none"> • <code>SC_GOOD_STATUS</code> (Task is completed successfully) • <code>SC_CHECK_CONDITION</code> (Some error, sense data provides more information) • <code>SC_BUSY_STATUS</code> (LUN is busy, cannot accept command) • <code>SC_RESERVATION_CONFLICT</code> (Violation of existing SCSI reservation.) • <code>SC_COMMAND_TERMINATED</code> (The device ended the command.) • <code>SC_QUEUE_FULL</code> (The device queue is full.) • <code>SC_ACA_ACTIVE</code> (The device is in Auto Contingent Allegiance state.) • <code>SC_TASK_ABORTED</code> (The device stopped the command.) <p>Note: All possible values are not defined. Hence, <code>SC_SCSI_ERROR</code> can have a value that might not match any of the built-in values. You can look up the corresponding SCSI command response code.</p>

Table 48. The `__diskcmd` built-in variable members (continued)

Member name	Type	Description
<code>adapter_status</code>	int	<p>If the <code>status_validity</code> field is set to <code>SC_ADAPTER_ERROR</code>, this field provides more information about the error. It can match one of the following built-in constant values:</p> <ul style="list-style-type: none"> • <code>ADAP_HOST_IO_BUS_ERR</code> (Host I/O bus error) • <code>ADAP_TRANSPORT_FAULT</code> (transport layer error) • <code>ADAP_CMD_TIMEOUT</code> (I/O command was timed out) • <code>ADAP_NO_DEVICE_RESPONSE</code> (no response from the device) • <code>ADAP_HDW_FAILURE</code> (adapter hardware failure) • <code>ADAP_SFW_FAILURE</code> (adapter microcode failure) • <code>ADAP_TRANSPORT_RESET</code> (adapter detected an external SCSI bus reset) • <code>ADAP_TRANSPORT_BUSY</code> (transport layer is busy) • <code>ADAP_TRANSPORT_DEAD</code> (transport layer is inoperative) • <code>ADAP_TRANSPORT_MIGRATED</code> (transport layer is migrated) • <code>ADAP_FUSE_OR_TERMINAL_PWR</code> (adapter blown fuse or bad electrical termination)

`__iopath` built-in variable

You can use the `__iopath` special built-in variable to get various information about the I/O path for the current operation. It is available in probes of sub type `disk` for `iostart` and `iodone` events only. Its member elements can be accessed by using the `__iopath->member` syntax.

Note: Whenever the actual value cannot be obtained, the value, which is marked as `Invalid Value`, is returned. There might be following reasons for getting this value:

- Page fault context is required, but the current `probevctrl` tunable value `num_pagefaults` is either 0 or not sufficient.
- The memory location that contains the value is paged out.
- Any other severe system error such as invalid pointer or corrupted memory.

`__iopath` has following members:

Member name	Type	Description	Invalid Value
<code>path_id</code>	int	The ID of the current path (starting from 0).	-1
<code>scsi_id</code>	unsigned long long	The SCSI ID of the target on this path.	0xFFFFFFFFFFFFFFFF
<code>lun_id</code>	unsigned long long	The Logical Unit Number (LUN) on this path.	0xFFFFFFFFFFFFFFFF
<code>ww_name</code>	unsigned long long	The worldwide name of the target port on this path.	0
<code>cmds_out</code>	int	The number of I/O commands outstanding on this path.	-1

`__j2info` built-in variable

The `__j2info` is a special built-in variable that you can use to get various information about JFS2 file system operation. It is available in probes of sub type `jfs2`. Its member elements can be accessed by using the `__j2info->member` syntax.

Note: Whenever the actual value cannot be obtained, the value, which is marked as `Invalid Value` is returned. There might be following reasons for getting this value:

- Page fault context is required, but the current `probevctrl` tunable value `num_pagefaults` is either 0 or not sufficient.
- The memory location that contains the value is paged out.
- Any other severe system error such as invalid pointer or corrupted memory.

`__j2info` has the following members:

Member name	Type	Description	Invalid Value
<code>inode_id</code>	unsigned long long	A system-wide unique number that is associated with the file of current operation. Note: It is different from the file inode number.	0
<code>f_type</code>	int	Type of the file. The <code>__file->f_type</code> description provides possible values.	-1
<code>mount_path</code>	char *	The path where the file system is mounted.	null string.
<code>devnum</code>	unsigned long long	The device number of the underlying block device of the file system. It has both major number and minor number embedded.	0

Table 50. The `__j2info` built-in variable members (continued)

Member name	Type	Description	Invalid Value
major_num	int	The major number of the underlying block device of the file system.	-1
minor_num	int	The minor number of the underlying block device of the file system.	-1
l_blknum	unsigned long long	The logical block number for this file operation.	0xFFFFFFFFFFFFFFFF
l_bcount	unsigned long long	The requested byte count between the logical blocks in this operation.	0xFFFFFFFFFFFFFFFF
child_bufid	unsigned long long	The bufid of the I/O request buffer that is sent down to the lower layer (for example, LVM). In that layer, it appears as <code>__iobuf->bufid</code> .	0
child_blknum	unsigned long long	The block number of the I/O request buffer that is sent down to the lower layer (for example, LVM). In that layer, it appears as <code>__iobuf->blknum</code> .	0xFFFFFFFFFFFFFFFF
child_bcount	unsigned long long	The byte count of the I/O request buffer that is sent down to the lower layer (for example, LVM). In that layer, it appears as <code>__iobuf->bcount</code> .	0xFFFFFFFFFFFFFFFF
child_bflags	unsigned long long	The flags of the I/O request buffer that is sent down to the lower layer (for example, LVM). In that layer, it appears as <code>__iobuf->bflags</code> .	0

Example scripts for I/O probe manager

1. Script to trace any write operation to the `/etc/passwd` file:

```
int write(int, char *, int);
@@BEGIN {
    target_inodeid = fpath_inodeid("/etc/passwd");
}
@@syscall:*:write:entry {
    if (fd_inodeid(__arg1) == target_inodeid) {
        printf("write on /etc/passwd: timestamp=%A, pid=%lld, pname=[%s],
uid=%lld\n",
                timestamp(), __pid, __pname, __uid);
    }
}
```

```

If the scripts is in a VUE file, names etc_passwd.e. The script can be run as:
# probevue etc_passwd.e
In another terminal, if the user (root) runs:
# mkuser user1
Then probevue displays an output similar to the following example:
write on /etc/passwd: timestamp=Mar/03/15 16:10:07, pid=14221508, pname=[mkuser], uid=0

```

2. Script to find the maximum and minimum I/O operation time for a disk (for example, hdisk0) in a period. Also, find the block number, requested byte count, time of operation and type of operation (read or write) corresponding to the maximum or minimum time.

```

long long min_time, max_time;
@@BEGIN {
    min_time = max_time = 0;
}
@@io:disk:entry:*:hdisk0 {
    ts_entry[__iobuf->bufid] = (long long)timestamp();
}
@@io:disk:exit:*:hdisk0 {
    if (ts_entry[__iobuf->bufid]) { /* only if we recorded entry time */
        ts_now = timestamp();
        op_type = (__iobuf->bflags & B_READ) ? "READ" : "WRITE";
        dt = (long long)diff_time(ts_entry[__iobuf->bufid], ts_now, MICROSECONDS);
        if (min_time == 0 || dt < min_time) {
            min_time = dt;
            min_blknum = __iobuf->blknum;
            min_bcount = __iobuf->bcount;
            min_ts = ts_now;
            min_optype = op_type;
        }
        if (max_time == 0 || dt > max_time) {
            max_time = dt;
            max_blknum = __iobuf->blknum;
            max_bcount = __iobuf->bcount;
            max_ts = ts_now;
            max_optype = op_type;
        }
        ts_entry[__iobuf->bufid] = 0;
    }
}
@@END {
    printf("Maximum and minimum IO operation time for [hdisk0]:\n");
    printf("Max: %lld usec, block=%lld, byte count=%lld, operation=%s, time of
operation=[%A]\n",
        max_time, max_blknum, max_bcount, max_optype, max_ts);
    printf("Min: %lld usec, block=%lld, byte count=%lld, operation=%s, time of
operation=[%A]\n",
        min_time, min_blknum, min_bcount, min_optype, min_ts);
}

```

```

Let this script be in a VUE file named disk_min_max_time.e. It can be executed as:
# probevue disk_min_max_time.e
Let there be some IO activity on hdisk0 (dd command can be used).
Then after a few minutes, if the above command is terminated (by pressing CTRL-C), then it will print output similar to:
^CMaximum and minimum IO operation time for [hdisk0]:
Max: 48174 usec, block=6927976, byte count=4096, operation=READ, time of operation=[Mar/04/15 03:31:07]
Min: 133 usec, block=6843288, byte count=4096, operation=READ, time of operation=[Mar/04/15 03:31:03]

```

Network probe manager

Network probe manager tracks incoming and outgoing network packets in a system (packet information as interpret by the bpf module in AIX). Probe specification allows the user to specify Berkeley Packet Filter (BPF) filters, similar to tcpdump filter expression for granular tracking.

You can use built-in variables to collect packet header and payload information for Internet protocols. For example, Ethernet, Internet Protocol Version 4/Version 6 (IPv4/v6), Transmission Control Protocol (TCP), User Datagram Protocol (UDP), Internet Control Message Protocol (ICMP), Internet Group Message Protocol (IGMP), and Address Resolution Protocol (ARP) protocols.

Network probe manager reports critical protocol-specific events (TCP state changes, round-trip times, retransmissions, UDP buffer overflows).

The network probe manager addresses following primary use cases:

- Provide the following packet-specific information according to the bpf module based on IP address and ports:
 - Track the incoming and outgoing bytes for a connection.
 - Use following built-ins to gather protocol header and payload information.
 - TCP flags (SYN, FIN), TCP sequence and acknowledgment number.
 - IPv4/IPv6 (IP addresses, protocol types: tcp, udp, icmp, igmp, and so on).
 - ICMP (packet type: ECHO REQUEST, ECHO RESPONSE, and so on).
- Provide access to complete RAW network packet for probe script processing.
- Report the following protocol-related events:
 - Track TCP sender and receiver buffer full events.
 - TCP connection state changes from SYN-SENT state to ESTABLISHED state or from ESTABLISHED state to CLOSE state.
 - Monitor delta time between state changes (for example, time that is taken from SYN-SENT state to ESTABLISHED state).
 - Identify the listener (connection information) that discarded connections because the listener's queue is full.
 - Identify retransmissions (second and further retransmission for a packet) for TCP connections.
 - Identify the UDP socket that dropped packets because of insufficient receiving buffer.

Probe specification

Probe specification for network probe manager contains three or five tuples that are separated by : (colon). First tuple is always @@net.

Network probe manager supports two major categories of specifications: One category gathers packet-specific information and another category gathers protocol-specific information.

- Format to gather packet specific information:

@@net:bpf:<interface1>|<interface 2>|... :<protocol>:<Filter>

- Format to gather protocol specific information

@@net:tcp:<event_name>

@@net:udp:<event_name>

Probe sub type

The second tuple signifies the sub type of the probe that indicates which layer of AIX network stack contains the probe. This tuple can have one of the following values (it cannot be *):

Second Tuple (sub type)	Description
bpf	This probe starts at network interface layer when a packet matches the specific filter.
tcp	This probe starts for TCP protocol-specific events.
udp	This probe starts for UDP protocol-specific events.

Probe network event or gather network packet information

The third tuple is specific to particular sub type (specified in second tuple). It cannot have a value of *.

bpf-based probes

The specification contains 5 tuples for bpf-based probes that are described in the following table:

Table 52. bpf-based probes: Tuple specification

Second tuple (Sub type)	Subsequent tuples	Description
bpf	Third tuple: interface names	This tuple specifies an interface or a list of interfaces for which the packet information can be captured. Possible values are enX (for example, en0,en1) and lo0. The * value is not supported for this tuple. You can specify one or more interfaces at a time by using as delimiter.
	Fourth tuple: protocol	<p>This tuple specifies the network protocol to start the probe. Possible values are ether, arp, rarp, ipv4, ipv6, tcp, udp, icmp4, icmp6 and igmp. Protocol-specific built-ins are populated for access in Vue script. For example, a protocol value of ipv4 populates __ip4hdr built-ins.</p> <p>The * value for this tuple indicates that the probe starts for all protocol types that match the specified filter. When the protocol is *, none of the built-in values that are supported by network probe manager are available to Vue scripts. You can access the raw packet data of requested size by using the Vue function copy_kdata () and map to corresponding protocol headers.</p> <p>Note: Specifying * as a value can be performance intensive as the probe is started for all incoming and outgoing packets on the specified interfaces that match the filter. There are also copies involved when the packet information is spanned across multiple packet buffers.</p>
bpf	Fifth tuple: bpf filter string	This tuple specifies the bpf filter expression (filter expressions as specified in tcpdump command). Filter expression must be provided in the double quotation marks. Filter expression and protocol that is specified in the fourth tuple must be compatible. The * value is not supported in this tuple.

Examples

1. Specification format to access the built-in variables that are related to Ethernet header (`__etherhdr`), IP header (`__ip4hdr` or `__ip6hdr`), and TCP header (`__tcphdr`) information from the Vue script when interface `en0` receives or sends packet on port 23 (filter string "port 23"):

```
@@net:bpf:en0:tcp:"port 23"
```

2. Specification format to access the built-in variables related to Ethernet header (`__etherhdr`), IP header (`__ip4hdr` or `__ip6hdr`), and UDP header (`__udphdr`) information from the Vue script when system receives or sends packet from host `example.com` (filter string "example.com") on `en0` and `en1` interfaces:

```
@@net:bpf:en0|en1:udp:"host example.com"
```

3. Specification format to access the raw packet information when system receives or sends packet from or to "host `example.com`":

```
@@net:bpf:en1:*:"host example.com"
```

Note: Each bpf probe specification uses a bpf device. These devices are shared by ProbeVue, `tcpdump`, and any other application that uses the `libpcap` or `bpf` services for packet capture and injection. The number of bpf probes depends on the number of available bpf devices in the system.

When a bpf probe is started, the `__mdata` variable contains the raw packet data. You can access the raw data of requested size by using the Vue function `copy_kdata()` and map to the `ether_header`, `ip_header`, and so on. Use the following structures to find out the header and payload data information.

Example

VUE script to access the raw packet data when the "*" is specified as the protocol.

```
/* Define the ether header structure */
struct ether_header {
    char ether_dhost[6];
    char ether_shost[6];
    short ether_type;
};

/* ProbeVue script to access and interpret the data from RAW packet */
@@net:bpf:en0:*:"port 23"
{
    /* define the script local variables */
    __auto struct ether_header eth;
    __auto char *mb;

    /* __mdata contains the address of packet data */
    mb =(char *) __mdata;
    printf("Network probevue\n");

    /*
    * Use already available "copy_kdata(...)" VUE function to copy data of
    * requested size (size of ether_header) from mbuf data pointer to eth
    * (ether_header) variable.
    */
    copy_kdata (mb, eth);
    printf("Ether Type from raw data :%x\n",eth.ether_type);
}
}
```

TCP probes

The specification contains three tuples for TCP probes as described in the following table:

Table 53. TCP probes: Tuple specification

Second tuple (Sub type)	Events (Third tuple) The * value is not supported in this tuple.	Description
tcp	state_change	This probe is started whenever the TCP state changes.
	send_buf_full	This probe is started whenever the send buffer full event occurs.
	recv_buf_full	This probe is started whenever the receive buffer full event occurs.
	retransmit	This probe is started whenever the re-transmission of packet happens for TCP connection.
	listen_q_full	This probe is started whenever a server (listener socket) discards the new connection requests due to listener's queue being full.

`__proto_info` built-in variable provides the TCP connection (four tuple) information (local IP, remote IP, local port, and remote port) whenever the TCP-related event occurs. Remote port and IP address contains a value of NULL for the `listen_q_full` event.

Example

Probe specifications for TCP protocol state changes:

```
@@net:tcp:state_change
```

udp probes

For udp probes the specification contains three tuples as described in the following table:

Table 54. udp second tuple: Third tuple values

Second tuple (Sub type)	Events (third tuple) The * value is not supported in this tuple.	Description
udp	sock_recv_buf_overflow	This probe is started whenever the datagram or the UDP socket's receive buffer overflows.

The `__proto_info` built-in variable provides the UDP protocol related data (source IP and destination IP addresses, source and destination port numbers) whenever socket receive buffer overflow event occurs.

```
@@net:udp:sock_recv_buf_overflow
```

Example

Probe specifications for UDP socket's receive buffer overflow:

```
@@net:udp:sock_recv_buf_overflow
```

Network probe-related built-in variables for Vue scripts

Network related events can be probed using following built-in variables.

__etherhdr built-in variable

The `__etherhdr` variable is a special built-in variable to get ether header information from filtered packet. This built-in variable is available when you probe the packet information at interface layer with any one of these protocols: “ether”, “ipv4”, “ipv6”, “tcp”, “udp”, “icmp4”, “icmp6”, “igmp”, “arp”, and “rarp”. This variable is available in probes of sub type bpf. Its member elements can be accessed by using the syntax `__etherhdr->member`.

The `__etherhdr` built-in value has the following members:

<i>Table 55. The <code>__etherhdr</code> built-in variable members</i>		
Member name	Type	Description
src_addr	mac_addr_t	Source MAC address. The data type <code>mac_addr_t</code> is used to store the MAC address. Use format specifier “M” to print the MAC address.
dst_addr	mac_addr_t	Destination MAC address. The data type <code>mac_addr_t</code> is used to store the MAC address. Use format specifier “M” to print the MAC address.
ether_type	unsigned short	This name indicates the protocol encapsulated in the payload of an Ethernet frame. Protocols can be IPv4, IPv6, ARP, and REVARP. It can match one of the following built-in constant values for <code>ether_type</code> : <ul style="list-style-type: none"> • <code>ETHERTYPE_IP</code> • <code>ETHERTYPE_IPV6</code> • <code>ETHERTYPE_ARP</code> • <code>ETHERTYPE_REVARP</code> Refer the header files <code>/usr/include/netinet/if_ether.h</code> and <code>/usr/include/netinet/if_ether6.h</code> for <code>ether_type</code> values.

Note: The `__etherhdr` built-in variable is applicable only for Ethernet interfaces and not for loopback interfaces.

__ip4hdr built-in variable

The `__ip4hdr` variable is a special built-in variable to get the IPv4 header information from filtered packet. This variable is available when you probe the packet information at interface layer with any one of the protocols: “ipv4”, “tcp”, “udp”, “icmp4”, and “igmp”. And, it has valid data when IP version is IPv4. This variable is available in probes of sub type bpf. Its member elements can be accessed by using the syntax `__ip4hdr->member`.

This built-in variable has the following members:

Table 56. The `__ip4hdr` built-in variable members

Member name	Type	Description
<code>src_addr</code>	<code>ip_addr_t</code>	Source IP address. The data type <code>ip_addr_t</code> is used to store the IP address. Use format specifier "I" to print the IP address in dotted decimal format and use format specifier "H" to print the host name. Host name printing is a costly operation.
<code>dst_addr</code>	<code>ip_addr_t</code>	Destination IP address. The data type <code>ip_addr_t</code> is used to store the IP address. Use format specifier "I" to print the IP address in dotted decimal format and use format specifier "H" to print host name. Host name printing is a costly operation.
<code>protocol</code>	unsigned short	This member name indicates the protocol that is used in the data portion of the IP datagram. Protocols can be TCP, UDP, ICMP, IGMP, FRAGMENTED, and so on. It can match one of the following built-in constant values for protocol. <pre> IPPROTO_HOPOPTS, IPPROTO_ICMP, IPPROTO_IGMP, IPPROTO_TCP, IPPROTO_UDP, IPPROTO_ROUTING, IPPROTO_FRAGMENT, IPPROTO_NONE, IPPROTO_LOCAL </pre> Refer the header file <code>/usr/include/netinet/in.h</code> for protocol values.
<code>ttl</code>	unsigned short	Time to live or hop limit.
<code>cksum</code>	unsigned short	IP header checksum.
<code>id</code>	unsigned short	Identification number. This member is used for uniquely identifying the group of fragments of a single IP datagram.
<code>total_len</code>	unsigned short	Total length. This value is entire packet (fragment) size, including IP header and data in bytes.
<code>hdr_len</code>	unsigned short	Size of the IP header.

Table 56. The `__ip4hdr` built-in variable members (continued)

Member name	Type	Description
tos	unsigned short	Type of service.
frag_offset	unsigned short	<p>Fragment offset.</p> <p>This value specifies the offset of particular fragment, relative to beginning of the original unfragmented IP datagram. The first fragment has an offset of zero.</p> <p>It can match one of the built-in constant <code>frag_offset</code> flag values. The flag values must be bitwise and with the built-in constant flag value to validate the presence of the particular flag:</p> <ul style="list-style-type: none"> • <code>IP_DF</code> (No fragment flag) • <code>IP_MF</code> (more fragments flag) <p>Refer the header file <code>/usr/include/netinet/ip.h</code> for flag values.</p>

`__ip6hdr` built-in variable

The `__ip6hdr` variable is a special built-in variable to get the IPv6 header information from filtered packet. This variable is available when user probes the packet information at interface layer. This variable with any one of the protocols (“`ipv6`”, “`tcp`”, “`udp`” and “`icmp6`”) has valid data when IP version is IPv6. This variable is available in probes of sub type `bpf`. Its member elements can be accessed by using the syntax `__ip6hdr->member`.

This built-in variable has the following members:

Table 57. The `__ip6hdr` built-in variable members

Member name	Type	Description
src_addr	ip_addr_t	<p>Source IP address.</p> <p>The data type <code>ip_addr_t</code> is used to store the IP address. Use format specifier “<code>I</code>” to print the IP address and use format specifier “<code>H</code>” to print the host name. Host name printing is a costly operation.</p>
dst_addr	ip_addr_t	<p>Destination IP address.</p> <p>The data type <code>ip_addr_t</code> is used to store the IP address. Use format specifier “<code>I</code>” to print the IP address and use format specifier “<code>H</code>” to print host name. Host name printing is costly operation.</p>

Table 57. The `__ip6hdr` built-in variable members (continued)

Member name	Type	Description
protocol	unsigned short	This value indicates the protocol that is used in the data portion of the IP datagram. Protocols can be TCP, UDP, and ICMPV6, and so on. It can match one of the following built-in constant values for protocol: IPPROTO_TCP, IPPROTO_UDP, IPPROTO_ROUTING, IPPROTO_ICMPV6, IPPROTO_NONE, IPPROTO_DSTOPTS, IPPROTO_LOCAL Refer the header file <code>/usr/include/netinet/in.h</code> for protocol values.
hop_limit	unsigned short	Hop limit (time to live).
total_len	unsigned short	Total length (payload length). The size of the payload including any extension headers.
next_hdr	unsigned short	Specifies the type of the next header. This field usually specifies the transport layer protocol that is used by a packet's payload. When extension headers are present in the packet, this field indicates which extension header follows. The values are shared with those used for the IPv4 protocol field.
flow_label	unsigned int	Flow label.
traffic_class	unsigned int	Traffic class.

`__tcphdr` built-in variable

The `__tcphdr` variable is a special built-in variable to get the tcp header information from filtered packet. This variable is available when you probe the packet information at interface layer with tcp protocol. It is available in probes of sub type bpf. Its member elements can be accessed by using the syntax `__tcphdr->member`.

The `__tcphdr` built-in variable has the following members:

Table 58. The `__tcphdr` built-in variable members

Member name	Type	Description
src_port	unsigned short	Source port of the packet.
dst_port	unsigned short	Destination port of the packet.

Table 58. The `__tcp_hdr` built-in variable members (continued)

Member name	Type	Description
flags	unsigned short	<p>These values are the control bits and are set to indicate the communication of control information. 1 bit for each flag.</p> <p>It can match one of the built-in constant flag values. The flag values must be bitwise and with the built-in constant flag value to validate the presence of the particular flag.</p> <ul style="list-style-type: none"> • TH_FIN (No more data from sender) • TH_SYN (request to establish the connection) • TH_RST (Reset the connection) • TH_PUSH (Push function. Asks to push the buffered data to the receiving application) • TH_ACK (Indicates that this packet contains acknowledgment) • TH_URG (Indicates that the urgent pointer field is significant) <p>Refer TCP documentation for detailed information about these flags and refer the header file <code>/usr/include/netinet/tcp.h</code> for flag values.</p>
seq_num	unsigned int	Sequence number.
ack_num	unsigned int	Acknowledgment number.
hdr_len	unsigned int	TCP header length information
cksum	unsigned short	Checksum.
window	unsigned short	Window size.
urg_ptr	unsigned short	Urgent pointer.

__udphdr built-in variable

The `__udphdr` is a special built-in variable that is used to get the udp header information from filtered packet. This built-in is available when user probes the packet information at interface layer with udp as protocol. It is available in probes of sub type `bpf`. Its member elements can be accessed by using the syntax `__udphdr->member`.

`__udphdr` built-in variable has the following members:

Table 59. The `__udphdr` built-in variable members

Member name	Type	Description
<code>src_port</code>	unsigned short	Source port of the packet.
<code>dst_port</code>	unsigned short	Destination port of the packet.
<code>length</code>	unsigned short	UDP header and data length information.
<code>cksum</code>	unsigned short	Checksum.

`__icmp` built-in variable

The `__icmp` is a special built-in variable that is used to get the icmp header information from filtered packet. This built-in is available when user probes the packet information at interface layer with icmp protocol. It is available in probes of sub type bpf. Its member elements can be accessed by using the syntax `__icmp->member`.

This built-in variable has the following members:

Table 60. The `__icmp` built-in variable members

Member name	Type	Description
<code>type</code>	unsigned short	<p>Type of ICMP message.</p> <p>For example: 0 - echo reply, 8 - echo request, 3 - destination unreachable. Look in for all the types. For more information, refer to the standard network documentation.</p> <p>It can match one of the following built-in constant values for of ICMP message types:</p> <pre> ICMP_ECHOREPLY, ICMP_UNREACH ICMP_SOURCEQUENCH, ICMP_REDIRECT, ICMP_ECHO, ICMP_TIMXCEED, ICMP_PARAMPROB, ICMP_TSTAMP, ICMP_TSTAMPREPLY, ICMP_IREQ, ICMP_IREQREPLY, ICMP_MASKREQ, ICMP_MASKREPLY </pre> <p>Refer the header file <code>/usr/include/netinet/ip_icmp.h</code> for protocol values.</p> <p>Note: All possible message type values are not defined, and hence there can be other options present in the value.</p>

Table 60. The `__icmp` built-in variable members (continued)

Member name	Type	Description
code	unsigned short	<p>Subtype of ICMP message.</p> <p>For each type of message, several different codes and subtypes are defined. For example, no route to destination, communication with destination administratively prohibited, not a neighbor, address unreachable, port unreachable. For more information, refer to the standard network documentation.</p> <p>It can match one of the following built-in constant values for ICMP sub types:</p> <p>ICMP_UNREACH_NET ICMP_UNREACH_HOST ICMP_UNREACH_PROTOCOL ICMP_UNREACH_PORT ICMP_UNREACH_NEEDFRAG ICMP_UNREACH_SRCFAIL ICMP_UNREACH_NET_ADMIN_PROHIBITED ICMP_UNREACH_HOST_ADMIN_PROHIBITED</p> <p>Subtype values for type 4</p> <p>The subtype values for type 4 are as follows:</p> <pre>ICMP_REDIRECT_NET ICMP_REDIRECT_HOST ICMP_REDIRECT_TOSNET ICMP_REDIRECT_TOSHOST</pre> <p>Subtype values for type 6</p> <p>The subtype values for type 6 are as follows:</p> <pre>ICMP_TIMXCEED_INTRANS ICMP_TIMXCEED_REASS</pre> <p>Subtype values for type 7</p> <p>The subtype values for type 7 are as follows:</p> <p>ICMP_PARAMPROB_PTR ICMP_PARAMPROB_MISSING</p> <p>Refer the header file <code>/usr/include/netinet/ip_icmp.h</code> for message subtype values.</p> <p>Note: Not all possible message sub types values are defined, and hence there might be other options present in the message sub type value.</p>
cksum	unsigned short	Checksum.

__icmp6 built-in variable

`__icmp6` is a special built-in variable that is used to get the `icmpv6` header information from filtered packet. This is available when user probes the packet information at interface layer with `icmp6` protocol. It is available in probes of sub type `bpf`. Member elements of this built-in variable can be accessed using syntax “`__icmp6->member`”.

`__icmp6` has the following members:

Table 61. The `__icmp6` built-in variable members

Member name	Type	Description
type	unsigned short	<p>Type of ICMPV6 message.</p> <p>This specifies the type of message, which determines the format of the remaining data.</p> <p>It can match one of the following built-in constant values for ICMPV6 types.</p> <pre>ICMP6_DST_UNREACH ICMP6_PACKET_TOO_BIG ICMP6_TIME_EXCEEDED ICMP6_PARAM_PROB ICMP6_INFOMSG_MASK ICMP6_ECHO_REQUEST ICMP6_ECHO_REPLY</pre> <p>Refer the header file <code>/usr/include/netinet/icmp6.h</code> for protocol values.</p> <p>Note: Not all possible message type values are defined, and hence there might be other options present in the value.</p>
code	unsigned short	<p>Subtype of ICMPV6 message.</p> <p>This value depends on the message type. It provides an extra level of message granularity.</p> <p>It can match one of the following built-in constant values for ICMPV6 sub types.</p> <pre>ICMP6_DST_UNREACH_NOROUTE ICMP6_DST_UNREACH_ADMIN ICMP6_DST_UNREACH_ADDR ICMP6_DST_UNREACH_BEYONDSCOPE ICMP6_DST_UNREACH_NOPORT</pre> <p>Refer the header file <code>/usr/include/netinet/icmp6.h</code> for message subtype values.</p> <p>Note: Not all possible message sub type values are defined, and hence there might be other options present in the value.</p>
cksum	unsigned short	Checksum.

__igmp built-in variable

`__igmp` is a special built-in variable that is used to get the `igmp` header information from filtered packet. This is available when user probes the packet information at interface layer with `igmp` protocol. This is available in probes of sub type `bpf`. Its member elements can be accessed using syntax “`__igmp->member`”.

__igmp built-in has the following members:

Table 62. The __igmp built-in variable members		
Member name	Type	Description
type	unsigned short	<p>Type of IGMP message.</p> <p>For example: Membership Query (0x11), Membership Report (IGMPv1: 0x12, IGMPv2: 0x16, IGMPv3: 0x22), Leave Group (0x17)</p> <p>For more information, refer to the standard or Network documentation.</p> <p>It can match one of the following built-in constant values for IGMP Message types.</p> <pre>IGMP_HOST_MEMBERSHIP_QUERY IGMP_HOST_MEMBERSHIP_REPORT IGMP_DVMRP IGMP_HOST_NEW_MEMBERSHIP_REPORT IGMP_HOST_LEAVE_MESSAGE IGMP_HOST_V3_MEMBERSHIP_REPORT IGMP_MTRACE IGMP_MTRACE_RESP IGMP_MAX_HOST_REPORT_DELAY</pre> <p>Refer the header file /usr/include/netinet/igmp.h for protocol values.</p> <p>Note: Not all possible message type values are defined, and hence there could be other options present in the value.</p>
code	unsigned short	<p>Subtype of IGMP type.</p> <p>It can match one of the following built-in constant values for IGMP Message subtypes.</p> <p>Subtype values for type no 3.</p> <pre>DVMPP_PROBE 1 DVMRP_REPORT 2 DVMRP_ASK_NEIGHBORS 3 DVMRP_ASK_NEIGHBORS2 4 DVMRP_NEIGHBORS 5 DVMRP_NEIGHBORS2 6 DVMRP_PRUNE 7 DVMRP_GRAFT 8 DVMRP_GRAFT_ACK 9 DVMRP_INFO_REQUEST 10 DVMRP_INFO_REPLY 11</pre> <p>Note: Not all possible message sub type values are defined, and hence there could be other options present in the value.</p>

Table 62. The `__igmp` built-in variable members (continued)

Member name	Type	Description
<code>cksum</code>	unsigned short	IGMP Checksum value.
<code>group_addr</code>	<code>ip_addr_t</code>	<p>Group address that is reported or queried.</p> <p>This address is the multicast address that is queried when you are sending a Group-Specific or Group-and-Source-Specific Query. The field has a value of zero when you are sending a General Query.</p> <p>The data type <code>ip_addr_t</code> is used to store the group IP address. Use format specifier "I" to print the IP address.</p>

`__arphdr` built-in variable

The `__arphdr` variable is a special built-in variable that is used to get the `arphdr` header information from filtered packet. This variable is available when user probes the packet information at interface layer with `arp` or `rarp` protocol. It is available in probes of sub type `bpf`. The `__arphdr` member elements can be accessed by using the syntax `__arphdr->member`.

The `__arphdr` built-in variable has following members:

Table 63. The `__arphdr` built-in variable members

Member name	Type	Description
<code>hw_addr_type</code>	unsigned short	<p>Format of the hardware address type. This field identifies the specific data-link protocol that is being used.</p> <p>It can match one of the following built-in constant values for data link protocol:</p> <pre>ARPHRD_ETHER, ARPHRD_802_5, ARPHRD_802_3, and ARPHRD_FDDI</pre> <p>Refer the header file <code>/usr/include/net/if_arp.h</code> for protocol values.</p>

Table 63. The `__arphdr` built-in variable members (continued)

Member name	Type	Description
<code>protocol_type</code>	unsigned short	<p>Format of the protocol address type. This field identifies the specific network protocol that is being used.</p> <p>It can match one of the following built-in constant values for network protocol:</p> <pre>SNAP_TYPE_IP, SNAP_TYPE_AP, SNAP_TYPE_ARP, VLAN_TAG_TYPE</pre> <p>Refer the header file <code>/usr/include/net/nd_lan.h</code> for the protocol values.</p>
<code>hdr_len</code>	unsigned short	Mac or hardware address length.
<code>proto_len</code>	unsigned short	Protocol or IP address length.
<code>operation</code>	unsigned short	<p>Specifies the operation that the sender is performing: 1 for request, 2 for reply.</p> <p>It can match one of the following built-in constant values for network protocol:</p> <pre>ARPOP_REQUEST, ARPOP_REPLY</pre> <p>Refer the header file <code>/usr/include/net/if_arp.h</code> for protocol values.</p>
<code>src_mac_addr</code>	<code>mac_addr_t</code>	<p>Sender or source MAC address.</p> <p>Sender hardware or mac address is stored in <code>mac_addr_t</code> data type. The format specifier “%M” is used to print sender MAC or hardware address.</p>
<code>dst_mac_addr</code>	<code>mac_addr_t</code>	<p>Target or Destination MAC address.</p> <p>Target hardware or MAC address is stored in <code>mac_addr_t</code> data type. The format specifier “%M” is used to print target MAC or hardware address.</p>

Table 63. The `__arphdr` built-in variable members (continued)

Member name	Type	Description
<code>src_ip</code>	<code>ip_addr_t</code>	Source or sender IP address. Sender IP address is stored in <code>ip_addr_t</code> data type. The format specifier “%I” is used to print sender IP address.
<code>dst_ip</code>	<code>ip_addr_t</code>	Target or Destination IP address. Target IP address is stored in <code>ip_addr_t</code> data type. The format specifier “%I” is used to print target IP address.

Example

Vue script to probe packet header information for packets received or sent over port 23. Provides the source and destination node information and also tcp header length information

```
@@net:bpf:en0:tcp:"port 23"
{
    printf("src_addr:%I and dst_addr:%I\n", __ip4hdr->src_addr, __ip4hdr->dst_addr);
    printf("src port:%d\n", __tcphdr->src_port);
    printf("dst port:%d\n", __tcphdr->dst_port);
    printf("tcp hdr_len:%d\n", __tcphdr->hdr_len);
}
```

```
Output:
# probevue bpf_tcp.e
src_addr:10.10.10.12 and dst_addr:10.10.18.231
src port:48401
dst port:23
tcp hdr_len:20
.....
.....
```

`__proto_info` built-in variable

The `__proto_info` variable is a special built-in variable that is used to get the protocol (source and destination IP addresses and ports) information for TCP or UDP events. The `__proto_info` variable is available in probes of sub type `tcp` or `udp`. Its member elements can be accessed by using the syntax `__proto_info->member`.

The `__proto_info` built-in variable has the following members:

Table 64. The `__proto_info` built-in variable members

Member name	Type	Description
<code>local_port</code>	unsigned short	Local port
<code>remote_port</code>	unsigned short	Remote port
<code>local_addr</code>	<code>ip_addr_t</code>	Local address
<code>remote_addr</code>	<code>ip_addr_t</code>	Remote address

Additional information for TCP-specific events

The TCP state change events are described in the following table:

Table 65. TCP state change events

Name	Type	Description
__prev_state	short	Previous state information for connection.
__cur_state	short	Present state information for connection.
		<p>It can match one of the following built-in constant values for TCP states:</p> <ul style="list-style-type: none"> • TCPS_ESTABLISHED (connection established) • TPCS_CLOSED (Connection closed) • TPCS_LISTEN (Listening for connection) • TPCS_SYN_SENT (Sent SYN to remote end) • TCPS_SYN_RECEIVED (Received SYN from remote end) • TCPS_CLOSE_WAIT (Received Fin, waiting for close) • TCPS_FIN_WAIT_1 (are closed, sent fin) • TCPS_CLOSING (closed exchanged FIN, await FIN ACK) • TCPS_LAST_ACK (Had Fin and close , Await FIN ACK) • TCPS_FIN_WAIT_2 (are closed , Fin is Acked) • TCPS_TIME_WAIT (in 2*msl quiet wait after close) <p>The values are defined in exported header file /usr/include/netinet/tcp_fsm.h.</p>

Example:

The following Vue script provides state change information for a particular connection:

```

@@net:tcp:state_change
when(__proto_info->local_addr == "10.10.10.1" and __proto_info->remote_addr == 10.10.10.2"
    and __proto_info->local_port == "8000" and __proto_info->remote_port == "9000")
{
    printf("Previous state:%d and current_state:%d\n", __prev_state, __cur_state);
}
    
```

TCP retransmit event

Table 66. TCP retransmit event

Name	Type	Description
__nth_retransmit	unsigned short	Nth retransmission

Examples

1. Following example Identifies the listener which has discarded connections due to listener's queue is full.

```
@@net:tcp:listen_q_full
{
    printf("Listener IP address:%I and Port number is:%d\n",__proto_info->local_addr, __proto_info->local_port);
}
```

2. Following example Identifies connection which drop packets due to socket buffer overflows

```
@@net:udp:sock_recv_buf_overflow
{
    printf("Connection information which drops packet due to socket buffer overflows:\n");
    printf("Local IP address:%I and Remote IP address:%I\n",__proto_info->local_addr,__proto_info->remote_addr);
    printf("local port :%d and remote port:%d\n",__proto_info->local_port, __proto_info->remote_port);
}
```

3. Identify retransmissions (second & further retransmission for a packet) for TCP connections for particular connection.

```
@@net:tcp:retransmit
when (__proto_info->local_addr == "10.10.10.1" &&
    __proto_info->remote_addr == "10.10.10.2" &&
    __proto_info->local_port == "4000" &&
    __proto_info->remote_port == "5000")
{
    printf(" %d th re-transmission for this connection\n", _nth_retransmit);
}
```

4. Identify the connection information whenever sender buffer full event occurs .

```
@@net:tcp:send_buf_full
{
    printf("Connection information whenever send buffer full event occurs:\n");
    printf("Local IP address:%I and Remote IP address:%I\n",__proto_info->local_addr,__proto_info->remote_addr);
    printf("local port :%d and remote port:%d\n",__proto_info->local_port, __proto_info->remote_port);
}
```

Sysproc probe manager

The sysproc probe manager provides an infrastructure to users and administrators to dynamically trace process or thread related data without knowing internals of sysproc subsystem.

The aspects of sysproc subsystem for a user or administrator is divided into the following main categories:

- Process (or thread) creation or termination
- Signal generation and delivery
- Scheduler and dispatcher events
- DR and CPU binding events

Process (or thread) creation or termination

Information related to how a process or thread is created and destroyed is required to a system administrator to administer the resources of the system. The sysproc probe manager addresses the following important use-cases:

- Did a process exit naturally or because of an error?
- When a process or thread got created or terminated or exceed?
- How long did a process run?

- Track events when a thread receives or returns from an exception.

Signal generation and delivery

Signals decide the current state of a processor thread in a system. To understand misbehaving process or threads an administrator uses the state of signals and the current state of processes due to these signals. The important use-cases under signal generation and delivery category (but not limited) addressed by this probe manager follow:

- Signal source and signal information for a specific target.
- Signal delivery of asynchronous signals.
- Trace signal clears.
- Trace events when a signal handler other than default is installed.
- Signal target and signal information for a specific source.
- Trace signal handler entry or exit.

Scheduler and dispatcher events

Scheduler and dispatcher dictate how a process or thread runs in the system. Administrator analyzes system performance by using dynamic trace scheduler or dispatcher subsystem.

The dynamic trace scheduler or dispatcher subsystem helps discover the reasons for retention of threads.

Following are the important use-cases under scheduler and dispatcher events category (but not limited) addressed by sysproc probe manager.

- Trace thread or threads that are enqueued or dequeued from the run queue.
- Trace events when any thread in the system is preempted.
- Trace when a thread is being put to sleep over an event.
- Trace when a sleeping thread is being woken up.
- Track dispatches latency of a thread.
- Track virtual processor folding events.
- Trace change in any kernel thread priority.

Dynamic Reconfiguration (DR), and CPU binding events

This class of probes offer dynamic tracing capabilities to a user who tracks resources bound to a process.

Some of important use-cases (but not limited) under this category that is addressed by DR and CPU binding events probe manager follow:

- Track when a thread binding changes from one CPU to another.
- Track when the resources are attached or detached to a process.
- Track CPU binding events.
- Track start or end of a DR event.

Probe specification

The following format must be used in a Vue script to probe sysproc events:

```
@@sysproc:<sysproc_event>:<pid/tid/*>
```

First tuple @@sysproc indicates that this probe is specific to sysproc events.

Second tuple specifies the event to be probed.

Third tuple acts as a filter to isolate events that are specified through second tuple based on process or kernel thread id.

Note: Use of process or kernel thread id as filter in sysproc probes does not guarantee the event to occur in process or thread context. Sysproc probe manager uses process or thread id only as a filter. These events might be useful from a process or thread perspective despite the execution context of the probe event.

Signal send event, where either the process that is sending the signal or the one receiving it, can be useful. The following information specifies the appropriate filters for such probe events.

Probe points (events of interest)

A brief description of all events that can be probed through the sysproc probe manager is mentioned in the following table:

Probe (sysproc_event)	Description
forkfail	Track failures in fork interface.
execfail	Track failures in exec interface.
execpass	Track exec success.
exit	Track exit of a process.
threadcreate	Track creation of a kernel thread.
threadterminate	Track termination of a kernel thread.
threadexcept	Track process exceptions.
sendsig	Track signal sent to a process by external sources.
sigqueue	Tracks signals queued to a process
sigdispose	Tracks signal disposals.
sigaction	Track signal handler installations and reinstallations
sighandlestart	Track when a signal handler is about to be called.
sighandlefinish	Track when a signal handler completion
changepriority	Track when priority of a process changes
onreadyq	Track when a kernel thread gets on a ready queue.
offreadyq	Track when a kernel thread is moved out of ready queue.
dispatch	Track when the system dispatcher is called to schedule a thread
oncpu	Track when a kernel thread acquires CPU.
offcpu	Track when a kernel thread relinquishes CPU.
blockthread	Track when a thread is blocked from getting CPU.
foldcpu	Track folding of a CPU core.
bindprocessor	Track event when a process/thread is bound to a CPU
changecpu	Track events when a kernel thread changes CPU temporarily

Table 67. Sysproc probe events (continued)

Probe (sysproc_event)	Description
resourceattach	Track events when a resource is attached to another
resourcedetach	Track events when a resource is detached from another
drphasestart	Track when a drphase is getting initiated
drphasefinish	Track when a drphase completes

Method to access data at a probe-point

ProbeVue allows data access through built-in variables.

Built-ins values are of three types categorized based on accessibility:

1. Accessible at any probe point, irrespective of the probe manager. For example: `__curthread`.
2. Accessible throughout probes of a specific probe manager.
3. Accessible only at defined probes (events of interest)

The sysproc probe manager allows access of data through built-ins of type (1) and (3). The following table is used to indicate accessibility of built-ins of type (1). Special built-ins that are provided for sysproc probe manager are of type long long.

Following are the list of built values of type (1).

- `__trcid`
- `__errno__kernelmode`
- `__arg1` to `__arg7`
- `__curthread`
- `__curproc`
- `__mst`
- `__tid`
- `__pid`
- `__ppid`
- `__pgid`
- `__uid`
- `__euid`
- `__ublock`
- `__execname`
- `__pname`

The built-in variables are also classified as context specific and context independent. Context-specific built-ins provide data based on the execution context of the probe.

AIX kernel operates in thread or interrupt context. Context-specific probes produce correct result when probe is started at thread or process context.

Results that are obtained from context-specific built-ins in interrupt execution context might be unexpected. Context-independent built-ins do not depend on the execution context and can be accessed safely irrespective of probe execution environment.

Table 68. Context specific and independent built-in variables

Context specific built-in variables	Context independent built-in variables
__curthread	__trcid
__curproc	__errno
__tid	__kernelmode
__pid	__arg1 to __arg7
__ppid	__mst
__pgid	
__uid	
__euid	
__ublock	
__pname	
__execname	

Probe points

Probe points are the specific events for which a probe is fired. Following are the list of probe points.

forkfail

The forkfail probe starts when fork fails. This probe determines the reasons of fork failure.

Syntax: @@sysproc:forkfail:<pid/tid/*>

Special built-in supported

```

__forkfailinfo
{
fail_reason;
}

```

The fail_reason variable has one of the following values:

Reason	Description
FAILED_RLIMIT	Failed due to rlimit limitations
FAILED_ALLOCATIONS	Failed due to internal resource allocations
FAILED_LOADER	Failed at a loader stage
FAILED_PROCDUP	Failed at procdup

Other supported built-ins

__errno __kernelmode, __arg1 to __arg7, __curthread, __curproc, __mst, __tid, __pid, __ppid, __pgid, __uid, __euid, __ublock, __execname, __pname.

Execution environment

Runs in process environment.

Example

Following example shows how to monitor all fork failures because of rlimit in the system.

```

@@BEGIN
{
    x = 0;
}

@@sysproc:forkfail:*
when (__forkfailinfo->fail_reason == FAILED_RLIMIT)
{
    printf ("process %s with pid %llu failed to fork a child\n",__pname,__pid);
    x++;
}

@@END
{
    printf ("Found %d failures during this vue session\n",x);
}

```

execfail

The execfail probe starts when a exec function call fails. Use the execfail probe to determine the reasons for the failure.

Syntax: @@sysproc:execfail:<pid/tid/*>

Table 70. execfail probe: Failure reasons

Reason	Description
FAILED_PRIVILEGES	New process failed to acquire or inherit privileges
FAILED_COPYINSTR	New process failed to copy instruction
FAILED_V_USERACC	New process failed to discard v_useracc regions
FAILED_CLEARDATA	Failed during clearing data for new process
FAILED_PROCSEG	Failed to establish process private segment
FAILED_CH64	Failed to convert to a 64-bit process
FAILED_MEMATT	Failed to attach to a memory resource set
FAILED_SRAD	Failed to attach to a srad
FAILED_MSGBUF	Error message buffer length is zero
FAILED_ERRBUF	Failed to allocate error message buffer
FAILED_ENVAR	Failed to allocate environment variables
FAILED_CPYSTR	Copy string error
FAILED_ERRBUFCPY	Failed to copy the error messages from errmsg_buf
FAILED_TOOLNGENV	Env too long for allocated memory
FAILED_USRSTK	Failed to setup user stack
FAILED_CPYARG	Failed to copy arglist to stack
FAILED_INITPTRACE	Failed to init ptrace

Note: 64 is added to error value if loader error is encountered.

Other supported built-ins

__errno __kernelmode, __arg1 to __arg7, __curthread, __curproc, __mst, __tid, __pid, __ppid, __pgid, __uid, __euid, __ublock, __execname, __pname.

Execution environment

Runs in process environment.

exit

This probe starts when a process exits. Exit is also a system call manager and is traced through system call probe manager. Probing exit system call through sysproc probe manager explains nature and reasons of exit. It also explains reasons for a user thread termination in kernel space and not returned to user space.

Syntax: @@sysproc:forkfail:<pid/tid/*>

A program can exit because of the following reasons:

- On reaching a terminal condition when a user space program cannot proceed further.
- On receiving a terminal signal.

Special built-in supported

```
__exitinfo{
    signo;
    returnval;
    iscore;
}
```

Where, signo value signifies the signal number that caused process termination, returnval is the value that is returned by exit. Nonzero signo is valid only if the program is stopped by a signal.

The iscore variable is set when a core is generated as a result of process exit.

Other supported built-ins

__errno __kernelmode, __arg1 to __arg7, __curthread, __curproc, __mst, __tid, __pid, __ppid, __pgid, __uid, __euid, __ublock, __execname, __pname.

Execution Environment

Runs in process environment.

Example

Following example explains how to probe exit event

```
echo '@@sysproc:exit:* { printf (" %s    %llu    %llu\n", __pname, __pid,__exitinfo->returnval);} ' | probevue
```

Which will produce an output similar to the following.

```
ksh      5833042    0
telnetd  7405958    1
dumpctl  7405960    0
setmaps  7275006    0
termdef  7274752    0
hostname 7274754    0
id       8257976    0
id       8257978    0
uname   8257980    0
expr    8257982    1
```

threadcreate

threadcreate probe starts when a thread is created successfully.

Syntax: @@sysproc:threadcreate:<pid/tid/*>

Note: The specified pid or tid must be the process or thread ID of the process or thread that created the thread.

Special built-in supported

```
__threadcreateinfo
{
    tid;
    pri;
    policy;
}
```

where `tid` indicates the thread id of new thread that is created, and priority is the priority of the thread. Policy denotes the thread scheduling policy of the thread.

Policy	Description
SCHED_OTHER	default AIX scheduling policy
SCHED_FIFO	first in-first out scheduling policy
SCHED_RR	round robin scheduling policy
SCHED_LOCAL	local thread scope scheduling policy
SCHED_GLOBAL	global thread scope scheduling policy
SCHED_FIFO2	FIFO with RQHEAD after short sleep
SCHED_FIFO3	FIFO with RQHEAD all the time
SCHED_FIFO4	FIFO with weak preempt

Other supported built-ins

`__errno__kernelmode`, `__arg1` to `__arg7`, `__curthread`, `__curproc`, `__mst`, `__tid`, `__pid`, `__ppid`, `__pgid`, `__uid`, `__euid`, `__ublock`, `__execname`, `__pname`.

Execution environment

Runs in process environment (user or kproc).

Example

To continuously print all processes in the system creating a thread printing process name, creating process id , id of the newly created thread and creation time-stamp.

```
echo '@@sysproc:threadcreate:*
{ printf ("%s %llu %llu %A\n",__pname,__pid,__threadcreateinfo->tid,timestamp());}' |
probevue
```

An output similar to the following example is displayed.

```
nfssync_kproc 5439964 23921151 Feb/22/15 09:22:38
nfssync_kproc 5439964 24052201 Feb/22/15 09:22:38
nfssync_kproc 5439964 23920897 Feb/22/15 09:22:38
nfssync_kproc 5439964 22479285 Feb/22/15 09:22:55
nfssync_kproc 5439964 23920899 Feb/22/15 09:22:55
nfssync_kproc 5439964 22479287 Feb/22/15 09:22:55
```

threadterminate

The probe starts for a thread which is terminated.

Syntax: `@@sysproc:threadterminate:<pid/tid/*>`

Note: Specified process ID or thread ID must be corresponding to the process or thread currently getting stopped.

Special built-ins supported

None.

Other supported built-ins

`__errno__`, `kernelmode`, `__arg1` to `__arg7`, `__curthread`, `__curproc`, `__mst`, `__tid`, `__pid`, `__ppid`, `__pgid`, `__uid`, `__euid`.

Execution environment

Runs in process environment (user or kproc).

Example

To continuously print all processes in the system terminating a thread printing process name, creating process id, id of the newly created thread and creation time-stamp.

```
# echo '@@sysproc:threadterminate:* { printf ("%s %llu %llu %A\n", __pname, __pid, __tid, timestamp()); }' | probevue
```

An output similar to the following sample is displayed:

```
nfssync_kproc 5439964 23855555 Feb/22/15 09:59:30
nfssync_kproc 5439964 21758249 Feb/22/15 09:59:30
nfssync_kproc 5439964 23855557 Feb/22/15 09:59:30
```

threadexcept

This probe starts when a program exception occurs. A program exception is generated when system detects a condition in which a program cannot continue normally. Some exceptions are fatal (illegal instruction) while some can be recovered (address space change).

Syntax: `@@sysproc:threadexcept:<pid/tid/*>`

Special built-ins supported

```
__threadexceptinfo
{
    pid;
    tid;
    exception;
    excpt_address
}
```

where `pid` denotes process ID of the process that received exception, `tid` is the thread ID of the kernel thread that received exception, `excpt_address` is address that caused this exception while exception can assume one of the values as denoted in the table.

Exception	Description
EXCEPT_FLOAT	Floating point exception
EXCEPT_INV_OP	Invalid op-code
EXCEPT_PRIV_OP	Privileged op in user mode
EXCEPT_TRAP	Trap instruction
EXCEPT_ALIGN	Code or data alignment
EXCEPT_INV_ADDR	Invalid address
EXCEPT_PROT	Protection
EXCEPT_IO	Synchronous I/O
EXCEPT_IO_IOCC	I/O exception from IOCC
EXCEPT_IO_SGA	I/O exception from SGA
EXCEPT_IO_SLA	I/O exception from SLA

Table 72. Exception values for the threadexcept probe (continued)	
Exception	Description
EXCEPT_IO_SCU	I/O exception from SCU
EXCEPT_EOF	Reference beyond end-of-file (mmap)
EXCEPT_FLOAT_IMPRECISE	Imprecise floating point exception
EXCEPT_ESTALE_I	Stale text segment exception
EXCEPT_ESTALE_D	Stale data segment exception
EXCEPT_PT_WATCHP	Hit ptrace watchpoint

Other supported built-ins

__errno__ kernelmode, __arg1 to __arg7, __curthread, __curproc, __mst, __tid, __pid, __ppid, __pgid, __uid, __euid.

Execution environment

Runs in process or interrupt environment.

Note: Since this probe can start in an interrupt context, built-ins variables like __pid, __tid that depend upon the execution context might not indicate the process or thread id. Special built-in members for this probe guarantee correct process or thread id intended for the process or thread.

Example

Following example shows trace program exceptions generated by a probe event being traced by a debugger.

```
# cat threadexcept.e
@@sysproc:threadexcept:*
{
  printf ("PID = %llu TID= %llu EXCEPTION=%llu ADDRESS = %llu\n ", __threadexceptinfo->pid, __threadexceptinfo->tid, __threadexceptinfo->exception, __threadexceptinfo->excpct_address);
}

Run a debugging session on a program compiled with debugging support

# dbx a.out
Type 'help' for help.
Core file "core" is older than current program (ignored)
reading symbolic information ...
(dbx) stop in main
[1] stop in main
(dbx) r
[1] stopped in main at line 5
5          int a=5;
```

An output similar to the following sample is displayed:

```
PID = 6816134 TID= 24052015 EXCEPTION=131 ADDRESS = 268436372
```

sendsig

This probe is started when a signal is sent to a process through external sources (other process , process from user space, from kernel streams or Interrupt context)

Syntax:@@sysproc:sendsig:<pid/*>

```
__dispatchinfo{
  cpuid;          <- cpu id

  oldpid;         <- pid of the thread currently running
  oldtid;         <- thread id of the thread currently running
  oldpriority;    <- priority of the thread currently running
  newpid;         <- pid of the new process process selected for running
  newtid;         <- thread id of the thread selected for running
  newpriority;    <-priority of the thread selected for running
}
}
```

where pid is the process identifier of the target process receiving the signal. This probe does not allow specifying a thread identifier to filter results specific to a thread.

Special built-ins

```
_sigsendinfo{
  tpid;           ← target pid
  spid;           ← source pid
  signo;          ← signal sent
}
```

where tpid is the target source process identifier, spid identifies source of the signal. The spid is non-zero when signal is sent from user space or process context. Source process identifier is 0 if signal is sent from an exception or interrupt context. Signal number information is contained in signo.

Other supported built-ins

__errno__ kernelmode, __arg1 to __arg7, __curthread, __curproc, __mst, __tid, __pid, __ppid, __pgid, __uid, __euid.

Execution environment

Runs in process or interrupt environment.

Note: Since this probe can start in an interrupt context, it is possible that built-ins like __pid, __tid, which depend upon thread execution context might not indicate the process or thread id of interest. Special built-in members for this probe guarantee correct process or thread id intended for the process or thread.

When this probe starts in process context, built-in members that depend on execution context point to source process. built-in members like __pid, __tid, and __curthread provide information regarding the source process.

.

Example

```
To continuously print signal source signal target and signal number of all signals.
echo '@@sysproc:sendsig:* {printf ("Source=%llu Target=%llu sig=%llu\n",__sigsendinfo->spid,__sigsendinfo->tpid,__sigsendinfo->signo);} |
probevue
```

An output similar to the following sample is displayed:

```
Source=0 Target=6619618 sig=14
Source=0 Target=8257944 sig=20
Source=0 Target=8257944 sig=20
```

sigqueue

This probe starts when a queued signal is being sent to the process.

Syntax: @@sysproc:sigqueue:<pid/*>

Special built-ins

```
_sigsendinfo{
  tpid;           ← target pid
  spid;           ← source pid preprocess.cp
  signo;          ← signal sent
}
```

Since posix signals are queued to a process, specifying thread identifier is not allowed in this probe.

Other supported built-ins

__errno__ kernelmode, __arg1 to __arg7, __curthread, __curproc, __mst, __tid, __pid, __ppid, __pgid, __uid, __euid.

This probe starts in the context of the sending process. Hence, context-based built-ins refer to the sending process in this probe event.

Execution environment

This probe runs in process context.

Example

```
echo '@@sysproc:sigqueue:*{printf ("%llu %llu %llu\n",__sigsendinfo->spid,__sigsendinfo->tpid,__sigsendinfo->signo);} | probevue
```

An output similar to the following sample is displayed:

```
8258004 6095294 31
sigdispose
Syntax : @@sysproc:sigdispose:<pid/tid/*>
```

Probe starts when a signal is disposed to a target process. Specify process ID of the process which received this signal in the sysprobe specification to filter this probe.

Special built-ins

```
__sigdisposeinfo{
    tpid;          ← target pid
    ttid;         ← target tid
    signo;        ← signal whose action is being taken.
    fatal;        ← will be set if the process is going to be killed as part of signal
    action
}
```

Other supported built-ins

__errno __kernelmode, __arg1 to __arg7, __curthread, __curproc, __mst, __tid, __pid, __ppid, __pgid, __uid, __euid.

Execution environment

This probe can start from process or interrupt context. If started from interrupt context, this probe might not provide required value for context-based built-ins.

Example

Continuously print process identifier, thread identifier, signal number and indicate if this signal disposal will result in termination of the process for all processes in the system.

```
cat sigdispose.e
@@sysproc:sigdispose:*
{
    printf ("%llu %llu %llu %llu\n",__sigdisposeinfo->tpid,__sigdisposeinfo->ttid, __sigdisposeinfo->signo,__sigdisposeinfo->fatal);
}
```

An output similar to the following sample is displayed:

```
5964064 20840935 14 0
1 65539 14 0
4719084 19530213 14 0
```

sigaction

Syntax:@@sysproc:sigaction:<pid/tid/*>

This probe starts when a signal handler is installed or replaced.

Special built-ins

```
__sigactioninfo{
    old_sighandle;          ← old signal handler function address
    new_sighandle;         ←new signal handler function address
    signo;                 ← Signal number
    rpid;                  ← requester's pid
}
```

old_sighandle will be 0 if a signal handler is installed for the first time.

Other supported built-ins

__errno__ kernelmode, __arg1 to __arg7, __curthread, __curproc, __mst, __tid, __pid, __ppid, __pgid, __uid, __euid.

Execution environment

This probe starts in process environment.

Note: AIX kernel ensures that only one signal is delivered to a process or thread at a time. Another signal to that process or thread is only sent when signal delivery is finished.

Example

To track the beginning and finish of all signals in a system:

```
@@sysproc:sighandlestart:*
{
    signal[__tid] = __sighandlestartinfo->signo;
    printf ("Signal handler at address 0x%x invoked for thread id %llu to handle signal %llu\n", __sighandlestartinfo->
sighandle, __curthread->tid, __sighandlestartinfo->signo);
}

@@sysproc:sighandlefinish:*
{
    printf ("Signal handler completed for thread id %llu for signal %llu\n", __curthread->tid, signal[__tid]);
    delete (signal, __tid);
}
```

An output similar to the following sample can be displayed:

```
Signal handler at address 0x20001d58 invoked for thread id 19923365 to handle signal 20
Signal handler completed for thread id 19923365 for signal 20
Signal handler at address 0x10003400 invoked for thread id 20840935 to handle signal 14
Signal handler completed for thread id 20840935 for signal 14
Signal handler at address 0x10002930 invoked for thread id 19530213 to handle signal 14
Signal handler completed for thread id 19530213 for signal 14
Signal handler at address 0x300275d8 invoked for thread id 22348227 to handle signal 14
Signal handler completed for thread id 22348227 for signal 14
Signal handler at address 0x20001a3c invoked for thread id 65539 to handle signal 14
Signal handler completed for thread id 65539 for signal 14
```

sighandlefinish

This probe starts at signal handler completion.

Syntax: @@sysproc:sighandlestart:<pid/tid/*>

Special built-ins supported: None.

Other supported built-ins

__errno__ kernelmode, __arg1 to __arg7, __curthread, __curproc, __mst, __tid, __pid, __ppid, __pgid, __uid, __euid.

Execution environment

Runs in process environment. Protected, context switch is not allowed on executing CPU.

changepriority

This probe starts when the priority of a process is being changed. This event is not a scheduler or dispatcher-enforced.

Syntax: @@sysproc:changepriority:<pid/tid/*>

Note: The priority change might also be unsuccessful; success of priority change is not guaranteed.

Special built-ins supported

```
__chpriorityinfo{
    pid;
    old_priority;    <- current priority
```

```

    new_priority; <- new scheduling priority of the thread.
}

```

Execution Environment

This probe runs in process environment.

Other supported built-ins

__errno__kernelmode, __arg1 to __arg7, __curthread, __curproc, __mst, __tid, __pid, __ppid, __pgid, __uid, __euid, __ublock, __execname, __pname.

Example

To track all processes whose priority is being changed:

```

echo '@@sysproc:changepriority:* { printf ("%s priority changing from %llu to %llu\n",__pname,__chpriorityinfo->old_priority,__chpriorityinfo->new_priority);}' | probevue

```

An output similar to the following sample is displayed:

```

xmgc priority changing from 60 to 17
xmgc priority changing from 17 to 60
xmgc priority changing from 60 to 17
xmgc priority changing from 17 to 60
xmgc priority changing from 60 to 17

```

offreadyq

This probe starts when a thread is removed from a system run queue.

Syntax: @@sysproc:offreadyq:<pid/tid/*>

Special built-ins supported

```

__readyprocinfo{
    pid;          <- process id of thread becoming ready
    tid;          <- Thread id.
    priority;     <- priority of the thread
}

```

Other supported built-ins

__errno__kernelmode, __arg1 to __arg7, __curthread, __curproc, __mst, __tid, __pid, __ppid, __pgid, __uid, __euid.

Execution environment

Runs in process or interrupt environment.

Use case: Trace time taken by a thread that is performing I/O operation to get back to ready queue.

```

@@BEGIN
{
    printf ("          Pid      Tid      Time          Delta\n");
}

@@sysproc:offreadyq :*
{
    ready[__tid] = timestamp();
    printf ("offreadyq: %llu %llu %W\n",__readyprocinfo->pid,__readyprocinfo->tid,ready[__tid]);
}

@@sysproc:onreadyq :*
{
    if (diff_time(ready[__tid],0,MICROSECONDS))
    {
        auto:diff = diff_time (ready[__tid],timestamp(),MICROSECONDS);
        printf ("onreadyq : %llu %llu %W          %llu\n",__readyprocinfo->pid,__readyprocinfo->tid,ready[__tid],diff);
        delete (ready,__tid);
    }
}

```

An output similar to the following sample is displayed:

	Pid	Tid	Time	Delta
offreadyq:	7799280	20709717	5s 679697µs	
onreadyq :	7799280	20709717	5s 679697µs	6
offreadyq:	7799280	20709717	5s 908716µs	
onreadyq :	7799280	20709717	5s 908716µs	3
offreadyq:	7799280	20709717	6s 680186µs	
onreadyq :	7799280	20709717	6s 680186µs	5
offreadyq:	7799280	20709717	6s 710720µs	
onreadyq :	7799280	20709717	6s 710720µs	4
offreadyq:	7799280	20709717	6s 800720µs	
onreadyq :	7799280	20709717	6s 800720µs	2
offreadyq:	7799280	20709717	6s 882231µs	
onreadyq :	7799280	20709717	6s 882231µs	2
offreadyq:	7799280	20709717	6s 962313µs	
onreadyq :	7799280	20709717	6s 962313µs	2
offreadyq:	7799280	20709717	6s 980311µs	
onreadyq :	7799280	20709717	6s 980311µs	2

onreadyq

This probe starts when a thread is enqueued to system ready queue or its position in ready queue is modified.

Syntax:@@sysproc:offreadyq:<pid/tid/*>

Special built-ins supported

```
__readyprocinfo{
  pid;          <- process id of thread becoming ready
  tid;          <- Thread id.
  priority;     <- priority of the thread
}
```

Other supported built-ins

__errno__kernelmode, __arg1 to __arg7, __curthread, __curproc, __mst, __tid, __pid, __ppid, __pgid, __uid, __euid.

Execution environment

Runs in process or interrupt environment.

dispatch

This probe starts when system dispatcher is called to select a thread to run on a specific CPU.

Syntax:@@sysproc:dispatch:<pid/tid/*>

Special built-in supported

```
__dispatchinfo{
  cpuid;        <- CPU where selected thread will run.
  oldpid;       <- pid of the thread currently running
  oldtid;       <- thread id of the thread currently running
  oldpriority;  <- priority of the thread currently running
  newpid;       <- pid of the new process process selected for running
  newtid;       <- thread id of the thread selected for running
  newpriority;  <-priority of the thread selected for running
}
```

Other supported built-ins

__errno__kernelmode, __arg1 to __arg7, __curthread, __curproc, __mst, __tid, __pid, __ppid, __pgid, __uid, __euid.

Execution environment

Runs in interrupt environment only.

Example

```
print process thread id of old and selected thread on CPU '0' with dispatch time relative to start of the script
```

```
echo '@@sysproc:dispatch:* when (__cpuid == 0){printf ("%llu %llu %W\n",__dispatchinfo->oldtid,__dispatchinfo->newtid,timestamp());}' |
probevue
```

An output similar to the following example is displayed:

```
24641983 20709717 0s 48126µs
20709717 23593357 0s 48164µs
23593357 20709717 0s 48185µs
20709717 23593357 0s 48214µs
23593357 20709717 0s 48230µs
20709717 23593357 0s 48288µs
23593357 261 0s 48303µs
261 20709717 0s 48399µs
```

Example II

Time spent on CPU '0' by threads in between dispatch event.

```
@@BEGIN
{
    printf ("Thread cpu Time-Spent\n");
}

@@sysproc:dispatch:* when (__cpuid == $1)
{
    if (savetime[__cpuid] != 0)
        auto:diff = diff_time (savetime[__cpuid],timestamp(),MICROSECONDS);
    else
        diff = 0;
    savetime[__cpuid] = timestamp();
    printf ("%llu %llu %llu\n",__dispatchinfo->oldtid,__dispatchinfo->cpuid,diff);
}

# probevue cputime.e 6
Thread cpu Time-Spent
3146085 6 0
3146085 6 9995
3146085 6 10002
3146085 6 10008
3146085 6 99988
3146085 6 100006
3146085 6 99995
3146085 6 99989
3146085 6 100010
3146085 6 100001
3146085 6 100000
3146085 6 99998
```

As can be observed thread 3146085 is being re-dispatched on the CPU at an interval of 1sec in absence of any other thread competing for this CPU.

oncpu

This probe starts when a new process or thread acquires CPU.

Syntax:@@sysproc:oncpu:<pid/tid/*>

Where pid is process identifier and tid is thread identifier of process or thread that is acquiring the CPU.

Special built-ins supported

```
__dispatchinfo{
    cpuid;          <- CPU where selected thread will run.
    newpid;        <- pid of the new process process selected for running
    newtid;        <- thread id of the thread selected for running
    newpriority;   <-priority of the thread selected for running
}
```

Other supported built-ins

__errno__kernelmode, __arg1 to __arg7, __curthread, __curproc, __mst, __tid, __pid, __ppid, __pgid, __uid, __euid.

Execution environment

Runs in interrupt environment only.

Example

```
To print time spent by threads of sysncd on all CPU's
#!/usr/bin/probevue

@@BEGIN
{
    printf ("PROCESSID THREADID CPU TIME\n");
}

@@sysproc:oncpu:$1
{
    savetime[__cpuid] = timestamp();
}

@@sysproc:offcpu:$1
{
    if (savetime[__cpuid] != 0)
        auto:diff = diff_time (savetime[__cpuid],timestamp(),MICROSECONDS);
    else
        diff = 0;
    printf ("%llu %llu %llu %llu\n",
        __dispatchinfo->oldpid,
        __dispatchinfo->oldtid,
        __dispatchinfo->cpuid,
        diff);
}

# cputime.e `ps aux|grep sysncd| grep -v grep| cut -f 6 -d " "`
```

An output similar to the following example is displayed:

```
3735998 18612541 0 2
3735998 15663427 0 1
3735998 15073557 0 1
3735998 18743617 0 1
3735998 18874693 0 1
3735998 18809155 0 15
3735998 18940231 0 20
3735998 18547003 0 1
3735998 19267921 0 1
3735998 19071307 0 17
3735998 18678079 0 1
3735998 18481465 0 1
3735998 19202383 0 15
3735998 19005769 0 1
3735998 19136845 0 19
3735998 6160689 0 190
```

offcpu

This probe starts when a process or thread is dispatched from a CPU.

Syntax:@@sysproc:dispatch:<pid/tid/*>

Special built-ins supported

```
__dispatchinfo{
  cpuid;      <- CPU where selected thread will run.
  newpid;     <- pid of the new process process selected for running
  newtid;     <- thread id of the thread selected for running
  newpriority; <- priority of the thread selected for running
}
```

Other supported built-ins

__errno__kernelmode, __arg1 to __arg7, __curthread, __curproc, __mst, __tid, __pid, __ppid, __pgid, __uid, __euid.

Execution environment

Runs in interrupt environment only.

blockthread

This probe starts when a thread is blocked from running on a CPU. Blocking is a form of sleeping when a thread sleeps without holding any resources.

Syntax: @@sysproc:blockthread:*

Special built-ins supported

```
__sleepinfo{
  pid;
  tid;
  waitchan; <-- wait channel of this sleep.
}
```

Other supported built-ins

__errno__kernelmode, __arg1 to __arg7, __curthread, __curproc, __mst, __tid, __pid, __ppid, __pgid, __uid, __euid.

Execution environment

Runs in interrupt environment only.

foldcpu

This probe starts when a CPU core is about to be folded. This probe does not happen in process context and must not be filtered with a pid or tid.

Syntax: @@sysproc:foldcpu:*

Special built-ins supported

```
__foldcpuinfo{
  cpuid;      <- logical cpu id which triggers core folding
  gpcores;    <- general purpose (unfolded, non-exclusive) cores available.
}
```

Other supported built-ins

__errno__kernelmode, __arg1 to __arg7.

Example:

To track all CPU folding events in the system:

```
__foldcpuinfo{
  cpuid;      <- logical cpu id which triggers core folding
  gpcores;    <- general purpose (unfolded, non-exclusive) cores available.
}
```

bindprocessor

Syntax: @@sysproc:bindprocessor:<pid/tid/*>

This probe starts when a thread or process is bound to a CPU. Bindprocessor is a permanent event and must not be confused with temporary CPU switches.

Special built-ins supported

```
__bindprocessorinfo{
  ispid      <- 1 if cpu is bound to process; 0 for a thread
  id;        <- thread or process id.

  cpuid;

};
```

Other supported built-ins

__errno__kernelmode, __arg1 to __arg7, __curthread, __curproc, __mst, __tid, __pid, __ppid, __pgid, __uid, __euid.

Execution environment

Runs in process environment.

changecpu

This probe starts when a thread changes CPU temporarily. This event is more likely to be captured during CPU funneling events or intentional jumps of some kproc events to perform CPU related tasks (the xmgc process jumps to all CPUs to manage kernel heaps) special built-ins.

Syntax: @@sysproc:changecpu:*>

Special built-ins supported

```
__changecpuinfo
{
  oldcpuid;    <-source CPU
  newcpuid;    <- target CPU
  pid;
  tid;         <-Thread id
}
```

Other supported builtins

__errno__kernelmode, __arg1 to __arg7, __curthread, __curproc, __mst, __tid, __pid, __ppid, __pgid, __uid, __euid.

Execution environment

Runs in process environment.

Example

```
@@sysproc:changecpu:*
{
  printf ("changecpu PID=%llu TID=%llu old_cpuid=%d new_cpuid= %d \n",
    __changecpuinfo->pid,__changecpuinfo->tid,__changecpuinfo->oldcpuid,__changecpuinfo->newcpuid);
}
```

An output similar to the following example is displayed:

```

changecpu PID=852254 TID=1769787 old_cpuid=26 new_cpuid= 27
changecpu PID=852254 TID=1769787 old_cpuid=-1 new_cpuid= 0
changecpu PID=852254 TID=1769787 old_cpuid=0 new_cpuid= 1
changecpu PID=852254 TID=1769787 old_cpuid=1 new_cpuid= 2

```

resourceattach

This probe is fired when a resource is attached to another resource in the system.

Syntax: @@sysproc:resourceattach:*>

Special built-ins supported

```

__srcresourceinfo{
    type;
    subtype;
    id;          <- resource type identifier
    offset;     <-offset if a memory resource
    length;    <- length if a memory resource
    policy;
}
__tgtresourceinfo{
    type;
    subtype;
    id;          <- resource type identifier
    offset;     <-offset if a memory resource
    length;    <- length if a memory resource
    policy;
}

```

Where type and subtype could be have one of the following values.

<i>Table 73. The resourceattach probe: type and subtype values</i>	
Resource type	Description
R_NADA	Nothing - invalid specification
R_PROCESS	Process
R_RSET	Resource set
R_SUBRANGE	Memory range
R_SHM	Shared Memory
R_FILDES	File identified by an open file
R_THREAD	Thread
R_SRADID	SRAD identifier
R_PROCMEM	Process Memory

Other supported builtins

__errno__kernelmode, __arg1 to __arg7, __mst.

Execution environment

Runs in process environment.

resourcedetach

This probe is fired when a resource is detached from another resource in the system.

Syntax: @@sysproc:resourcedetach:*>

Special built-ins supported

```
__srcresourceinfo{
    type;
    subtype;
    id;          <- resource type identifier
    offset;     <-offset if a memory resource
    length;    <- length if a memory resource
    policy;
}

__tgtresourceinfo{
    type;
    subtype;
    id;          <- resource type identifier
    offset;     <-offset if a memory resource
    length;    <- length if a memory resource
    policy;
}
```

Where type and subtype could be have one of the following values.

Resource type	Description
R_NADA	Nothing - invalid specification
R_PROCESS	Process
R_RSET	Resource set
R_SUBRANGE	Memory range
R_SHM	Shared Memory
R_FILDES	File identified by an open file
R_THREAD	Thread
R_SRADID	SRAD identifier
R_PROCMEM	Process Memory

Other supported builtins

__errno__kernelmode, __arg1 to __arg7, __mst, __tid, __pname.

Execution environment

Runs in process environment.

drphasestart

This probe is fired when a dr handler is about to be called.

Syntax: @@sysproc:drphasestart:*

Special built-ins supported

```
__drphaseinfo{
    dr_operation;  <- dr operation
    dr_flags;
    dr_phase;
    handler_rc;   <- always 0 in drphasestart
}
```

dr_operation can have one of the following values:

- DR operation
- DR_RM_MEM_OPER

- DR_ADD_MEM_OPER
- DR_RM_CPU_OPER
- DR_ADD_CPU_OPER
- DR_CPU_SPARE_OPER
- DR_RM_CAP_OPER
- DR_ADD_CAP_OPER
- DR_RM_RESMEM_OPER
- DR_PMIG_OPER
- DR_WMIG_OPER
- DR_WMIG_CHECKPOINT_OPER
- DR_WMIG_RESTART_OPER
- DR_SOFT_RES_CHANGES_OPER
- DR_ADD_MEM_CAP_OPER
- DR_RM_MEM_CAP_OPER
- DR_CPU_AFFINITY_REFRESH_OPER
- DR_AME_FACTOR_OPER
- DR_PHIB_OPER
- DR_ACC_OPER
- DR_CHLMB_OPER
- DR_ADD_RESMEM_OPER

dr flags can be a combination of the following values:

- Flag
- DRP_FORCE
- DRP_RPDP
- DRP_DOIT_SUCCESS
- DRP_PRE_REGISTERED
- DRP_CPU DRP_MEM DRP_SPARE
- DRP_ENT_CAP
- DRP_VAR_WGT
- DRP_RESERVE
- DRP_PMIG DRP_WMIG
- DRP_WMIG_CHECKPOINT
- DRP_WMIG_RESTART
- DRP_SOFT_RES_CHANGES
- DRP_MEM_ENT_CAP
- DRP_MEM_VAR_WGT
- DRP_CPU_AFFINITY_REFRESH
- DRP_AME_FACTOR
- DRP_PHIB
- DRP_ACC_UPDATE
- DRP_CHLMB

Other supported builtins

`__errno__ kernelmode, __arg1 to __arg7, __tid`

Execution environment

Runs in process or interrupt environment.

Example

Shell scripts for ProbeVue

Shell scripts are useful when running ProbeVue.

The following are useful shell scripts:

sprobevue

Shell script that wraps all arguments in double quotes:

```
#!/usr/bin/ksh
#
# sprobevue:
#
#   Simple helper function for probevue
#   Wraps arguments to probevue in double quotes
#
#   Usage: sprobevue <probevue flags> <script> <args>
#   Doesn't support the -c and -A flags of probevue
#
usage()
{
    echo "Usage: sprobevue <probevue flags> <script> <args>" >&2
    echo "   Doesn't support the -c and -A flags of probevue" >&2
    exit 1
}

CMD=probevue
# Generate command to execute

while getopts 'c:A:I:s:o:t:X:' zargs
do
    case $zargs in
        I|s|o|t|X) CMD="$CMD -$zargs $OPTARG"      ;;
        ?) usage
    esac
done

shift $((OPTIND -1))

if [ -n "$1" ]
then
    CMD="$CMD $1"
    shift
fi

for i
do
    CMD="$CMD \"\$i\""
done

# Execute command
$CMD
```

prgrep

Shell script that prints process ID given process name:

```
#!/usr/bin/ksh
#
# prgrep:
#
#   Simple helper function for probevue
#   Prints all process IDs with given process name
#
#   Need options to print only one process
#   to print process belong to a certain UID
#
#   Usage: prgrep <process_name>
#   prgrep -p <processID>
```

```

#

usage()
{
    echo "Usage: prgrep <process_name>" >&2
    echo "          prgrep -p <process_ID>" >&2
    exit 1
}

[ -z "$1" ] && usage

if [ $1 = "-p" ]
then
    [ -z "$2" ] && usage
    pid=$2
    export pid
    ps -e | awk 'BEGIN {pid = ENVIRON["pid"]} {if ($1 == pid) print $4}'
else
    pname=$1
    export pname
    ps -e | awk 'BEGIN {pname = ENVIRON["pname"]} {if ($4 == pname) print $1}'
fi

```

ProbeVue error messages

As described earlier, running the **probevue** command requires privilege. If an ordinary user tries running the **probevue** command, the RBAC framework detects this and fails the execution of the command immediately.

```

$ probevue kernel.e
ksh: probevue: 0403-006 Execute permission denied.

```

The *Authorizations and privileges* section of the [Running ProbeVue](#) describes how to enable non-root users with the authorizations and privileges to issue the **probevue** command.

The ProbeVue compiler, which is built-in to the **probevue** command, prints detailed error messages during the compilation phase when it detects any syntax errors, semantic errors or type incompatibility errors. Consider the following script:

```

/* Syntax error example:
 * syntaxbug.e
 */
@@BEGIN
{
    int i, j, k;

    i = 4;
    j = 22;

    k = i _ z;

    printf("k = %d\n", k);

    exit();
}

```

The preceding script has a syntax error on line 11 at column 15, the assignment statement. Instead of a minus symbol (-) or an underscore symbol (_) was typed by mistake. On running the script, the ProbeVue compiler catches this error and generates an error message:

```

# probevue syntaxbug.e
syntaxbug.e: token between line 11: column 15 and line11: column 15: , expected
instead of this token

```

The ProbeVue compiler also invokes internal system calls to check if the probe specifications in the Vue script are valid. A common error is to pass an invalid process ID or the process ID of an exited process

in the probe point tuple. Another common error is to forget to pass a process ID as an argument on the command line when the script expects one. Consider the following script:

```
/* simpleprobe.e
 */
@@syscall:$1:read:entry
{
    printf("In read system call: thread ID = %d\n", __tid);
    exit();
}
```

The preceding script requires a process ID as an argument to replace the '\$1' variable in the probe point tuple at line 3. The kernel will return an error if you try to probe a process that has exited or does not exist. It also fails if the process ID indicates a kernel process or the init process. Further, you cannot probe a process that does not belong to you unless you have the required privileges to probe another user's processes. You can use the **prgrep** command with the **-p** flag to print the process name given a process ID.

Note: This command produces an empty output if the specified process ID does not exist.

```
# probevue simpleprobe.e 233
probevue: The process does not exist.
ERR-19: Line:3 Column:3 Invalid probe string
# prgrep -p 232
#
# probevue simpleprobe.e 1
ERR-19: Line:3 Column:3 Invalid probe string
# prgrep -p 1
init
# probevue simpleprobe.e
ERR-19: Line:3 Column:3 Invalid probe string
```

The **probevue** command can also detect if an unprivileged user tries to access kernel variables. Consider the `kernel.e` script from the sample programs section. The following example session shows what happens if you try running this as an unprivileged user:

```
$ probevue kernel.e
ERR-56: Line:93 Column:39 No authority to access kernel variable
ERR-56: Line:99 Column:23 No authority to access kernel variable
ERR-56: Line:100 Column:24 No authority to access kernel variable
ERR-56: Line:101 Column:25 No authority to access kernel variable
ERR-56: Line:102 Column:24 No authority to access kernel variable
ERR-102: Line:140 Column:13 Operation not allowed
ERR-46: Line:140 Column:9 Invalid Assignment, Type mismatch
```

After the Vue script has been compiled successfully, the **probevue** command invokes a system call to start a new ProbeVue session passing the intermediate code generated by the compiler. The system call will fail if the ProbeVue framework fails to initialize a new ProbeVue session. There can be several reasons for this. For example, starting the new session can cause memory resources for the user to exceed the administrator-specified limits. The session can need more memory resources than allowed for a single session. There can be unauthorized functions used in the interval probe manager. One of the processes being probed can have exited after the compilation phase checks were made. When the session cannot be started, the kernel fails the system call returning a unique 64-bit error.

The ProbeVue framework can abort a successfully started and active ProbeVue session if a severe or unrecoverable error is encountered while issuing the probe actions. Possible errors include exceeding session or user memory limits (memory requirements for thread-local variables and list variables can grow as the session progresses), exceeding temporary string or stack area limits, accessing out-of-array indexes, attempting to divide by zero, and so on. In all cases, the kernel will return a unique 64-bit error number while terminating the session.

When the session is failed whether at start or after it has been successfully started, the **probevue** command prints a generic error message including the unique 64-bit error number in hexadecimal format and exits. The following chart provides the meaning of some common 64-bit errors that could be returned by the kernel:

Kernel error	Meaning	Occurs at
0xEEEE00008C285034	Out of memory while allocating primary trace buffers.	Session start
0xEEEE00008C285035	Out of memory while allocating secondary trace buffers.	Session start
0xEEEE00008C52002B	Out of memory while allocating storage for probe specification strings.	Session start
0xEEEE000096284122	Out of memory while allocating storage for thread-local storage.	Session start
0xEEEE000081284049	Use of user-space access functions in the interval probe manager.	Session start
0xEEEE0000D3520022	Number of sessions limit for regular users.	Session start
0xEEEE000096284131	Illegal address passed to the get_userstring function.	Executing probe action
0xEEEE00008C520145	Maximum thread limit hit for thread local variables.	Executing probe action

RAS events functions

The "RAS events" functions are a privileged set of Vue functions provided for very specialized system or application debugging purposes.

The "RAS events" functions are a privileged set of Vue functions provided for very specialized system or application debugging purposes. They are not intended for general use. They provide system tracing and dumping facilities. Many of these functions are "pass through" functions that allow a Vue script to directly invoke kernel services, and hence there are risks involved with using them. You need special privileges to successfully invoke these functions in your Vue script: you must either be root or have the **aix.ras.probevue.rase** authorization.

To avoid the risk of these functions, pass the **-K** flag to the **probevue** command. Otherwise, these functions simply disappear from the Vue language completely.

Generating a trace record

The Vue functions for generating system trace (and LMT trace) records have the syntax that is similar to the kernel interfaces that are invoked by the Vue functions.

There are some restrictions as follows:

- If system trace is not started, or the *hookid* value is not being captured by system trace, these operations do not produce system trace records (LMT tracing to the common buffer for TRCHKLx traces will still be attempted, but LMT might also be disabled).
- You cannot generate a trace record from within a **@@systrace** Vue clause. Calls to the tracing functions only generate the LMT common buffer trace records for TRCHKLx traces in this case, assuming that LMT is enabled.
- You cannot probe these ProbeVue-generated trace events; only kernel and application generated tracing can be probed.
- You must be privileged, either as root or with the **aix.ras.probevue.rase** authorization.

The following Vue functions exist for writing system trace records. All data words are of type **long long integers**:

TRCHKLO(hookID)

Trace with no data words.

TRCHKL1(hookID, D1)

Trace with 1 data word.

TRCHKL2(hookID, D1,D2)

Trace with 2 data words.

TRCHKL3(hookID, D1,D2,D3)

Trace with 3 data words.

TRCHKL4(hookID, D1,D2,D3,D4)

Trace with 4 data words.

TRCHKL5(hookID, D1,D2,D3,D4,D5)

Trace with 5 data words.

void trcgenk(int channel, int hook_ID, unsigned long long data_word, int length, untyped buffer)

Trace a buffer.

These trace functions always append a timestamp to the event data. The *hookid* parameter to these functions is of the form 0xhhhh0000. This does not mean that the *hookid* value is required to be a constant, it just indicates how a *hookid* value is formed.

Note: Obsolete 12-bit *hookid* values will use the leftmost three hex digits, and the 4th digit will be zero.

With the **trcgenk** kernel service, the *buffer* parameter is a pointer to length bytes of data to trace, at most 4096 bytes. The *buffer* parameter can be an external variable like a kernel or application pointer to pinned data, or a script variable like a Vue string or structure instance. The "untyped" specification is a shorthand for this.

Note: The **trcgenk** kernel service only traces to the system trace, not to the LMT trace buffers.

You can use a non-zero channel number, but you must ensure that the specified channel is enabled for tracing. The return value from the **trace** command that started the trace of interest can be passed to the Vue script for this purpose. Using a disabled channel will result in no tracing.

These tracing functions do not return a value.

Stopping the trace

To freeze the system trace as soon as possible after a required event has occurred, you can use `void trcoff()` in a Vue script. This function disables channel zero tracing immediately. You must still stop the trace in the normal way, with the **trcstop** command external to ProbeVue, in order for trace processing to be completed normally.

You can immediately stop LMT and component traces so that ongoing tracing does not wrap data of interest. The corresponding resuming functions are needed because there is no command line equivalent available to restart these traces. There are following new Vue functions:

```
void mtrcsuspend()
void ctsuspend()
void mtrcresume()
void ctresume()
```

The **ctsuspend** routine stops all component tracing. You cannot use this routine for selective trace stop by component. It stops component trace only, not any other tracing that the CT_HOOKx macros might have requested, such as system and LMT trace recording.

You must use these trace control functions with caution, as there is no serialization of the kernel tracing code being affected. You must manually ensure that only one script or command will be affecting tracing at a time.

Stopping the system

You can terminate the system and take a full dump using the following routine:

```
void abend(long long code, long long data_word, ...)
```

This routine is similar to the **abend** kernel service, except that only up to 7 data parameters (which will be loaded into registers r3 through up to r10) are accepted here.

Untyped parameters

In function prototypes to follow, some parameters of the equivalent kernel functions are typed ambiguously. The Vue compiler generally performs type checking on all parameters passed to a Vue function, but the parameters designated as having an "untyped" type are exempted from type checking. For example, an optional string might be passed as NULL when using these kernel services directly in the kernel, but if the Vue function was defined as taking a parameter of type String, a NULL cannot be accepted. To avoid the inconvenience of having to pass an empty string instead and to let the Vue functions take the same parameters as the following kernel interface, these functions have been defined as taking untyped parameters. An untyped parameter provides us with the liberty of passing NULL instead of a real Vue string, but be careful when specifying values for "untyped" parameters, because the compiler will accept any type for the parameter.

Note: There is really no "untyped" variable specification in the Vue language. It is just used as a shorthand notation.

Related information

[trcgenk Kernel Service](#)

[Macros for recording trace events](#)

Taking a live dump

For creating the kernel live dump, ProbeVue services are used and the ProbeVue services are similar the corresponding kernel services.

An exception to this general similarity is the **ldmp_parms** structure, which is not exposed at the script level. Instead, the **ldmp_setupparms** built-in function owns a private instance of this structure, which is allocated and returned to the caller indirectly as a 64 bit cookie which must be passed to subsequent live dump services in its place. Only one session can use the private structure at a time. You can use the other live dump services to resemble the syntax of their kernel counterparts. Because of this hidden allocation (as well as hidden allocations made by the kernel live dump services themselves), it is necessary to call either the **ldmp_freeparms** kernel service or the **livedump** kernel service once the **ldmp_setupparms** kernel service has been called and returned successfully. Otherwise, the current session will continue to own the private structure, causing all future **ldmp_setupparms** calls to fail. After the private structure has been released, it can no longer be used by its former owner without another **ldmp_setupparms** call. Do not use the LDT_POST flag with the **ldmp_setupparms** kernel service, as that implies an unsupported future reference to the hidden structure.

The typical live dump application must hold the hidden structure for only a very brief interval, typically within a single **probevue** clause. The hidden structure is owned by the session, and can actually be used by any Vue clause in that session. The framework will release the private structure, and other kernel resources, with the **ldmp_freeparms** kernel service automatically when the ProbeVue session terminates.

As the **ldmp_parms** structure elements are not visible to ProbeVue, those that require or permit initialization by the caller are set using extra parameters passed to the ProbeVue version of the **ldmp_setupparms** kernel service instead.

```
long long ldmp_setupparms(String symptom,required symptom string
untyped title, dump title string or NULL
untyped prefix, dump file name prefix string or NULL
untyped func, failing function name string or NULL
long long errcode, error code
int flags, dump characteristics
int prio ) dump priority
```

The preceding **ldmp_setupparms** Vue function is an interface to the kernel service of the same name, except that the **ldmp_parms** structure is not visible to the calling Vue script. The value returned must be passed to the other live dump services as a substitute for the pointer to an **ldmp_parms** structure, although it is typed as a 64 bit integer.

The *symptom* string is a required String operand, while the *title*, *prefix*, and *func* strings are optional. Pass either a String or NULL for these three parameters. All String values must be local to the Vue script. The *flags* and *prio* parameters can be zero, or values from the kernel header file **sys/livedump.h**. The appropriate integer constants must be used here, although there is an alternative.

The following values are useful values for the *flags* parameter:

```
LDT_ONEPASS      0x02    limit dump to one pass
LDT_NOADDCOMPS  0x08    components can't be added by callbacks
LDT_NOLOG        0x10    no error is to be logged
LDT_FORCE        0x20    force this dump
```

Because the dump will be taken from ProbeVue's disabled internal environment, it must be a serialized, synchronous, one pass dump.

The following values are acceptable values for the *prio* parameter:

```
LDPP_INFO        1    informational dump
LDPP_CRITICAL    7    critical dump (this is the default)
```

If zero is specified for the *prio* parameter, LDPP_CRITICAL is defaulted by the **ldmp_setupparms** kernel service. Only a non-zero value will be stored in the hidden **ldmp_parms** structure to override this.

The return value upon success will be a positive cookie representing the ownership of the hidden **ldmp_parms** structure.

On any failure, the return value will be negative as follows:

Value	Description
EINVAL_EVM_COOKIE	Indicates that the private ldmp_parms structure is not available.
EINVAL_EVM_STRING	Indicates that a String valued parameter is not valid.

All of the subsequently described Vue functions return failure indications in a similar fashion, with a negative kernel error number:

Value	Description
EINVAL_EVM_COOKIE	Indicates that the caller did not correctly specify a cookie showing ownership of the private ldmp_parms structure.
EINVAL_EVM_STRING	Indicates that a String valued parameter is not valid.
EINVAL_EVM_EXTID	Indicates that the <i>extid</i> parameter is not supported, and must be zero.

Other kernel error numbers can be passed back by the following kernel services:

long long ldmp_freeparms (long long cookie)

After the **ldmp_setupparms** kernel service has returned successfully, the internal **ldmp_parms** structure has been allocated to the running Vue script. You must free this resource, plus other kernel-internal resources allocated by the services that add components to your dump, by either taking the dump by calling the **livedump** kernel service, or by calling the **ldmp_freeparms** kernel service. This releases the internal **ldmp_parms** structure for future use.

long long livedump (long long cookie)

After the **ldmp_setupparms** kernel service and at least one of the various services that add components (and pseudo-components) to a dump have been called, the dump is requested by the **livedump** service. This service produces the actual live dump in the **/var/adm/ras/livedump** file according to the specifications provided through the **ldmp_setupparms** kernel service and the other live dump services invoked. The *cookie* parameter is the cookie returned by the initial call to the

ldmp_setupparms kernel service. The return value will be zero if the dump was successfully taken, `EINVAL_EVM_COOKIE` if the cookie is not valid, and another kernel error number if an error occurs during kernel **livedump** processing.

long long dmp_compspec(long long flags, DCF_XXX flags defined in sys/dump.h untyped comp, component to be added (by ras_block_t, name, alias, and so on.) long long cookie, cookie returned by ldmp_setupparms long long extid, not supported – must be zero untyped p1, first possible component parameter ...); additional component parameters

You can add any component that supports live dump to the live dump by calling this service, which is identical in function to the kernel service with the same name except the following situations:

- The *extid* parameter, which allows a `dmp_extid_t` (long) to be returned in the kernel programming environment, is unsupported and must be zero. `EINVAL_EVM_EXTID` will be returned otherwise. There is no way to pass a pointer to ProbeVue memory to receive this value, which might then be used with the **dmp_compext** kernel service, which is therefore also not supported in ProbeVue. Instead, you can call the **dmp_compspec** service multiple times.
- The kernel service allows any number of parameters *p1*, *p2*, and so on where an additional NULL parameter must follow the last actual one to terminate the parameter list. The Vue function only accepts at most four parameters of *p1*, *p2*, and so on. The last must still be zero to tell the kernel service how many of these parameters there are, so in effect you can specify only up to 3 interesting values. The interface will automatically force the parameter following the last one of up to 3 variable parameters to zero, to ensure that this rule is followed.
- The *comp* parameter can be a long, a kernel `ras_block_t` address, or a String as appropriate. The type is not checked.
- The kernel **#define** flag values are not part of ProbeVue.

long long ras_block_lookup(String path)

This function locates the `ras_block_t` corresponding to the component path name parameter. This can be useful for calling the **dmp_ct** kernel service, which requires such an address, if you cannot more easily find the address in a kernel variable.

The return value from this function is either the kernel address of the requested `ras_block_t`, or NULL if the `ras_block_t` cannot be found.

The following functions are all simple "pass through" functions that allow a Vue script to directly invoke the corresponding kernel services. Some of the parameter lists have unused members for compatibility with the kernel, so you can use the kernel documentation directly. The value 0 must be passed for parameters shown as unused. You can use these services in the same way as their kernel counterparts, except that the address of an **ldmp_parms** structure is replaced by the cookie returned from the **ldmp_setupparms** kernel service.

As always, a negative return value indicates an error. This can be a kernel error number from the following kernel service, or from the interface routines if the cookie or a string is incorrect. The following interfaces provide most of the flexibility available to kernel or kernel extension driven live dumps.

```
long long dmp_context( long long flags, DCF_XXX flags from dump.h
long long cookie,      cookie returned by ldmp_setupparms
long long name,        unused by this function
long long ctx_type,    DMP_CTX_XXX flags from dump.h untyped p2)
parameter dependent on ctx_type (NULL, mst addr, cpuid, tid)
```

```
long long dmp_ct(      long long flags, DCF_XXX flags from dump.h
long long cookie,      cookie returned by ldmp_setupparms
long long name,        unused by this function
untyped rasb,          component's ras_block_t pointer
long long size)        amount of CT buffer to dump or 0 for all
```

```
long long dmp_eaddr(  long long flags, DCF_XXX flags from dump.h
long long cookie,      cookie returned by ldmp_setupparms
String name,          cdt name
untyped addr,         first address to dump
long long size)       number of bytes to dump
```

```

long long dmp_errbuf(long long flags, DCF_xxx flags from dump.h
long long cookie, cookie returned by ldmp_setupparms
long long name, unused by this function
long long erridx, 0 for global error log, or wpar id
long long p2) unused

```

```

long long dmp_mtrc(long long flags, DCF_xxx flags from dump.h
long long cookie, cookie returned by ldmp_setupparms
long long name, unused by this function
long long com_size, amount of LMT common data to dump
long long rare_size) amount of LMT rare data to dump

```

```

long long dmp_pid( long long flags, DCF_xxx flags from dump.h
long long cookie, cookie returned by ldmp_setupparms
long long name, unused by this function
long long pid, id of process to dump
long long p2) unused

```

```

long long dmp_systrace (long long flags, DCF_xxx flags from dump.h
long long cookie, cookie returned by ldmp_setupparms
long long name, unused by this function
long long size, amount to dump
long long p2) unused

```

```

long long dmp_tid( long long flags, DCF_xxx flags from dump.h
long long cookie, cookie returned by ldmp_setupparms
long long name, unused by this function
long long tid, id of thread to dump
long long p2) unused

```

Note: You must call the **ldmp_freeparms** kernel service after any failure in the preceding routines, assuming you then want to abandon the dump.

The following script is an example that takes a very small, simple live dump. The kernel symbol `dc_data` exports a structure from the kernel, whose actual shape and contents are of no importance to this example.

```

_kernel struct {int i1; int i2; int i3; int i4;} dc_data;

@@BEGIN
{
    long long ldmp_parms;
    long long rc;

    rc = ldmp_setupparms( "dc_data dump",
                        "My Sample Dump", /* dump title */
                        "pvdump", /* dump path prefix */
                        NULL, /* no function name */
                        0x1122334455667788LL, /* error code */
                        0x10, /* LDT_NOLOG flag */
                        0); /* default dump prio */
    printf("ldmp_setupparms rc = %016llx\n", rc);
    if (rc < 0) {
        exit();
    }

    ldmp_parms = rc; /* cookie for other livedump functions */

    /*
     * Add 16 bytes of kernel data to sample dump.
     * Note that "dc_data" passes the structure's address.
     */
    rc = dmp_eaddr(0, ldmp_parms, "dc_data", dc_data, sizeof(dc_data));
    if (rc) {
        printf("dump_eaddr failed: %llx\n", rc);
        ldmp_freeparms(ldmp_parms);
        exit();
    }

    /*
     * Take the sample live dump.
     */
    rc = livedump(ldmp_parms);
    if (rc) {

```

```

    printf("livedump failed: %llx\n", rc);
}
exit();
}

```

Using #define symbols for live dump flags

The following sample shell script, `probe.dump`, can be useful if you prefer to use the actual defined symbols for live dump flags rather than manually substituting from the header files. It captures the relevant definitions from the `livedump.h` and `dump.h` files, and uses the C preprocessor to substitute values for you before passing your script to ProbeVue. Your script must comply with the following rules:

- Must not begin with a `#!/usr/bin/probevue` comment.
- Must not use symbols beginning with `LDPP_`, `LDT_`, `DCF_`, or `DMP_` in conflict with the definitions in the header files.

Do not create files named `pvdump.*`, as the following script will overwrite them.

```

#!/bin/ksh
#
# Helper script for Vue scripts that need to pick up
# the values of the various flags used by livedump.
#
# The Vue script $1
# must not contain a "#!/usr/bin/probevue" comment because
# the C preprocessor doesn't like it.

sed -n \
-e '/(d' \
-e '/^#define LDPP_/p' \
-e '/^#define LDT_/p' \
-e '/^#define DCF_/p' \
-e '/^#define DMP_CTX_/p' \
/usr/include/sys/dump.h \
/usr/include/sys/livedump.h \
> pvdump.h

echo "#include \"pvdump.h\"" > pvdump.c
cat $1 >> pvdump.c
cc -P pvdump.c
/usr/bin/probevue -K pvdump.i
rm pvdump.[cih]

```

Related information

[livedump Kernel Service](#)

Vue functions

There are many Vue functions. For example, you can use Vue functions to append values to a list, generate a runtime stack trace, and return the sum of all elements in a list.

The Vue language supports the following list of functions:

Function	Description
<u>add_range</u>	Initialize the string range data type.
<u>append</u>	Appends a value to a list.
<u>atoi</u>	Returns the integer value of a string.
<u>avg</u>	Returns the average of all the elements in a list.
<u>commitdiscard</u>	Commits or discards the data in a tentative trace buffer.
<u>convert_ip4_addr</u>	Converts the IPv4 address (data) into the ProbeVue data-type <code>ip_addr_t</code> format.

Function	Description
<u>convert_ip6_addr</u>	Converts the IPv6 address (data) into the ProbeVue data-type ip_addr_t format.
<u>copy_kdata</u>	Copies data from kernel memory into a Vue script variable.
<u>copy_userdata</u>	Copies data from user memory into a Vue script variable.
<u>count</u>	Returns the count of the number of elements in a list.
<u>diff_time</u>	Returns the difference between two timestamps.
<u>eprintf</u>	Formats and prints data to the standard error.
<u>exit</u>	Terminates the Vue script.
<u>fd_fname</u>	Get the file name for a specific file descriptor.
<u>fd_fstype</u>	Get the file system type for a specific file descriptor.
<u>fd_ftype</u>	Get the file type for a specific file descriptor.
<u>fd_inodeid</u>	Get the inode ID for a specific file descriptor.
<u>fd_mpath</u>	Get the mount path of the file system for a specific file descriptor.
<u>fd_path</u>	Get the absolute file path for a specific file descriptor.
<u>fpath_inodeid</u>	Get the inode ID for a specific file path.
<u>get_function</u>	Returns the name of the function that is being probed.
<u>get_kstring</u>	Copies data from kernel memory into a String variable.
<u>get_location_point</u>	Returns the current probe location point.
<u>get_probe</u>	Returns the current probe point specification.
<u>get_stktrace</u>	Returns the runtime stack trace.
<u>get_kbytes</u>	Copies at most number of bytes of data specified, from kernel memory into a Vue script variable.
<u>get_ubytes</u>	Copies at most number of bytes of data specified, from user memory into a Vue script variable.
<u>get_userstring</u>	Copies data from user memory.
<u>list</u>	Creates and returns a new empty list.
<u>group_output_start</u>	Start output grouping for the clause that this function belongs to.
<u>group_output_end</u>	End output grouping for the clause that this function belongs to.
<u>lquantize</u>	Quantizes associative array values logarithmically and prints the key-value pairs in a graphical format.
<u>max</u>	Returns the maximum of all the elements in a list.
<u>min</u>	Returns the minimum of all the elements in a list.
<u>args</u>	Prints the name of the function being probed and its arguments
<u>print</u>	Prints the key-value pairs in an associative array
<u>printf</u>	Formats and copies data to the trace buffer.

Function	Description
<u>ptree</u>	Prints the process tree for the probed process
<u>quantize</u>	Quantizes associative array values linearly and prints the key-value pairs in a graphical format.
<u>qrange</u>	Finds the slot number of the range and adds it to the Associative array.
<u>round_trip_time</u>	Gets the smoothed round trip time for TCP connection for the specified socket descriptor.
<u>set_aso_print_options</u>	Specifies the sort-type, sort-by, and the <i>list-value</i> flags.
<u>set_date_format</u>	Updates the date format.
<u>set_range</u>	Initialize the linear and power range type.
<u>sockfd_netinfo</u>	Get the local and remote ports and IP addresses information for a specific socket descriptor.
<u>startend_tentative</u>	Indicates the start and the end of a tentative tracing section.
<u>stktrace</u>	Generates and prints a runtime stack trace.
<u>strstr</u>	Return a string inside another string.
<u>sum</u>	Returns the sum of all the elements in a list.
<u>timestamp</u>	Returns the current timestamp.
<u>trace</u>	Copies raw data to the trace buffer as hexadecimal characters.

add_range

Purpose

Initializes the string range data type and adds strings in a slot.

Syntax

```
add_range(range_t range_data, String S1, String S2, ..., String Sn);
```

Description

This routine initializes the `range_data` as a string range data type and also adds all the strings passed as arguments to the routine in one slot. If this routine is called the first time for a range data type, strings will be added in first slot. Otherwise, strings are added in the next slot.

Parameters

`range_data`

A `range_t` data type.

`S1, S2,...`

Strings to be added to the `range_data` parameter.

append

Purpose

Appends a value to a list.

Syntax

```
void append ( List listvar, long long val );
```

Description

The **append** function is the only list concatenation function available in Vue. It appends the value specified by the second parameter to the list variable specified by the first parameter. Each call to the **append** function adds a new value to the set of values already saved in the list variable and the size of the list variable grows. The append function also accepts another list as an argument, allowing you to join two lists.

Note: The value added to the list must be a parameter of integral or list type or this will result in a syntax error. The ProbeVue compiler will accept parameters that have any of the C-89 integer data types including both signed and unsigned types. No casts are needed.

The **append** function has no return value.

For more information about the list data type, see [“List type” on page 248](#). The preceding [list](#) section has an example script that uses the **append** function.

Parameters

Parameters	Description
<i>listvar</i>	Specifies a variable of type list .
<i>val</i>	Specifies a value or list to be appended.

atoi

Purpose

Returns the integer value of a string.

Syntax

```
int atoi( String str );
```

Description

The **atoi** function returns the integer whose value is represented by the string specified by the *str* parameter. It reads the string up to the first character that is not a numerical digit (0-9) and translates the scanned characters into the equivalent integer. Leading white-space characters are ignored, and an optional sign indicator can precede the digits.

The **atoi** function is useful in converting strings back to integers when running the sprobevue shell script that wraps double-quotes around all arguments. The following script is an example that captures a process forking faster than expected.

```
/* File: ffork.e
 *
 * Usage: sprobevue ffork.e processname delta
 *
 * Traces whenever a specified process is forking faster than
 * the "delta" value passed. Specify a process name and the time
 * in milliseconds as parameters.
 */

/* Ignore other parameters to execve */
int execve(char *path);

@@BEGIN
{
    int done;
    int pid;

    pname = $1;          /* name of process we are monitoring */

    /*
     * Since sprobevue is used, need to extract the integer value
     * from the string (double quotes around the delta).
     */
    delta = atoi($2);    /* minimum delta in millisecs between forks */
    printf("pname = %s, delta = %d\n", pname, delta);
}
}
```

```

@@syscall*:execve:entry
    when (done == 0)
    {
        __auto String exec[128];
        __thread int myproc;

        /* Find process being 'exec'ed */
        exec = get_userstring(__arg1, 128);

        /* Got it. Set a thread-local and reset 'done' so that we
         * avoid entering this probe from now on.
         */
        if (exec == pname) {
            pid = __pid;
            myproc = 1;
            done = 1;
            printf("Process name = %s, pid = %d\n", __pname, pid);
        }
    }

@@syscall*:fork:entry
    when (thread:myproc == 1)
    {
        /* old_ts is initialized to zero */
        probev_timestamp_t old_ts, new_ts;
        unsigned long long interval;

        /* Get current time stamp */
        new_ts = timestamp();

        /* Find time from last fork */
        if (old_ts != 0) {
            interval = diff_time(old_ts, new_ts, MILLISECONDS);

            /* if below the specified limit, trace that */
            if (interval < delta)
                printf("%s (%ld) forking too fast (%d milliseconds)\n",
                    pname, __pid, interval);
        }

        /* Save current fork timestamp */
        old_ts = new_ts;
    }

@@syscall*:exit:entry
    when (__pid == pid)
    {
        /* Catch process exit and terminate ourselves */
        printf("Process '%s' exited.\n", pname);
        exit();
    }

```

Parameter

Parameters	Description
<i>str</i>	Specifies the string to be converted.

avg

Purpose

Returns the average of all the elements in a list.

Syntax

```
long long avg ( List listvar );
```

Description

The **avg** function returns the average of all the elements that have been appended to the list variable specified by the *listvar* parameter.

Parameter

Parameters	Description
<i>listvar</i>	Specifies a variable of type list .

commit_tentative, discard_tentative

Purpose

Commits or discards the data in a tentative trace buffer.

Syntax

```
void commit_tentative( String bufID );  
void discard_tentative( String bufID );
```

Description

The **commit_tentative** function commits the trace data associated with the tentative trace buffer identified by the *bufID* parameter. This saves the data and makes it available to the trace consumer.

The **discard_tentative** function discards all the data in the tentative trace buffer indicated by the *bufID* parameter. This frees up the space in the trace buffers occupied by the tentative trace data.

When tentative trace data is being saved along with regular trace data, erstwhile tentative, but later committed trace data and regular trace data will be made available to the trace consumer in timestamp order. Thus, it is a good idea to commit or discard tentative data at the earliest opportunity to free up the trace buffers.

All tentative trace data that has not been committed is discarded when the ProbeVue session ends.

The [“Tentative tracing” on page 265](#) topic describes tentative tracing in more detail and includes an example Vue script that uses tentative tracing.

Parameter

Parameters	Description
<i>bufID</i>	Specifies a string constant that indicates the tentative trace buffer ID.

convert_ip4_addr

Purpose

Converts the IPV4 address (data) into ProbeVue IP address data-type format.

Syntax

```
ip_addr_t convert_ip4_addr (unsigned int ipv4_data);
```

Description

The **convert_ip4_addr** function converts the IPv4 address in the **in_addr** structure defined in the `/usr/include/netinet/in.h` file to ProbeVue IP address data-type **ip_addr_t**. This function returns the converted **ip_addr_t** value.

Parameters

ipv4_data	Specifies the ipv4 address data that needs to be converted into the ip_addr_t format.
------------------	--

convert_ip6_addr

Purpose

Converts the IPv6 address (data) into ProbeVue IP address data-type format.

Syntax

```
ip_addr_t convert_ip6_addr (int *ipv6_data);
```

Description

The **convert_ip6_addr** function converts the IPv6 address in the **in6_addr** structure defined in the **/usr/include/netinet/in.h** file to ProbeVue IP address data-type of type **ip_addr_t**. This function returns the converted **ip_addr_t** value.

Parameters

ipv6_data

Specifies the ipv6 address data that needs to be converted into the **ip_addr_t** format.

The following script is an example that prints the information about to whom the probed process is sending the data.

```
/* Declare the Function prototype */
int sendto(int s, char * uap_buf, int len, int flags, char * uap_to, int tolen);

typedef unsigned int in_addr_t;

/* Structure Declarations */

/* Declare the in_addr structure */
struct in_addr {
    in_addr_t      s_addr;
};

/* Declare the sockaddr_in structure */
struct sockaddr_in {
    unsigned char   sin_len;
    unsigned char   sin_family;
    unsigned short  sin_port;
    struct in_addr  sin_addr;
    unsigned char   sin_zero[8];
};

/* Declare the in6_addr structure */
struct in6_addr {
    union {
        int s6_addr32[4];
        unsigned short s6_addr16[8];
        unsigned char s6_addr8[16];
    } s6_addr;
};

/* Declare the sockaddr_in6 structure */
struct sockaddr_in6 {
    unsigned char   sin6_len;
    unsigned char   sin6_family;
    unsigned short  sin6_port;
    unsigned int    sin6_flowinfo;
    struct in6_addr sin6_addr;
    unsigned int    sin6_scope_id; /* set of interfaces for a scope */
};

/* Print the information about to whom it is sending data */
@@syscall:*.sendto:entry
{
    struct sockaddr_in6 in6;
    struct sockaddr_in in4;
    ip_addr_t ip;

    /* Copy the arg5 data into sockaddr_storage variable */
    /* using copy_userdata( ) Vue function */
    copy_userdata(__arg5, in4);

    /*
     * Verify whether the destination address is IPv4 or IPv6 and based on that call the
     * corresponding IPv4 or IPV6 conversion routine.
     */
    if (in4.sin_family == AF_INET)
    {
        /* Copy the ipv4 data into sockaddr_in structure using copy_userdata routine
        */
        copy_userdata(__arg5, in4);

        /* Convert Ipv4 data into ip_addr_t format */
        ip = convert_ip4_addr(in4.sin_addr.s_addr);

        /* Print the destination address and hostname using %H and %I format
        specifier */
        printf("It is sending the data to node %H(%I)\n",ip,ip);
    }
}
```

```

    }
    else if(in4.sin_family == AF_INET6)
    {
        /* Copy the ipv6 data into sockaddr_in6 structure using copy_userdata routine
*/
        copy_userdata(__arg5, in6);

        /* Convert Ipv6 data into ip_addr_t format */
        ip = convert_ip6_addr(in6.sin6_addr.s6_addr.s6_addr32);

        /* Print the destination address and hostname using %H and %I format
specifier */
        printf("It is sending the data to node %H(%I)\n", ip,ip);
    }
}

```

count

Purpose

Returns the number of elements in a list.

Syntax

```
long long count ( List listvar );
```

Description

The **count** function returns the number of elements that have been appended to the list variable specified by the *listvar* parameter.

For more information about the list data type, see [Data types in Vue](#). The preceding [list](#) section has an example script that uses the **count** function.

Parameter

Parameters	Description
<i>listvar</i>	Specifies a variable of type list .

copy_kdata

Purpose

Copies data from kernel memory into a Vue script variable.

Syntax

```
void copy_kdata( <type> *kaddr, <type> svar );
```

Description

The **copy_kdata** function reads data from the kernel memory into a Vue script variable. The variable might be of any of C-89 types that are supported by Vue except for pointer types. The length that is copied is equal to the size of the variable. For example, 4 bytes are copied if the target Vue script variable is of the type `int`, 8 bytes are copied if it is of type `long long`, and 48 bytes are copied if it is an array of 12 integers or `int[12]`.

Data in kernel space must be copied before it can be used in expressions or passed as parameters to a Vue function.

If an exception occurs while the process is running this function, for example, when a bad kernel address is passed to the function, the ProbeVue session is aborted with an error message.

Parameter

<i>kaddr</i>	Specifies the address of the kernel space data.
---------------------	---

svar

Specifies the script variable in which kernel data is copied. The script variable type can be of the type kernel data.

copy_userdata**Purpose**

Copies data from user memory into a Vue script variable.

Syntax

```
void copy_userdata( <type> *uaddr, <type> svar );
```

Description

The **copy_userdata** function reads data from user memory into a Vue script variable. The variable might be of any of the C-89 types that are supported by Vue . The length that is copied is equal to the size of the variable type. For example, 4 bytes are copied if the target Vue script variable is of the type int, 8 bytes are copied if it is of type long long, and 48 bytes are copied if it is an array of 12 integers or int[12].

Data in user space must be copied in before it can be used in expressions or passed as parameters to a Vue function.

If an exception occurs while the process is running this function, as for example when a bad user address is passed to the function, the ProbeVue session is aborted with an error message.

Parameter**uaddr**

Specifies the address of the user space data.

svar

Specifies the script variable in which user data is copied. The script variable type must be of the type user data.

diff_time**Purpose**

Returns the difference between two timestamps.

Syntax

```
unsigned long long diff_time( probev_timestamp_t ts1, probev_timestamp_t ts2, int format );
```

Description

The **diff_time** function returns the time difference between two timestamps that were recorded using the **timestamp** function. This function can return time difference in microseconds or milliseconds as specified by the *format* parameter.

The [get_location_point](#) and [list](#) sections have example scripts that use the **diff_time** function.

Parameter

Parameters	Description
<i>ts1</i>	Indicates the earlier timestamp.
<i>ts2</i>	Indicates the later timestamp.

Parameters	Description
<i>format</i>	Sets to one of the following values: MILLISECONDS Returns the time difference to the closest millisecond. MICROSECONDS Returns the time difference to the closest microsecond. You cannot pass a variable for this parameter.

eprintf

Purpose

Formats and prints data to the standard error.

Syntax

```
void eprintf ( String format[ , data, ... ] );
```

Description

The **eprintf** function is similar to **printf** function except that the output is sent to the standard error. The **eprintf** function converts, formats, and copies the data parameter values to the trace buffer under the control of the format parameter. As indicated by the syntax, a varying list of arguments can be passed as data parameters to the **eprintf** function. Vue supports all the conversion specifiers supported by the **printf** subroutine that is provided by the C library except for the %p conversion specifier.

The **eprintf** function cannot be used to print variables of list type. However, a variable of string type can be printed by using the %s conversion specifier. A variable of `probev_timestamp_t` type is printed in numerical form by using the %lld or %16llx specifier. A variable of type `probev_timestamp_t` is printed in the date format by using the %A or %W specifier.

Parameter

format

One string that contains plain characters that are directly copied to the trace buffer without any change. Another or more conversion specifiers that provide indication on how to format the data parameters.

data

Specifies zero or more arguments that correspond to the conversion specifiers in the format parameter.

Note: Tentative tracing is not allowed with the **eprintf** function.

exit

Purpose

Terminates the Vue script.

Syntax

```
void exit();
```

Description

The **exit** function terminates the Vue script. This disables all the probes enabled in the dynamic tracing session, discards any tentative trace data, issues the actions indicated in @@END probe and flushes all the captured trace data to the trace consumer. After the trace consumer prints any output traced in the @@END probe, the tracing session is terminated and the **probevue** process exits.

The same effect can be obtained by typing a Ctrl-C on the terminal where the **probevue** command was issued if it is running as the foreground task. Alternatively, you can directly send a **SIGINT** to the **probevue** process using the **kill** command or the **kill** system call.

The [list](#) section has an example script that uses the **exit** function. The [atoi](#) section has an example script for how to exit at the same time as the process that you are probing terminates.

Parameter

The **exit** function does not take any parameters, unlike the **exit** subroutine that is provided by the C library.

fd_fname

Purpose

Returns the file name for a particular file descriptor.

Syntax

```
char * fd_fname(int fd);
```

Description

This function gets the name of the file for a specific file descriptor. This returns the same value as that of the `__file->fname` (refer to `__file` built-in of I/O probe manager) for the same file.

Note: This function requires the `num_pagefaults` tunable value of the **probevctrl** command to be greater than 0. If it is 0 (or insufficient), then this function returns a null string as the file name.

Parameters

fd

File or socket descriptor value

fd_fstype

Purpose

Returns the file system type for a particular file descriptor.

Syntax

```
int fd_fstype(int fd);
```

Description

This function gets the type of the file system that the file of the specific file descriptor belongs to. It returns the same values as that of `__file->fs_type` (refer to `__file` built-in of I/O probe manager).

Note: This function requires the `num_pagefaults` tunable of **probevctrl** command to be greater than 0. If it is 0 (or insufficient), then this function returns -1 as the type of the file system.

Parameters

fd

File descriptor value

fd_ftype

Purpose

Returns the file type for a specific file descriptor.

Syntax

```
int fd_fstype(int fd);
```

Description

This function gets the file type for a specific file descriptor. It returns the same values as that of `__file->f_type` (refer to `__file` built-in of I/O probe manager).

Note: This function requires the `num_pagefaults` tunable of `probevctrl` command to be greater than 0. If it is 0 (or insufficient), then this function returns -1 as the file type.

Parameters

fd
File descriptor value

fd_inodeid

Purpose

Returns the inode ID for a specific file descriptor.

Syntax

```
unsigned long long fd_inodeid(int fd);
```

Description

This function returns the inode ID for the file that is related to a specific file descriptor. The inode ID is a system-wide unique `unsigned long long` value (it is different from the file system inode number and can change value if the system is rebooted). This value matches with the value that is returned by the `fpath_inodeid()` function for the same file.

Note: This function requires the `num_pagefaults` tunable of `probevctrl` command to be greater than 0. If it is 0 (or insufficient), then this function returns 0 as the inode ID.

Parameters

fd
File descriptor value

fd_mpath

Purpose

Get the mount path of the file system for a specific file descriptor.

Syntax

```
char * fd_mpath(int fd);
```

Description

This function gets mount path of the file system that the file of a specific file descriptor belongs to. It returns the same value as that of `__file->mount_path` (refer to `__file` built-in of I/O probe manager) for the same file.

Note: This function requires the `num_pagefaults` tunable of `probevctrl` command to be greater than 0. If it is 0 (or not sufficient), then this function returns a null string as the mount path.

Parameters

fd
File descriptor value

fd_path

Purpose

Returns the absolute path of the file for a specific file descriptor.

Syntax

```
path_t fd_path(int fd);
```

Description

This function returns the absolute path of the file for a specific file descriptor. The return value is of type `path_t`. It returns the same value as that of `__file->path` (refer to `__file` built-in of I/O probe manager) for the same file.

Note: This function requires the `num_pagefaults` tunable of `probevctrl` command to be greater than 0. If it is 0 (or insufficient), then this function returns null path, which when printed with the `printf(%p)` function prints a null string.

Parameters

fd

File descriptor value

fpath_inodeid

Purpose

Returns the inode ID for a specific file path.

Syntax

```
unsigned long long fpath_inodeid(String file_path);
```

Description

This function returns the inode ID for a specific file path. The inode ID is a system-wide unique unsigned long long value (it is different from the file system inode number and can change value if the system is restarted). If the file path does not exist, then the Vue script is rejected by the **probevue** command. The inode ID value remains the same as provided by the `__file->inode_id` in `vfs` probe events for the same file (refer to `__file` built-in of I/O probe manager).

Note: This function can be used anywhere in a Vue script (wherever Vue functions are allowed).

Parameters

file_path

A double-quoted literal string that represents an existing file. For example, `"/usr/lib/boot/unix_64"`. It cannot be a variable.

get_function

Purpose

Returns the name of the function that encloses the current probe. When the **get_function** function is called from `interval`, `systrace`, `BEGIN`, and `END` clause, the function returns an empty string.

Syntax

```
String get_function ( );
```

Description

The **get_function** function returns the name of the probed function, that is the function that encloses the current probe point. Generally, the name of the probed function is the probe point tuple field that precedes the location field.

The preceding `get_probe` section has an example script that uses the **get_function** function.

The **get_function** function returns an empty string when called from the interval probe manager.

Parameters

The **get_function** function does not take any parameters.

get_kstring

Purpose

Copies data from kernel memory into a String variable.

Syntax

```
String get_kstring( char *kaddr, int len);
```

Description

The **get_kstring** function reads data in kernel memory into a variable of type String.

Strings in kernel space must be copied in before they can be used in expressions or passed as parameters to a Vue function.

The target of the **get_kstring** function must always be a variable of type String. If a value of *-1* is specified for the **len** parameter, data is copied from kernel memory until a NULL byte is read (NULL bytes are used to stop text strings in the C language). If the length of the string is larger than the declared size of the target String variable, only the string characters up to the size of the variable is copied to it. However, the string in its entirety that is until a NULL byte is read it is initially copied into the temporary string buffer area. Users of the function must be careful that the kernel address points to a NULL-terminated string to avoid temporary string buffer area overflows, which cause the ProbeVue session to be stopped.

The maximum length of the string to be read in from kernel memory can be fixed by specifying a non-negative value for the **len** parameter. In this case, the copy proceeds until either a NULL byte is read or the specified number of bytes are read. This feature allows long strings in kernel memory to be copied more safely as the copy is limited by the value of the **len** parameter and does not cause ProbeVue's internal temporary string buffer to overflow.

If an exception occurs when this function is running, as for example when a bad kernel address is passed to the function, the ProbeVue session is aborted with an error message.

Parameters

addr

Specifies the address of the kernel space data.

len

Specifies the number of bytes of the kernel data to be copied. A value of *-1* indicates that the kernel data be treated as a "C" string and the copy to continue until a null byte (the '\0' character) is read. Be careful when you specify a *-1* as the value of the **len** parameter.

get_location_point

Purpose

Returns the current probe location point.

Syntax

```
int get_location_point ( );
```

Description

The **get_location_point** function returns the current probe location point as an offset from the enclosing function entry point. Specifically, it will return **FUNCTION_ENTRY** or zero if the probe point is at the entry point of the function and **FUNCTION_EXIT** or *-1* if it is at any exit point; else, it returns the actual address offset.

The following script shows an example of using the **get_location_point** function:

```
@@syscall:$1:read:entry, @@syscall:$1:read:exit
{
    probev_timestamp_t ts1, ts2;
    int diff;
```

```

    if (get_location_point() == FUNCTION_ENTRY)
        ts1 = timestamp();
    else if (get_location_point() == FUNCTION_EXIT) {
        ts2 = timestamp();
        diff = diff_time(ts1, ts2, MICROSECONDS);
        printf("Time for read system call = %d\n", diff);
    }
}

```

This function is not supported when called from the interval probe manager.

Parameters

The **get_location_point** function does not take any parameters.

get_probe

Purpose

Returns the current probe point specification.

Syntax

```
String get_probe ( );
```

Description

The **get_probe** function returns the internal representation of the current probe point specification. When saved internally, the probe point does not include the initial @@ prefix or the process ID, if any.

The following script shows an example of using the **get_probe** function:

```

#cat get_probe.e
@@uft:312678:*:run_prog:entry
{
    printf("function '%s' probe '%s'\n", get_function(), get_probe());
}

#probevue get_probe.e
function 'run_prog' probe 'uft:*:run_prog:entry'

```

Parameters

The **get_probe** function does not take any parameters.

get_stktrace

Purpose

Returns current stack trace.

Syntax

```
stktrace_t get_stktrace(int level);
```

Description

The **get_stktrace** function returns stack trace of the current thread. This stack trace is either stored in a **stktrace_t** type variable or printed when you use the **%t** or **%T** specifier in the ProbeVue built-in **printf** function. The level parameter indicates the number of levels up to which the stack trace is to be printed. The behavior of the **get_stktrace** function that is used inside the **printf** function is similar to the **stktrace** built-in function. The only difference, is that the symbol with address is printed when the **%t** specifier is used for the running thread, otherwise raw address is printed. Also, it prints the whole CPU stack by traversing all the machine states.

The following script shows an example of using the **get_stktrace** function:

```

t1 = get_stktrace(3)           // Returns the current stack trace & stores in
stktrace_t                   // type variable t1.
printf('%t\n', get_stktrace(4)); // Prints the current stack trace up to level

```

```

4.                                     // Prints symbol with addresses.

printf("%T\n", get_stktrace(4)); // Prints the current stack trace up to level 4
                                // Prints raw addresses.

```

Parameters

Parameters	Description
<i>level</i>	<p>Indicates the number of levels up to which the stack trace is to be saved in the stktrace_t type variable. A value of -1 indicates that the stack back chain is to be traversed to the extent possible. The default value of 0 tracks back to 2 levels and saves 2 entries. Any other positive integer value decides the number of levels to be saved in variable. The max value of level can be 240.</p> <p>Note: If entries from multiple msts are printed, the mst boundary is separated by a line that consists of '-' character. This line is also considered 1 level. It means that the entries that are printed are the level parameter that is provided minus the number of separator lines (unless the level parameter is -1).</p>

get_ubytes

Purpose

Copies at most **len** number of bytes of data, from user memory into a Vue script variable.

Syntax

```
int get_ubytes(void *src, char *dest, int len);
```

Description

The **get_ubytes** function reads **len** number of bytes of data from user memory into a Vue script variable. The variable can be of character pointer type. The length of the data that is copied is equal to the **len** parameter passed to this function or size of the destination Vue variable. If the **len** value is larger than the size of the destination Vue variable, length of the data that is copied is equal to the size of the destination variable. Data in user space must be copied to a Vue variable before the data can be used in expressions or passed as parameters to a Vue function.

Return value

The **get_ubytes** function returns the total number of bytes that were successfully copied and the function returns -1 on failure. If an exception occurs when the function is running, for example, if a bad user address is passed to the function, the ProbeVue session is ended with an error message.

get_kbytes

Purpose

Copies at most **len** number of bytes of data from kernel memory into a Vue script variable.

Syntax

```
int get_kbytes(void *src, char *dest, int len);
```

Description

The **get_kdata** function reads **len** number of bytes of data from kernel memory into a Vue script variable. The variable can be of character pointer type. The length of that data that is copied is equal to the **len** parameter passed to this function or size of the destination Vue variable. If the **len** value is larger than the size of the destination Vue variable, length of the data that is copied is equal to the size of the destination variable. Data in user space must be copied to a Vue variable before the data can be used in expressions or passed as parameters to a Vue function.

Return values

The `get_kbytes` function returns the total number of bytes that were successfully copied and the function returns -1 on failure. If an exception occurs when the function is running, for example, if a bad kernel address is passed to the function, the ProbeVue session is ended with an error message.

get_userstring

Purpose

Copies data from user memory.

Syntax

```
String get_userstring( char * addr, int len );
```

Description

The `get_userstring` function reads data in user memory into a variable of type **String**.

Data in user space must be copied in before you can use it in expressions or pass it as parameters to a Vue function. The target of the `get_userstring` function is generally a variable of type **String**. If a value of -1 is specified for the `len` parameter, data will be copied from user memory until a NULL byte is read (NULL bytes are used to terminate text strings in the C language). If the length of the string is larger than the declared size of the target **String** variable, only the string characters up to the size of the variable will be copied to it. However, the string in its entirety, that is until a NULL byte is read will initially need to be copied into the temporary string buffer area. Users of the function need to be careful that the user address points to a NULL-terminated string to avoid temporary string buffer area overflows which can cause the ProbeVue session to be aborted.

The actual length of the string to be read in from user memory can be fixed by specifying a value for the `len` parameter. In this case, the copy continues till either a NULL byte is read or the specified number of bytes are read. This feature allows non-string type data to be copied into a **String** variable, which can be printed out later by using the trace function.

Note: Vue does not treat a NULL byte as a string terminator, so real strings must not generally be copied in with this mechanism.

This function is only allowed in user space probes (like the `uft` probe type) or probes provided by the `syscall` probe manager. If a page fault occurs while copying the data, the copy operation is terminated and the String variable will contain only the data that was successfully copied. If an exception occurs while issuing this function, as for example when a bad user address is passed to the function, the ProbeVue session is aborted with an error message.

“ProbeVue dynamic tracing facility” on page 184 has an example script that uses the `get_userstring` function.

Note: You can modify the data type of the target of the copy operation using casts, although this will generate a warning message from the compiler. Thus, you can use the `get_userstring` function to copy not only strings, but also data laid out as structures and in other data formats from user space into ProbeVue space. The following script is an example of how to do this kind of data manipulation:

```
/* File: string2int.e
 *
 * Reads an integer passed in a pointer using get_userstring()
 *
 */
int get_file_sizep(int *fd);

@@BEGIN
{
    int i;
}

@@uft:$1:*.get_file_sizep:entry
{
    i = *(int *) (get_userstring(__arg1, 4));
```



```
    printf("fd = %d\n", i);  
}
```

Note: The target of the copy operation must be a String variable whose declared length is large enough to accept the copied in data or the ProbeVue session can be aborted. The **get_userstring** will accept any value for the size of the data to be copied in, but the maximum length that can be copied is limited by the memory limits of the ProbeVue session.

Parameters

Parameters	Description
<i>addr</i>	Specifies the address of the user space data.
<i>len</i>	Specifies the number of bytes of the user data to be copied. A value of -1 indicates that the user data must be considered as a string of characters and the copy operation must continue until a null byte is read (until the '\0' character is read). Use caution while specifying -1 as the value of the <i>len</i> parameter.

group_output_start

Purpose

Group the output of the clause to which this function belongs.

Syntax

```
void group_output_start();
```

Description

The `group_output_start()` function groups the output of a clause. Output messages that are generated by the VUE statements of a clause after the `group_output_start()` function call and before the `group_output_end()` function call are not interleaved by messages that are generated from clauses that might be running simultaneously in other CPUs.

group_output_end

Purpose

End grouping the output of the clause to which this function belongs.

Syntax

```
void group_output_end();
```

Description

Stops grouping of the output of the clause. Any subsequent VUE statements, if any, will not have a grouped output and can be interleaved by output messages from clauses that are running simultaneously in other CPUs. Call to the `group_output_end()` function must be preceded by a `group_output_start()` call. After the `group_output_start()` function is called, calling the `group_output_end()` function is optional. If the `group_output_end()` function is not called then the output is grouped till the end of clause.

list

Purpose

Creates and returns an empty list.

Syntax

```
List list ( );
```

Description

The **list** function is the constructor function for the list data type. It returns an empty list and auto-declares the target to be a list data type. There is no explicit way to declare a variable to be a list data type. A list variable is always created as a variable of global class.

The **list** function can be invoked from any clause. If you specify an existing list name when invoking the **list** function, the existing list is cleared.

A list variable can be used to collect values that are of integral type. Any value stored in the list is automatically promoted to a **long long** (or 64-bit integer) data type.

The following script shows an example of using the **list** function. It assumes that the sprobevue shell program, which encloses double quotes around each argument, invokes the Vue script.

```
/* File: list.e
 *
 * Collect execution time for read system call statistics
 *
 * Usage: sprobevue list.e <-s|-d>
 *
 * Pass -s for summary and -d for detailed information
 */

int read(int fd, void *buf, int n);

@@BEGIN
{
    String s[10];
    int detail;
    times = list();          /* declare and create a list */

    /* Check for parameters */
    s = $1;
    if (s == "-d")
        detail = 1;
    else if (s == "-s")
        detail = 0;
    else {
        printf("Usage: sprobevue list.e <-s|-d>\n");
        exit();
    }
}

@@syscall:*:read:entry
{
    /*
     * Save entry time in a thread-local to ensure that
     * in the read:exit probe point we can access our thread's value for
     * entry timestamp. If we use a global, then the variable can be
     * overlaid by the next thread to enter read and this can give
     * wrong values when we try to find the difference at read:exit
     * time since we use this later value instead of the original value.
     */

    __thread probev_timestamp_t t1;
    t1 = timestamp();
}

@@syscall:*:read:exit
{
    when (thread:t1 != 0)
    {
        __auto t2;
        __auto long diff;

        /* Get exit time */
        t2 = timestamp();
        diff = diff_time(t1, t2, MICROSECONDS);

        /* Append read time to list */
        append(times, diff);

        /* print detail data if "-d" was passed to script */
        if (detail)
            printf("%s (%ld) read time = %d micros\n", __pname, __pid, diff);
    }
}

@@interval:*:clock:10000
```

```

{
    /* Print statistics every 10 seconds */
    printf("Read calls so far = %d, total time = %d, max time = %d, " +
          "min = %d, avg = %d\n",
          count(times),
          sum(times),
          max(times),
          min(times),
          avg(times));
}

```

Parameters

The **list** function does not take any parameters.

lquantize

Purpose

Prints the keys and associated values of an associative array in graphical format by quantizing the values in logarithmic scale.

Syntax

```
void lquantize( aso-name ,int num-of-entries, int flags, sort_key_ind)
```

Description

This function displays the entries of an associative array in a graphical format based on the logarithmic value of the values of the associative array. If you want to print only elements having particular set of keys, the keys can be specified along with the associative array variable name in the first argument. To restrict only a certain dimension keys and allow any value for other keys in, the ANY keyword can be used. For example see the refer `print()` function section.

The first parameter is required and all other parameters are optional. If you do not specify optional parameters, then default print options are used.

Parameters

aso-name

The name of the associative array variable that you want to print. You can also specify keys for all dimensions in brackets. You can use the ANY keyword to match all keys in a key dimension.

num-of-entries

Specifies how many entries to print. This parameter is optional. Specify 0 to display all the entries. If no value is specified, the default print option for the session is used. Any negative value is equivalent to 0.

flags

Specifies the sort-type, sort-by and *list-value* flags. This parameter is optional. Flags sort-type, sort-by and list-value are described under the Associative Array Type section. If you specify 0, the default print option for the session is used.

sort_key_ind

The index of the key (key dimension) using which the output is sorted. If you specify -1, the first key is used for sorting. If the first key is not a sortable type, the output is not sorted.

max

Purpose

Returns the maximum of all the elements in a list.

Syntax

```
long long max ( List listvar );
```

Description

The **max** function returns the maximum of all the elements that have been appended to the list variable specified by the *listvar* parameter.

For more information about the list data type, refer the [“ProbeVue dynamic tracing facility” on page 184](#) topic. The preceding *listvar* section has an example script that uses the **max** function.

Parameters

listvar: Specifies a variable of type **list**.

min

Purpose

Returns the minimum of all the elements in a list.

Syntax

```
long long min ( List listvar );
```

Description

The **min** function returns the minimum of all the elements that have been appended to the list variable specified by the *listvar* parameter.

The preceding *listvar* section has an example script that uses the **min** function.

Parameters

listvar: Specifies a variable of type **list**.

print_args

Purpose

Prints current function and its argument values.

Syntax

```
void print_args();
```

Description

The **print_args** function prints the function name followed by comma separated function arguments enclosed in round brackets. The argument values are printed based on the argument type information available in the traceback table of the function. This routine is allowed in entry probes of uft/uftxc++ and syscall/syscallx probes. This is useful in probes where the probe location is specified as a regular expression.

Parameters

The **print_args** function does not take any parameters.

Note: The **print_args** routine do not generate any output, if the traceback table of the routine has been paged out, and no free page fault context is available. The number of page faults to be handled can be increased using the `probevctl` command, and the script can be retried.

print

Purpose

Prints the keys and associated values of an associative array.

Syntax

```
void print ( aso-name , int num-of-entires , int flags , int sort_key_ind );
```

Description

This function prints the elements of the associative array specified by the variable `aso-name`. To print only elements that have a particular set of keys, you can specify the keys with the associative array variable name in the first argument. To restrict only a certain dimension keys and allow any value for all other keys, use the ANY keyword.

Example

```
print(aso_var[0][ANY][ANY]); // print all elements having first key as 0 (other keys can
be anything)
print(aso_var[ANY][ANY][ANY]); // print all; equivalent to print(aso_var);
```

The first parameter is required and all other parameters are optional. If you do not specify optional parameters, the default print options are used.

Note: Tentative tracing is not allowed with the print function.

For associative arrays with multi-dimensional keys, the keys are printed as a '|' separated list and the value is printed on the same line. If a key produces multi-line output, key is printed in a separate line and the value is printed in a new line. The following example displays a script that has an associative array with 3 key dimensions types of int and a value type of integer:

```
aso1[0]["a"][2.5] = 100;
aso1[1]["b"][3.5] = 101;
print(aso1);
The output from previous print() function follows:
[key1 | key2 | key3 | value
0 | a | 2.5000000 | 100
1 | b | 3.5000000 | 101
```

The following example uses an associative array with two key dimensions of type `int` and `stktrace_t` has a value type of string:

```
aso2[0][get_stktrace(-1)] = "abc";
print(aso2);
The output from print() function above will be similar to:
[key1 | key2 | value
0 | |
| 0x100001b8
| 0x10003328
| 0x1000166c
| 0x10000c30
| .read+0x288
| sc_entry_etrc_point+0x4
| .kread+0x0
| abc
```

Parameters

aso-name

The name of the associative array variable that you want to print. You can also specify keys for all dimensions in brackets. You can use the ANY keyword to match all keys in a key dimension.

num-of-entries

Specifies how many entries to print. This parameter is optional. Specify 0 to display all the entries. If no value is specified, the default print option for the session is used. Any negative value is equivalent to 0.

flags

Specifies the sort-type, sort-by and the *list-value* flags. This parameter is optional. Flags sort-type, sort-by and list-value are described under the Associative Array Type section. If you specify 0, the default print option for the session is used.

sort_key_ind

The key index (key dimension) is used for sorting the output. If you specify -1, the first key is used for sorting. If the first key is not a sortable type, then sorting of output is not sorted.

printf

Purpose

Formats and copies data to the trace buffer.

Syntax

```
void printf ( String format[ , data, ... ] );
```

Description

The **printf** function converts, formats, and copies the *data* parameter values to the trace buffer under the control of the format parameter. As indicated by the syntax, a varying list of arguments can be passed as *data* parameters to the **printf** function. Vue supports all the conversion specifiers supported by the **printf** subroutine provided by the C library except for the *%p* conversion specifier.

Apart from the C library's `printf()` specifiers, Vue language supports two more specifiers: *%A* and *%W*.

%A

Prints the `probev_timestamp_t` 't' in the default date format. This format can be changed by using the `set_date_format()` function.

%W

prints the `probev_timestamp_t` 't', in seconds and microseconds, relative to the start of the probevue session.

%p

prints the string corresponding to the absolute file path of the specified `path_t` value.

%M

prints the MAC address of the specified `mac_addr_t` value.

%I

prints the IP address in dotted decimal format for ipv4 address and dotted hex format for IPV6 address of the specified `ip_addr_t` value

%H

prints the host name in string or dotted decimal or hex format of the specified `ip_addr_t` value.

Note: If IP address is resolved by Domain Name System (DNS), the `printf` function displays corresponding host name. Otherwise, it displays the IP address in dotted decimal or hex format.

The **printf** function cannot be used to print variables of type **list**. However, a variable of type **string** can be printed by using the *%s* conversion specifier. A variable of type `probev_timestamp_t` type is printed in numerical form by using the *%lld* or *%16llx* specifier. `probev_timestamp_t` is printed in the date format by using the *%A* or *%W* specifier.

The following script demonstrates some examples of using the **printf** function:

```
@@BEGIN
{
    String s[128];
    int i;
    float f;
    f = 2.3;

    s = "Test: %d, float = %f\n";
    i = 44;

    printf(s, i, f);

    s = "Note:";
    printf("%s Value of i (left justified) = %-12d and right-justified = %12d\n",
        s, i, i);
}
```

```

printf("With a long format string that may span multiple lines, you " +
      "can use the '+' operator to concatenate the strings " +
      "in multiple lines into a single string: 0x%08x\n", i);

exit();
}

```

Parameters

format

A string that contains plain characters that are directly copied to the trace buffer without any change and one or more conversion specifiers that provide indication on how to format the *data* parameters.

data

Specifies zero or more arguments that correspond to the conversion specifiers in the *format* parameter.

ptree

Purpose

To print the process tree for the probed process.

Syntax

```
void ptree ( int depth );
```

Description

The **ptree** function prints the process tree for the probed process. This functions prints both the parent and child hierarchy. The depth passed as a parameter can help in controlling the depth of child processes that need to be printed. This function cannot be used in BEGIN or END or systrace probes. Also, this function can be used in the interval probes only if the PID is specified.

Note: This function is not immediately executed in kernel when it is called from a probe, but instead it is scheduled to run later in the user space. Hence, if the process tree changes in between, the output of the **ptree** function might not match the tree structure when the function was actually called.

Sample output

The sample output of the process tree follows as:

```

PID          CMD
1            init
             |
             V
3342460      srcmstr
             |
             V
3539052      inetd
             |
             V
7667750      telnetd
             |
             V
6881336      ksh
             |
             V
5112038      probevue
             |
             V
7930038      tree      <=====
6553782      | \--tree
4849828      | \--tree
6422756      | \--tree
3408074      | \--tree
5963846      | \--tree
7864392      | \--tree
7799006      | \--tree

```

Parameters

Parameters	Description
<i>depth</i>	Specifies the maximum depth that the ptree function traverses while printing the children information for the process. If -1 is passed, it prints all the children of the process.

quantize

Purpose

Prints the keys and associated values of an associative array in graphical format by quantizing the values in linear scale.

Syntax

```
void quantize ( aso-name, int num-of-entries, int flags, int sort_key_ind)
```

Description

This function displays the entries of an associative array in a graphical format based on the linear value of the values of the associative array. To print only elements that have a particular set of keys, the keys can be used with the variable name of the associative array in the first argument. To restrict only a certain dimension keys and allow all keys in the remaining dimensions, you can use the ANY keyword.

Apart from the first parameter, the others are optional parameters. If these optional parameters are not given, then default print options are used.

Parameters

aso-name

The name of the associative array variable that you want to print. You can also specify keys for all dimensions in brackets. You can use the ANY keyword to match all keys in a key dimension.

num-of-entries

Specifies how many entries to print. This parameter is optional. Specify 0 to display all the entries. If no value is specified, the default print option for the session is used. Any negative value is equivalent to 0.

flags

Specifies the sort-type, sort-by and *list-value* flags. This parameter is optional. Flags sort-type, sort-by and list-value are described under the Associative Array Type section. If you specify 0, the default print option for the session is used.

sort_key_ind

The index of the key (key dimension) using which the output is sorted. If you specify -1, the first key is used for sorting. If the first key is not a sortable type, then output is not sorted.

qrange

Purpose

This routine gets the slot number for ranges and adds the range data type as value type for an associative array.

Syntax

```
void qrange(aso[key], range_t range_data, int value);  
void qrange(aso[key], range_t range_data, String value);
```

Description

The `qrange` routine can find the slot number for both integral and String range types. If the range type is Integral type, then the third argument type should be an integer otherwise for the String range

data type, the third argument should be a String type. The `qrange` routine will find the slot number where the passed value falls. The count for that slot number will be increased for the range type stored in an associative array as a value.

Parameters

aso[key]

An associative array with specified key.

range_data

A **range_t** data type.

value

value can be either integer or can be of String type.

round_trip_time

Purpose

Returns the smoothed round-trip time for TCP connection for a specific socket descriptor.

Syntax

```
int round_trip_time(int sock_fd);
```

Description

The `round_trip_time` function gets the smoothed round-trip time (srtt) for a specific socket descriptor. It provides the valid round-trip value for stream socket descriptor and returns -1 as a smooth round-trip time value for invalid or non-stream sock descriptor. This function is available only in `uft` and `syscall` probe managers.

Note: This function requires the `num_pagefaults` tunable value of `probevctrl` command to be greater than 0. If it is 0, then this function returns -1 as the round-trip time.

Parameters

fd

File or socket descriptor value.

set_aso_print_options

Purpose

Changes the default print options for associative arrays.

Syntax

```
void set_aso_print_options( int num-of-entries, int flags);
```

Description

The `set_aso_print_options()` function changes the default print options for associative arrays. The print options that can be provided by the user and their initial values are listed under the Associative Array type section. This function is allowed only in BEGIN probe.

Parameters

num-of-entries

Specifies to print the first 'n' key or value pairs. If it is 0, all the entries are displayed.

flags

Specifies the **sort-type**, **sort-by**, **list-value** and **stack-raw** flags. These flags are described under the Associative Array Type section. These parameters are optional.

set_range

Purpose

Initializes the linear and power range type data.

Syntax

```
void set_range(range_t range_data, LINEAR, int min, int max, int step);  
void set_range(range_t range_data, POWER, 2);
```

Description

There are two different variants of *set_range*. To initialize a range data as a linear range the flag `LINEAR` with *min*, *max* and *step* will be passed as arguments. For initialize a Power range the flag `POWER` with two will be passed as arguments. This routine will initialize the range type either as Linear or as Power based on arguments passed. The linear range type data will be initialized with the *min*, *max* and *step* value passed, while power range type data will be initialized the power value as 2.

Parameters

Parameters (for Linear range type):

range_data

A **range_t** data type.

LINEAR

An integer consistent flag indicating that distribution of **range_data** is linear.

min

Indicates the lower bound of **range_data**.

max

Indicates the upper bound of **range_data**.

step

Indicates the size of the specified range of the value for each row of **range_data**. The type of min, max & step can be only integral (int, short, long, long, long). No other types will be allowed.

Parameters (for Power range type):

range_data

A **range_t** data type.

POWER

A integer constant flag indicating that the distribution of values is a POWER distribution.

A constant indicating the value of power. As of now only power of two is supported.

set_date_format

Purpose

Updates the date format that is used for printing the `probev_timestamp_t` data type.

Syntax

```
void set_date_format(String s);
```

Description

Updates the date format.

This function supports all the conversion specifiers that are supported by the C library's `strftime()` for the date format. Any specifier, which is not supported by `strftime()`, is invalid and default format is used.

Default format

MM:DD:YYYY hh:mm:ss TZ

MM

Month of the year as a decimal number (01 to 12).

DD

Day of the month as a decimal number (01 to 31).

YYYY

Year as a decimal number (for example, 1989).

hh

24-hour clock hour as a decimal number (00 to 23).

mm

Minutes of the hour as a decimal number (00 to 59).

ss

Seconds of the minute as a decimal number (00 to 59).

TZ

Time-zone name if one can be determined (for example, CDT).

Note: The `set_date_format()` function is called only in the @@ BEGIN probe. Constant string must be passed as *format*.

Parameters

S - string that holds the date format.

sockfd_netinfo**Purpose**

Gets the local and remote ports and IP addresses information for a specific socket descriptor.

Syntax

```
void sockfd_netinfo(int sock_fd, net_info_t ninfo);
```

Description

The `sockfd_netinfo` function gets the local IP address, remote IP address, local port number, and remote port number information for input socket descriptor. This function gets the valid local and remote port numbers and IP addresses information for valid socket descriptor. It gets 0 for invalid descriptor or if the descriptor is a not of socket type.

Note: This function requires the **num_pagefaults** tunable value of `probevctl` command to be greater than 0 and preferably 2 or more. If it is 0, then this function gets the invalid (0) as local and remote ports and IP addresses information.

Parameters**fd**

File or socket descriptor value.

ninfo

Specifies the script variable `net_info_t` in which network four tuples (local and remote IP addresses and port numbers) information for a specific file descriptor is copied.

start_tentative, end_tentative**Purpose**

Indicates the start and the end of a tentative tracing section.

Syntax

```
void start_tentative( String bufID );  
void end_tentative( String bufID );
```

Description

These functions indicate the start and end of a tentative tracing section within a Vue clause. Trace data generated by trace output functions enclosed within the tentative tracing section is saved on a tentative basis until a **commit_tentative** or a **discard_tentative** function is called to commit or discard this data. The **end_tentative** function is optional and if it is not specified, the end of the Vue clause is implicitly assumed to indicate the end of the tentative tracing section.

The generated tentative trace data is identified by the *bufID* parameter, which needs to be a string constant or literal and not a variable. Tentative trace data can be collected under different IDs simultaneously which can then each be committed or discarded as a separate block. ProbeVue supports up to 16 tentative trace buffers in the same dynamic tracing session, so up to 16 different trace IDs can be used in a Vue script. A single Vue clause can contain more than one tentative tracing section with different IDs.

Parameters

bufID

Specifies a string constant that indicates the tentative trace buffer ID.

stktrace

Purpose

Generates and prints a runtime stack trace.

Syntax

```
void stktrace ( int flags, int levels );
```

Description

The **stktrace** function prints the stack trace at the current probe point. By default, the stack trace is generated in compact form with only call chain addresses for only up to two levels. You can use the *flags* and the *levels* parameters to modify the format and contents of the stack trace. ProbeVue cannot read paged out data, so the stack trace is truncated if a page fault is encountered when accessing the stack.

The **stktrace** function does not return any values.

Parameters

Parameters

Description

flags

Either sets to 0 to specify default behavior or specifies one or more of the following flags:

PRINT_SYMBOLS

Prints symbol names instead of addresses.

GET_USER_TRACE

By default, the stack trace is stopped at the system call boundary if the probe location is in kernel space. This flag indicates to trace all the way into user space also up to the number of levels specified by the *levels* parameter.

GET_ALL_MSTS

By default, the stack trace is collected for only one context (machine state) where the probe was started. If this flag is specified, then stack trace is printed for all the chained contexts for that CPU.

If you want to pass multiple flags, the different flags must be 'or'ed using the OR operator, that is the '|' operator. You cannot pass a variable for this parameter.

levels

Indicates the number of levels up to which the stack trace is to be printed. A value of -1 indicates that the stack back chain is to be traversed to the extent possible. Default value of 0 tracks back to 2 levels.

Note: If entries from multiple msts are printed, the mst boundary is separated by a line that consists of '-' character. This line is also considered 1 level. It means that the entries that are printed are the level parameter that is provided minus the number of separator lines (unless the level parameter is -1).

strstr

Purpose

Return a string inside another string.

Syntax

```
String strstr( String s1, String s2 );
```

Description

The **strstr** function finds the first occurrence of the string specified by the *s2* parameter in the string specified by the *s1* parameter and returns a new string that contains the characters in the *s1* parameter starting from this location. Neither the *s1* parameter nor the *s2* parameter are modified by this operation. If the sequence of characters specified by the *s2* parameter does not appear even once in the *s1* parameter, then this function returns an empty string.

Note: This function's behavior is not the same as the **strstr** subroutine in the C library.

Parameters

Parameters	Description
<i>s1</i>	Specifies the string within which to search.
<i>s2</i>	Specifies the string to be searched for

sum

Purpose

Returns the sum of all the elements in a list.

Syntax

```
long long sum ( List listvar );
```

Description

The **sum** function returns the sum of all the elements that have been appended to the list variable specified by the *listvar* parameter.

Parameters

listvar	Specifies a variable of type list .
----------------	--

timestamp

Purpose

Returns the current timestamp.

Syntax

```
probev_timestamp_t timestamp( );
```

Description

The **timestamp** function returns the current timestamp in the **probev_timestamp_t** abstract data type. Although abstract, the value has the following properties:

- It returns an equal or close value when called simultaneously from different CPUs.
- If the **timestamp** function is invoked twice and the second call can be architecturally guaranteed to have happened later in time, the value returned for the second call is greater than or equal to the value returned by the first call (provided the system has not been rebooted in between the calls).

There is no relationship between the values returned by the **timestamp** function on two different systems. Although the compiler will let you treat the returned value as a 64-bit integer, doing so can introduce compatibility problems.

Note: The **lbolt** kernel variable, whose value indicates the number of ticks since boot-up, or the **time** kernel variable whose value indicates the number of seconds since epoch (Jan 1, 1970) can both be used instead of this function if a lower-resolution timestamp is acceptable.

```
typedef long long time_t;
__kernel time_t lbolt;      /* number of ticks since last boot */
__kernel time_t time;      /* memory mapped time in secs since epoch */
```

Parameters

The **timestamp** function does not take any parameters.

trace

Purpose

Copies raw data to the trace buffer in hexadecimal format.

Syntax

```
void trace ( data );
```

Description

The **trace** function takes a single parameter which must be a variable. The **trace** function does not accept expressions.

The **trace** function copies the value of the passed-in argument to the trace buffer. The argument can be of any data type and the size of the data copied to the trace buffer is based on its innate size. Thus, four bytes are copied for an integer argument, four or eight bytes for pointers (depending upon whether the execution is in 32-bit or 64-bit mode) and the size of the structure for arguments of type **struct**. For a variable of type **String**, the number of bytes copied over is the declared string length (which is not the same as the length of the string contained in the variable). A variable of type **probev_timestamp_t** is at least 8 bytes long.

The trace reporter displays the hexadecimal data written by the **trace** function in groups of four characters without any additional formatting.

Note: The **trace** function also accepts a variable of type **list** as a parameter, but the output is not useful in this case.

Parameter

data

Data argument to be copied to trace buffer.

Process attributes

You can obtain attributes of a process different than the running context in AIX 7.2 TL3. ProbeVue provides the following interfaces to access the attributes.

```
get_proc(<process identifier>, PROCESS_ATTRIBUTE_NAME, <variable>);
```

process identifier

The process identifier of the process for which you are collecting attribute information about.

PROCESS_ATTRIBUTE_NAME

The process attribute name of the process for which you are collecting attribute information about.

variable

The ProbeVue variable that holds the results of the operation.

```
get_thread(<thread identifier>, THREAD_ATTRIBUTE_NAME, <variable>);
```

thread identifier

The thread identifier of the thread whose attributes are requested.

THREAD_ATTRIBUTE_NAME

The thread attribute name of the thread for which you are collecting attributes about.

variable

The ProbeVue variable that holds the results of the operation.

```
get_ublock(<process/thread identifier>, UBLOCK_ATTRIBUTE_NAME,<variable>);
```

process/thread identifier

The process/thread identifier.

UBLOCK_ATTRIBUTE_NAME

The ublock attribute name.

The ProbeVue variable which holds the results of the operation.

```
get_cwd(<process> identifier , <variable>);
```

process identifier

The process identifier.

variable

The ProbeVue variable that holds the results of the operation.

The ProbeVue variable must be of type *string* with length greater than or equal to 1024.

Note: All ublock attributes are of type *string* except for the CWD attribute. All other attributes are of type *long long*.

Return values

These interfaces return 0 when successful and the following error codes during error scenarios.

SL No.	Error	Description
1	ERR_INTR_UBLOCK_ACCESS	Ublock cannot be accessed in the current environment.
2	ERR_INTR_CWD_ACCESS	CWD cannot be accessed in this environment.
3	ERR_INVALID_PROC	Invalid process identifier provided
4	ERR_INVALID_THREAD	Invalid thread identifier provided
5	ERR_PROC_INVALID_STATE	Process is in an invalid state
6	ERR_THREAD_INVALID_STATE	Thread is in an invalid state
7	ERR_VMATTACH	Error encountered while attaching a VM segment
8	ERR_INVALID_STORAGE	Invalid storage provided
9	ERR_INVALID_ATTR	Invalid attribute requested
10	ERR_SHORT_STRING	String provided is short
11	ERR_GET_PATH	Error encountered while accessing <i>PATH</i>
12	ERR_INSUFF_SORAGE	Internal EVM storage is insufficient

Notes:

- get_ublock and get_cwd are not allowed in the Java interval probe manager or the system trace probe manager.
- get_proc and get_thread are not allowed in the Java Probe Manager (JPM).

Process attributes

You can use the following process attributes to access the interfaces.

Table 75.

process_attribute_name	Description
PID	Process identifier
PPID	Parent process identifier
PGID	Process group identifier
PUID	Real user identifier
SUID	Saved user identifier
PRI	Process priority
NICE	Process Nice Value
CPU	Processor Usage
ADSPACE	Process Address space
MAJFLT	I/O Page fault
MINFLT	Non-I/O page fault
SIZE	Size of image in pages
SIGPEND	Signals pending in the process
SIGIGNORE	Signals ignored by the process
SIGCATCH	Signals being caught by the process
FORKTIME	Creation time of the process

Thread attributes

You can use the following thread attributes to access the interfaces.

Table 76.

thread_attribute_name	Description
TPID	Process identifier
POLICY	Thread scheduling policy
TPRI	Thread priority
CPUUSAGE	CPU usage by thread
CPUID	CPU identifier where thread is running
SIGMASK	Signals blocked on the thread
LOCKCOUNT	Number of kernel locks taken by the thread
PTID	Pthread identifier of the thread
HOMESRAD	HOMESRAD of the thread
HOMECPU	HOMECPU of the thread

ublock attributes

You can use the following ublock attributes to access the interfaces.

Table 77.

ublock attribute name	Description
TEXT	Start of process text
TSIZE	Text size of the process
DATA	Start of the data for a process
SDATA	Current data size
MDATA	Maximum data size of the process
STACK	Start of STACK
STKMAX	Maximum size of stack
EUID	Effective user identifier
UUID	Real user identifier
EGID	Effective group identifier
GID	Group identifier
UTIME	Process user resource usage time in seconds
STIME	Process system resource usage time in seconds
MAXFD	MAX FD value in user
IS64U	Is a 64-bit process.

Multithreaded programming

This section provides guidelines for writing multithreaded programs using the threads library (**libpthreads.a**).

The AIX threads library is based on the X/Open Portability Guide Issue 5 standard. For this reason, the following information presents the threads library as the AIX implementation of the XPG5 standard.

Parallel programming uses the benefits of multiprocessor systems, while maintaining a full binary compatibility with existing uniprocessor systems. The parallel programming facilities are based on the concept of threads.

The advantages of using parallel programming instead of serial programming techniques are as follows:

- Parallel programming can improve the performance of a program.
- Some common software models are well-suited to parallel-programming techniques.

Traditionally, multiple single-threaded processes have been used to achieve parallelism, but some programs can benefit from a finer level of parallelism. Multithreaded processes offer parallelism within a process and share many of the concepts involved in programming multiple single-threaded processes.

The following information introduces threads and the associated programming facilities. It also discusses general topics concerning parallel programming:

Note: In this topic collection, the word *thread* used alone refers to *user threads*. This also applies to user-mode environment programming references, but not to topics related to kernel programming.

Understanding threads and processes

A *thread* is an independent flow of control that operates within the same address space as other independent flows of controls within a process.

Traditionally, thread and process characteristics are grouped into a single entity called a *process*. In other operating systems, threads are sometimes called *lightweight processes*, or the meaning of the word *thread* is sometimes slightly different.

The following sections discuss the differences between a thread and a process.

In traditional single-threaded process systems, a process has a set of properties. In multithreaded systems, these properties are divided between processes and threads.

Threads have some limitations and cannot be used for some special purposes that require multi-processed programs.

Related concepts

Contention scope and concurrency level

The *contention scope* of a user thread defines how it is mapped to a kernel thread

Process duplication and termination

Because all processes have at least one thread, creating (that is, duplicating) and terminating a process implies the creation and the termination of threads.

Process properties

A process in a multithreaded system is the changeable entity.

It must be considered as an execution frame. It has traditional process attributes, such as:

- Process ID, process group ID, user ID, and group ID
- Environment
- Working directory

A process also provides a common address space and common system resources, as follows:

- File descriptors
- Signal actions
- Shared libraries
- Inter-process communication tools (such as message queues, pipes, semaphores, or shared memory)

Thread properties

A thread is the schedulable entity.

It has only those properties that are required to ensure its independent control of flow. These include the following properties:

- Stack
- Scheduling properties (such as policy or priority)
- Set of pending and blocked signals
- Some thread-specific data

An example of thread-specific data is the **errno** error indicator. In multithreaded systems, **errno** is no longer a global variable, but usually a subroutine returning a thread-specific **errno** value. Some other systems may provide other implementations of **errno**.

Threads within a process must not be considered as a group of processes. All threads share the same address space. This means that two pointers having the same value in two threads refer to the same data. Also, if any thread changes one of the shared system resources, all threads within the process are affected. For example, if a thread closes a file, the file is closed for all threads.

Initial thread

When a process is created, one thread is automatically created. This thread is called the *initial thread*.

It ensures the compatibility between the old processes with a unique implicit thread and the new multithreaded processes. The initial thread has some special properties, not visible to the programmer, that ensure binary compatibility between the old single-threaded programs and the multithreaded operating system. It is also the initial thread that executes the **main** routine in multithreaded programs.

Modularity

Programs are often modeled as a number of distinct parts interacting with each other to produce a desired result or service.

A program can be implemented as a single, complex entity that performs multiple functions among the different parts of the program. A more simple solution consists of implementing several entities, each entity performing a part of the program and sharing resources with other entities.

By using multiple entities, a program can be separated according to its distinct activities, each having an associated entity. These entities do not have to know anything about the other parts of the program except when they exchange information. In these cases, they must synchronize with each other to ensure data integrity.

Threads are well-suited entities for modular programming. Threads provide simple data sharing (all threads within a process share the same address space) and powerful synchronization facilities, such as mutexes (mutual exclusion locks) and condition variables.

Software models

This section describes the different software models.

All these models lead to modular programs. Models may also be combined to efficiently solve complex tasks.

These models can apply to either traditional multi-process solutions, or to single process multi-thread solutions, on multithreaded systems. In the following descriptions, the word *entity* refers to either a single-threaded *process* or to a single *thread* in a multithreaded process.

The following common software models can easily be implemented with threads:

Controller/Worker Model

In the controller/worker (sometimes called boss/worker) model, a controller entity receives one or more requests, then creates worker entities to execute them. Typically, the controller controls the number of workers and what each worker does. A worker runs independently of other workers.

An example of this model is a print job spooler controlling a set of printers. The spooler's role is to ensure that the print requests received are handled in a timely fashion. When the spooler receives a request, the controller entity chooses a printer and causes a worker to print the job on the printer. Each worker prints one job at a time on a printer, while it also handles flow control and other printing details. The spooler may support job cancellation or other features that require the controller to cancel worker entities or reassign jobs.

Divide-and-Conquer Models

In the divide-and-conquer (sometimes called *simultaneous computation* or *work crew*) model, one or more entities perform the same tasks in parallel. There is no master entity; all entities run in parallel independently.

An example of a divide-and-conquer model is a parallelized **grep** command implementation, which could be done as follows. The **grep** command first establishes a pool of files to be scanned. It then creates a number of entities. Each entity takes a different file from the pool and searches for the pattern, sending

the results to a common output device. When an entity completes its file search, it obtains another file from the pool or stops if the pool is empty.

Producer/Consumer Models

The producer/consumer (sometimes called *pipelining*) model is typified by a production line. An item proceeds from raw components to a final item in a series of stages.

Usually a single worker at each stage modifies the item and passes it on to the next stage. In software terms, an AIX command pipe, such as the **cpio** command, is an example of this model.

For example, a reader entity reads raw data from standard input and passes it to the processor entity, which processes the data and passes it to the writer entity, which writes it to standard output. Parallel programming allows the activities to be performed concurrently: the writer entity may output some processed data while the reader entity gets more raw data.

Kernel Threads and User Threads

A kernel thread is the schedulable entity, which means that the system scheduler handles kernel threads.

These threads, known by the system scheduler, are strongly implementation-dependent. To facilitate the writing of portable programs, libraries provide *user* threads.

A *kernel thread* is a kernel entity, like processes and interrupt handlers; it is the entity handled by the system scheduler. A kernel thread runs within a process, but can be referenced by any other thread in the system. The programmer has no direct control over these threads, unless you are writing kernel extensions or device drivers. For more information about kernel programming, see *Kernel Extensions and Device Support Programming Concepts*.

A *user thread* is an entity used by programmers to handle multiple flows of controls within a program. The API for handling user threads is provided by the *threads library*. A user thread only exists within a process; a user thread in process *A* cannot reference a user thread in process *B*. The library uses a proprietary interface to handle kernel threads for executing user threads. The user threads API, unlike the kernel threads interface, is part of a POSIX-standards compliant portable-programming model. Thus, a multithreaded program developed on an AIX system can easily be ported to other systems.

On other systems, user threads are simply called *threads*, and *lightweight process* refers to kernel threads.

Thread models and virtual processors

User threads are mapped to kernel threads by the threads library. The way this mapping is done is called the *thread model*.

There are three possible thread models, corresponding to three different ways to map user threads to kernel threads.

- M:1 model
- 1:1 model
- M:N model

The mapping of user threads to kernel threads is done using *virtual processors*. A virtual processor (VP) is a library entity that is usually implicit. For a user thread, the VP behaves like a CPU. In the library, the VP is a kernel thread or a structure bound to a kernel thread.

In the M:1 model, all user threads are mapped to one kernel thread; all user threads run on one VP. The mapping is handled by a library scheduler. All user-threads programming facilities are completely handled by the library. This model can be used on any system, especially on traditional single-threaded systems.

In the 1:1 model, each user thread is mapped to one kernel thread; each user thread runs on one VP. Most of the user threads programming facilities are directly handled by the kernel threads. This model is the default model.

In the M:N model, all user threads are mapped to a pool of kernel threads; all user threads run on a pool of virtual processors. A user thread may be bound to a specific VP, as in the 1:1 model. All unbound user threads share the remaining VPs. This is the most efficient and most complex thread model; the user threads programming facilities are shared between the threads library and the kernel threads. This model can be set by setting the AIXTHREAD_SCOPE environment variable to **P**.

Threads Library API

This section provides general information about the threads library API.

Although the following information is not required for writing multithreaded programs, it can help the programmer understand the threads library API.

Object-Oriented Interface

The threads library API provides an object-oriented interface. The programmer manipulates opaque objects using pointers or other universal identifiers.

This ensures the portability of multithreaded programs between systems that implement the threads library and also allows implementation changes between two releases of AIX, necessitating only that programs be recompiled. Although some definitions of data types may be found in the threads library header file (**pthread.h**), programs should not rely on these implementation-dependent definitions to directly handle the contents of structures. The regular threads library subroutines must always be used to manipulate the objects.

The threads library essentially uses the following kinds of objects (opaque data types): threads, mutexes, rwlocks, and condition variables. These objects have attributes that specify the object properties. When creating an object, the attributes must be specified. In the threads library, these creation attributes are themselves objects, called *threads attributes objects*.

The following pairs of objects are manipulated by the threads library:

- Threads and thread-attributes objects
- Mutexes and mutex-attributes objects
- Condition variables and condition-attributes objects
- Read-write locks

An attributes object is created with attributes having default values. Attributes can then be individually modified by using subroutines. This ensures that a multithreaded program will not be affected by the introduction of new attributes or by changes in the implementation of an attribute. An attributes object can thus be used to create one or several objects, and then destroyed without affecting objects created with the attributes object.

Using an attributes object also allows the use of object classes. One attributes object may be defined for each object class. Creating an instance of an object class is done by creating the object using the class attributes object.

Naming Convention for the Threads Library

The identifiers used by the threads library follow a strict naming convention. All identifiers of the threads library begin with **pthread_**.

User programs should not use this prefix for private identifiers. This prefix is followed by a component name. The following components are defined in the threads library:

Component	Description
pthread_	Threads themselves and miscellaneous subroutines
pthread_attr	Thread attributes objects
pthread_cond	Condition variables

Component	Description
pthread_condattr	Condition attributes objects
pthread_key	Thread-specific data keys
pthread_mutex	Mutexes
pthread_mutexattr	Mutex attributes objects

Data type identifiers end with **_t**. Subroutine and macro names end with an **_** (underscore), followed by a name identifying the action performed by the subroutine or the macro. For example, **pthread_attr_init** is a threads library identifier (**pthread_**) that concerns thread attributes objects (**attr**) and is an initialization subroutine (**_init**).

Explicit macro identifiers are in uppercase letters. Some subroutines may, however, be implemented as macros, although their names are in lowercase letters.

pthread implementation files

This section describes the **pthread** implementation files.

The following AIX files provide the implementation of pthreads:

Implementaion	Description
/usr/include/pthread.h	C/C++ header with most pthread definitions.
/usr/include/sched.h	C/C++ header with some scheduling definitions.
/usr/include/unistd.h	C/C++ header with pthread_atfork() definition.
/usr/include/sys/limits.h	C/C++ header with some pthread definitions.
/usr/include/sys/pthdebug.h	C/C++ header with most pthread debug definitions.
/usr/include/sys/sched.h	C/C++ header with some scheduling definitions.
/usr/include/sys/signal.h	C/C++ header with pthread_kill() and pthread_sigmask() definitions.
/usr/include/sys/types.h	C/C++ header with some pthread definitions.
/usr/lib/libpthreads.a	32-bit/64-bit library providing UNIX98 and POSIX 1003.1c pthreads.
/usr/lib/libpthreads_compat.a	32-bit only library providing POSIX 1003.1c Draft 7 pthreads.
/usr/lib/profiled/libpthreads.a	Profiled 32-bit/64-bit library providing UNIX98 and POSIX 1003.1c pthreads.
/usr/lib/profiled/libpthreads_compat.a	Profiled 32-bit only library providing POSIX 1003.1c Draft 7 pthreads.

Threadsafe and threaded libraries in AIX

This section describes the thread libraries in the AIX.

By default, all applications are now considered "threaded," even though most are of the case "single threaded." These threadsafe libraries are as follows:

Threadsafe libraries		
libbsd.a	libc.a	libm.a

Threadsafe libraries		
libsvid.a	libtli.a	libxti.a
libnetvc.a		

POSIX threads libraries

The following POSIX threads libraries are available:

libpthreads.a POSIX threads library

The **libpthreads.a** library is based on the POSIX 1003.1c industry standard for a portable user threads API. Any program written for use with a POSIX thread library can be ported for use with another POSIX threads library; only the performance and very few subroutines of the threads library are implementation-dependent. To enhance the portability of the threads library, the POSIX standard made the implementation of several programming facilities optional. For more information about checking the POSIX options, see [Threads Library Options](#).

libpthread_compat.a POSIX draft 7 threads library

AIX provides binary compatibility for existing multi-threads applications that were coded to Draft 7 of the POSIX thread standard. These applications will run without relinking. The **libpthread_compat.a** library is provided only for compatibility with earlier versions of applications written by using the Draft 7 of the POSIX Thread Standard. All new applications must use the **libpthreads.a** library, which supports both 32-bit and 64-bit applications. The **libpthread_compat.a** library supports only 32-bit applications. Beginning with AIX 5.1, the **libpthreads.a** library supports the Single UNIX Specification, Version 2, which includes the final POSIX 1003.1c Pthread Standard.

Related concepts

[Benefits of threads](#)

Multithreaded programs can improve performance compared to traditional parallel programs that use multiple processes. Furthermore, improved performance can be obtained on multiprocessor systems using threads.

[Threads library options](#)

This section describes special attributes of threads, mutexes, and condition variables.

[Developing multithreaded programs](#)

Developing multithreaded programs is similar to developing programs with multiple processes. Developing programs also consists of compiling and debugging the code.

Creating threads

Thread creation differs from process creation in that no parent-child relation exists between threads.

All threads, except the *initial thread* automatically created when a process is created, are on the same hierarchical level. A thread does not maintain a list of created threads, nor does it know the thread that created it.

When creating a thread, an entry-point routine and an argument must be specified. Every thread has an entry-point routine with one argument. The same entry-point routine may be used by several threads.

A thread has attributes, which specify the characteristics of the thread. To control thread attributes, a thread attributes object must be defined before creating the thread.

Thread attributes Object

The thread attributes are stored in an opaque object, the *thread attributes object*, used when creating the thread. It contains several attributes, depending on the implementation of POSIX options. The object is accessed through a variable of type **pthread_attr_t**. In AIX, the **pthread_attr_t** data type is a pointer to a structure; on other systems, it may be a structure or another data type.

Creating and destroying the thread attributes object

o

The thread attributes object is initialized to default values by the **pthread_attr_init** subroutine. The attributes are handled by subroutines. The thread attributes object is destroyed by the **pthread_attr_destroy** subroutine. This subroutine can release storage dynamically allocated by the **pthread_attr_init** subroutine, depending on the implementation of the threads library.

In the following example, a thread attributes object is created and initialized with default values, then used and finally destroyed:

```
pthread_attr_t attributes;
    /* the attributes object is created */
if (!pthread_attr_init(&attributes)) {
    /* the attributes object is initialized */
    ...
    /* using the attributes object */
    ...
    pthread_attr_destroy(&attributes);
    /* the attributes object is destroyed */
}
```

The same attributes object can be used to create several threads. It can also be modified between two thread creations. When the threads are created, the attributes object can be destroyed without affecting the threads created with it.

Detachstate attribute

The following attribute is always defined:

Detachstate

Specifies the detached state of a thread.

The value of the attribute is returned by the **pthread_attr_getdetachstate** subroutine; it can be set by the **pthread_attr_setdetachstate** subroutine. Possible values for this attributes are the following symbolic constants:

PTHREAD_CREATE_DETACHED

Specifies that the thread will be created in the detached state

PTHREAD_CREATE_JOINABLE

Specifies that the thread will be created in the joinable state

The default value is **PTHREAD_CREATE_JOINABLE**.

If you create a thread in the joinable state, you must call the **pthread_join** subroutine with the thread. Otherwise, you may run out of storage space when creating new threads, because each thread takes up a significant amount of memory. For more information on the **pthread_join** subroutine, see [Calling the pthread_join Subroutine](#).

Other threads attributes

AIX also defines the following attributes, which are intended for advanced programs and may require special execution privilege to take effect. Most programs operate correctly with the default settings. The use of the following attributes is explained in [Using the inheritsched Attribute](#).

Contention Scope

Specifies the contention scope of a thread

Inheritsched

Specifies the inheritance of scheduling properties of a thread

Schedparam

Specifies the scheduling parameters of a thread

Schedpolicy

Specifies the scheduling policy of a thread

The use of the following stack attributes is explained in [Stack Attributes](#).

Stacksize

Specifies the size of the thread's stack

Stackaddr

Specifies the address of the thread's stack

Guardsize

Specifies the size of the guard area of the thread's stack

Creating a thread using the `pthread_create` subroutine

A thread is created by calling the `pthread_create` subroutine. This subroutine creates a new thread and makes it runnable.

Using the thread attributes object

When calling the `pthread_create` subroutine, you may specify a thread attributes object. If you specify a `NULL` pointer, the created thread will have the default attributes. Thus, the following code fragment:

```
pthread_t thread;
pthread_attr_t attr;
...
pthread_attr_init(&attr);
pthread_create(&thread, &attr, init_routine, NULL);
pthread_attr_destroy(&attr);
```

is equivalent to the following:

```
pthread_t thread;
...
pthread_create(&thread, NULL, init_routine, NULL);
```

Entry point routine

When calling the `pthread_create` subroutine, you must specify an entry-point routine. This routine, provided by your program, is similar to the `main` routine for the process. It is the first user routine executed by the new thread. When the thread returns from this routine, the thread is automatically terminated.

The entry-point routine has one parameter, a void pointer, specified when calling the `pthread_create` subroutine. You may specify a pointer to some data, such as a string or a structure. The creating thread (the one calling the `pthread_create` subroutine) and the created thread must agree upon the actual type of this pointer.

The entry-point routine returns a void pointer. After the thread termination, this pointer is stored by the threads library unless the thread is detached. For more information about using this pointer, see [Returning Information from a Thread](#).

Returned information

The `pthread_create` subroutine returns the thread ID of the new thread. The caller can use this thread ID to perform various operations on the thread.

Depending on the scheduling parameters of both threads, the new thread may start running before the call to the `pthread_create` subroutine returns. It may even happen that, when the `pthread_create` subroutine returns, the new thread has already terminated. The thread ID returned by the `pthread_create` subroutine through the `thread` parameter is then already invalid. It is, therefore, important to check for the `ESRCH` error code returned by threads library subroutines using a thread ID as a parameter.

If the `pthread_create` subroutine is unsuccessful, no new thread is created, the thread ID in the `thread` parameter is invalid, and the appropriate error code is returned. For more information, see [Example of a Multi-Threaded Program](#).

Handling Thread IDs

The thread ID of a newly created thread is returned to the creating thread through the `thread` parameter. The current thread ID is returned by the `pthread_self` subroutine.

A thread ID is an opaque object; its type is `pthread_t`. In AIX, the `pthread_t` data type is an integer. On other systems, it may be a structure, a pointer, or any other data type.

To enhance the portability of programs using the threads library, the thread ID should always be handled as an opaque object. For this reason, thread IDs should be compared using the **pthread_equal** subroutine. Never use the C equality operator (**==**), because the **pthread_t** data type may be neither an arithmetic type nor a pointer.

Related concepts

Developing multithreaded programs

Developing multithreaded programs is similar to developing programs with multiple processes.

Developing programs also consists of compiling and debugging the code.

Terminating threads

A thread automatically terminates when it returns from its entry-point routine.

A thread can also explicitly terminate itself or terminate any other thread in the process, using a mechanism called *cancelation*. Because all threads share the same data space, a thread must perform cleanup operations at termination time; the threads library provides cleanup handlers for this purpose.

Exiting a thread

A process can exit at any time when a thread calls the **exit** subroutine. Similarly, a thread can exit at any time by calling the **pthread_exit** subroutine.

Calling the **exit** subroutine terminates the entire process, including all its threads. In a multithreaded program, the **exit** subroutine should only be used when the entire process needs to be terminated; for example, in the case of an unrecoverable error. The **pthread_exit** subroutine should be preferred, even for exiting the initial thread.

Calling the **pthread_exit** subroutine terminates the calling thread. The *status* parameter is saved by the library and can be further used when joining the terminated thread. Calling the **pthread_exit** subroutine is similar, but not identical, to returning from the thread's initial routine. The result of returning from the thread's initial routine depends on the thread:

- Returning from the initial thread implicitly calls the **exit** subroutine, thus terminating all the threads in the process.
- Returning from another thread implicitly calls the **pthread_exit** subroutine. The return value has the same role as the *status* parameter of the **pthread_exit** subroutine.

To avoid implicitly calling the **exit** subroutine, to use the **pthread_exit** subroutine to exit a thread.

Exiting the initial thread (for example, by calling the **pthread_exit** subroutine from the **main** routine) does not terminate the process. It terminates only the initial thread. If the initial thread is terminated, the process will be terminated when the last thread in it terminates. In this case, the process return code is 0.

The following program displays exactly 10 messages in each language. This is accomplished by calling the **pthread_exit** subroutine in the **main** routine after creating the two threads, and creating a loop in the **Thread** routine.

```
#include <pthread.h>    /* include file for pthreads - the 1st */
#include <stdio.h>      /* include file for printf()          */

void *Thread(void *string)
{
    int i;
    for (i=0; i<10; i++)
        printf("%s\n", (char *)string);
    pthread_exit(NULL);
}

int main()
{
    char *e_str = "Hello!";
    char *f_str = "Bonjour !";

    pthread_t e_th;
    pthread_t f_th;
```

```

int rc;

rc = pthread_create(&e_th, NULL, Thread, (void *)e_str);
if (rc)
    exit(-1);
rc = pthread_create(&f_th, NULL, Thread, (void *)f_str);
if (rc)
    exit(-1);
pthread_exit(NULL);
}

```

The **pthread_exit** subroutine releases any thread-specific data, including the thread's stack. Any data allocated on the stack becomes invalid, because the stack is freed and the corresponding memory may be reused by another thread. Therefore, thread synchronization objects (mutexes and condition variables) allocated on a thread's stack must be destroyed before the thread calls the **pthread_exit** subroutine.

Unlike the **exit** subroutine, the **pthread_exit** subroutine does not clean up system resources shared among threads. For example, files are not closed by the **pthread_exit** subroutine, because they may be used by other threads.

Canceling a thread

The thread cancelation mechanism allows a thread to terminate the execution of any other thread in the process in a controlled manner. The target thread (that is, the one that is being canceled) can hold cancelation requests pending in a number of ways and perform application-specific cleanup processing when the notice of cancelation is acted upon. When canceled, the thread implicitly calls the **pthread_exit((void *)-1)** subroutine.

The cancelation of a thread is requested by calling the **pthread_cancel** subroutine. When the call returns, the request has been registered, but the thread may still be running. The call to the **pthread_cancel** subroutine is unsuccessful only when the specified thread ID is not valid.

Cancelability state and type

The cancelability state and type of a thread determines the action taken upon receipt of a cancelation request. Each thread controls its own cancelability state and type with the **pthread_setcancelstate** and **pthread_setcanceltype** subroutines.

The following possible cancelability states and cancelability types lead to three possible cases, as shown in the following table.

Cancelability State	Cancelability Type	Resulting Case
Disabled	Any (the type is ignored)	Disabled cancelability
Enabled	Deferred	Deferred cancelability
Enabled	Asynchronous	Asynchronous cancelability

The possible cases are described as follows:

- *Disabled cancelability.* Any cancelation request is set pending, until the cancelability state is changed or the thread is terminated in another way.

A thread should disable cancelability only when performing operations that cannot be interrupted. For example, if a thread is performing some complex file-save operations (such as an indexed database) and is canceled during the operation, the files may be left in an inconsistent state. To avoid this, the thread should disable cancelability during the file save operations.

- *Deferred cancelability.* Any cancelation request is set pending, until the thread reaches the next cancelation point. It is the default cancelability state.

This cancelability state ensures that a thread can be cancelled, but limits the cancelation to specific moments in the thread's execution, called *cancelation points*. A thread canceled on a cancelation point leaves the system in a safe state; however, user data may be inconsistent or locks may be held by the canceled thread. To avoid these situations, use cleanup handlers or disable cancelability within critical regions. For more information, see [Using Cleanup Handlers](#).

- *Asynchronous cancelability*. Any cancelation request is acted upon immediately.

A thread that is asynchronously canceled while holding resources may leave the process, or even the system, in a state from which it is difficult or impossible to recover. For more information about *async-cancel safety*, see [Async-Cancel Safety](#).

Async-cancel safety

A function is said to be *async-cancel safe* if it is written so that calling the function with asynchronous cancelability enabled does not cause any resource to be corrupted, even if a cancelation request is delivered at any arbitrary instruction.

Any function that gets a resource as a side effect cannot be made *async-cancel safe*. For example, if the **malloc** subroutine is called with asynchronous cancelability enabled, it might acquire the resource successfully, but as it was returning to the caller, it could act on a cancelation request. In such a case, the program would have no way of knowing whether the resource was acquired or not.

For this reason, most library routines cannot be considered *async-cancel safe*. It is recommended that you use asynchronous cancelability only if you are sure only to perform operations that do not hold resources and only to call library routines that are *async-cancel safe*.

The following subroutines are *async-cancel safe*; they ensure that cancelation will be handled correctly, even if asynchronous cancelability is enabled:

- **pthread_cancel**
- **pthread_setcancelstate**
- **pthread_setcanceltype**

An alternative to asynchronous cancelability is to use deferred cancelability and to add explicit cancelation points by calling the **pthread_testcancel** subroutine.

Cancelation points

Cancelation points are points inside of certain subroutines where a thread must act on any pending cancelation request if deferred cancelability is enabled. All of these subroutines may block the calling thread or compute indefinitely.

An explicit cancelation point can also be created by calling the **pthread_testcancel** subroutine. This subroutine simply creates a cancelation point. If deferred cancelability is enabled, and if a cancelation request is pending, the request is acted upon and the thread is terminated. Otherwise, the subroutine simply returns.

Other cancelation points occur when calling the following subroutines:

- **pthread_cond_wait**
- **pthread_cond_timedwait**
- **pthread_join**

The **pthread_mutex_lock** and **pthread_mutex_trylock** subroutines do not provide a cancelation point. If they did, all functions calling these subroutines (and many functions do) would provide a cancelation point. Having too many cancelation points makes programming very difficult, requiring either lots of disabling and restoring of cancelability or extra effort in trying to arrange for reliable cleanup at every possible place. For more information about these subroutines, see [Using Mutexes](#).

Cancelation points occur when a thread is executing the following functions:

Function

aio_suspend	close
creat	fcntl
fsync	getmsg
getpmsg	lockf

Function

mq_receive	mq_send
msgrcv	msgsnd
msync	nanosleep
open	pause
poll	pread
pthread_cond_timedwait	pthread_cond_wait
pthread_join	pthread_testcancel
putpmsg	pwrite
read	readv
select	sem_wait
sigpause	sigsuspend
sigtimedwait	sigwait
sigwaitinfo	sleep
system	tcdrain
usleep	wait
wait3	waitid
waitpid	write
writev	

A cancelation point can also occur when a thread is executing the following functions:

Function Function Function

catclose	catgets	catopen
closedir	closelog	ctermid
dbm_close	dbm_delete	dbm_fetch
dbm_nextkey	dbm_open	dbm_store
dlclose	dlopen	endgrent
endpwent	fwprintf	fwrite
fwscanf	getc	getc_unlocked
getchar	getchar_unlocked	getcwd
getdate	getgrent	getgrgid
getgrgid_r	getgrnam	getgrnam_r
getlogin	getlogin_r	popen
printf	putc	putc_unlocked
putchar	putchar_unlocked	puts
pututxline	putw	putwc
putwchar	readdir	readdir_r
remove	rename	rewind
endutxent	fclose	fcntl

Function Function Function

fflush	fgetc	fgetpos
fgets	fgetwc	fgetws
fopen	fprintf	fputc
fputs	getpwent	getpwnam
getpwnam_r	getpwuid	getpwuid_r
gets	getutxent	getutxid
getutxline	getw	getwc
getwchar	getwd	rewinddir
scanf	seekdir	semop
setgrent	setpwent	setutxent
strerror	syslog	tmpfile
tmpnam	ttyname	ttyname_r
fputwc	fputws	fread
freopen	fscanf	fseek
fseeko	fsetpos	ftell
ftello	ftw	glob
iconv_close	iconv_open	ioctl
lseek	mkstemp	nftw
opendir	openlog	pclose
perror	ungetc	ungetwc
unlink	vwprintf	vfwprintf
vprintf	wprintf	wprintf
wscanf		

The side effects of acting upon a cancelation request while suspended during a call of a function is the same as the side effects that may be seen in a single-threaded program when a call to a function is interrupted by a signal and the given function returns [EINTR]. Any such side effects occur before any cancelation cleanup handlers are called.

Whenever a thread has cancelability enabled and a cancelation request has been made with that thread as the target and the thread calls the **pthread_testcancel** subroutine, the cancelation request is acted upon before the **pthread_testcancel** subroutine returns. If a thread has cancelability enabled and the thread has an asynchronous cancelation request pending and the thread is suspended at a cancelation point waiting for an event to occur, the cancelation request will be acted upon. However, if the thread is suspended at a cancelation point and the event that it is waiting for occurs before the cancelation request is acted upon, the sequence of events determines whether the cancelation request is acted upon or whether the request remains pending and the thread resumes normal execution.

Cancelation example

In the following example, both "writer" threads are canceled after 10 seconds, and after they have written their message at least five times.

```
#include <pthread.h>    /* include file for pthreads - the 1st */
#include <stdio.h>      /* include file for printf()          */
#include <unistd.h>     /* include file for sleep()          */
```

```

void *Thread(void *string)
{
    int i;
    int o_state;

    /* disables cancelability */
    pthread_setcancelstate(PTHREAD_CANCEL_DISABLE, &o_state);

    /* writes five messages */
    for (i=0; i<5; i++)
        printf("%s\n", (char *)string);

    /* restores cancelability */
    pthread_setcancelstate(o_state, &o_state);

    /* writes further */
    while (1)
        printf("%s\n", (char *)string);
    pthread_exit(NULL);
}

```

```

int main()
{
    char *e_str = "Hello!";
    char *f_str = "Bonjour !";

    pthread_t e_th;
    pthread_t f_th;

    int rc;

    /* creates both threads */
    rc = pthread_create(&e_th, NULL, Thread, (void *)e_str);
    if (rc)
        return -1;
    rc = pthread_create(&f_th, NULL, Thread, (void *)f_str);
    if (rc)
        return -1;

```

```

    /* sleeps a while */
    sleep(10);

    /* requests cancelation */
    pthread_cancel(e_th);
    pthread_cancel(f_th);

    /* sleeps a bit more */
    sleep(10);
    pthread_exit(NULL);
}

```

Timer and sleep subroutines

Timer routines execute in the context of the calling thread. Thus, if a timer expires, the watchdog timer function is called in the thread's context. When a process or thread goes to sleep, it relinquishes the processor. In a multithreaded process, only the calling thread is put to sleep.

Using cleanup handlers

Cleanup handlers provide a portable mechanism for releasing resources and restoring invariants when a thread terminates.

Calling cleanup handlers

Cleanup handlers are specific to each thread. A thread can have several cleanup handlers; they are stored in a thread-specific LIFO (last-in, first-out) stack. Cleanup handlers are all called in the following cases:

- The thread returns from its entry-point routine.
- The thread calls the **pthread_exit** subroutine.
- The thread acts on a cancelation request.

A cleanup handler is pushed onto the cleanup stack by the **pthread_cleanup_push** subroutine. The **pthread_cleanup_pop** subroutine pops the topmost cleanup handler from the stack and optionally executes it. Use this subroutine when the cleanup handler is no longer needed.

The cleanup handler is a user-defined routine. It has one parameter, a void pointer, specified when calling the **pthread_cleanup_push** subroutine. You can specify a pointer to some data that the cleanup handler needs to perform its operation.

In the following example, a buffer is allocated for performing some operation. With deferred cancelability enabled, the operation can be stopped at any cancellation point. In that case, a cleanup handler is established to release the buffer.

```
/* the cleanup handler */
cleaner(void *buffer)
{
    free(buffer);
}

/* fragment of another routine */
...
myBuf = malloc(1000);
if (myBuf != NULL) {

    pthread_cleanup_push(cleaner, myBuf);

    /*
     *   perform any operation using the buffer,
     *   including calls to other functions
     *   and cancellation points
     */

    /* pops the handler and frees the buffer in one call */
    pthread_cleanup_pop(1);
}
```

Using deferred cancelability ensures that the thread will not act on any cancellation request between the buffer allocation and the registration of the cleanup handler, because neither the **malloc** subroutine nor the **pthread_cleanup_push** subroutine provides any cancellation point. When popping the cleanup handler, the handler is executed, releasing the buffer. More complex programs may not execute the handler when popping it, because the cleanup handler should be thought of as an "emergency exit" for the protected portion of code.

Balancing the push and pop operations

The **pthread_cleanup_push** and **pthread_cleanup_pop** subroutines should always appear in pairs within the same lexical scope; that is, within the same function and the same statement block. They can be thought of as left and right parentheses enclosing a protected portion of code.

The reason for this rule is that on some systems these subroutines are implemented as macros. The **pthread_cleanup_push** subroutine is implemented as a left brace, followed by other statements:

```
#define pthread_cleanup_push(rtm,arg) { \
    /* other statements */
```

The **pthread_cleanup_pop** subroutine is implemented as a right brace, following other statements:

```
#define pthread_cleanup_pop(ex) \
    /* other statements */ \
}
```

Adhere to the balancing rule for the **pthread_cleanup_push** and **pthread_cleanup_pop** subroutines to avoid compiler errors or unexpected behavior of your programs when porting to other systems.

In AIX, the **pthread_cleanup_push** and **pthread_cleanup_pop** subroutines are library routines, and can be unbalanced within the same statement block. However, they must be balanced in the program, because the cleanup handlers are stacked.

Subroutine

[pthread_attr_destroy](#)

[pthread_attr_getdetachstate](#)

[pthread_attr_init](#)

[pthread_cancel](#)

[pthread_cleanup_pop](#)

[pthread_cleanup_push](#)

[pthread_create](#)

[pthread_equal](#)

[pthread_exit](#)

[pthread_self](#)

[pthread_setcancelstate](#)

[pthread_setcanceltype](#)

[pthread_testcancel](#)

Description

Deletes a thread attributes object.

Returns the value of the detachstate attribute of a thread attributes object.

Creates a thread attributes object and initializes it with default values.

Requests the cancelation of a thread.

Removes, and optionally executes, the routine at the top of the calling thread's cleanup stack.

Pushes a routine onto the calling thread's cleanup stack.

Creates a new thread, initializes its attributes, and makes it runnable.

Compares two thread IDs.

Terminates the calling thread.

Returns the calling thread's ID.

Sets the calling thread's cancelability state.

Sets the calling thread's cancelability type.

Creates a cancelation point in the calling thread.

Related concepts

[One-time initializations](#)

Some C libraries are designed for dynamic initialization, in which the global initialization for the library is performed when the first procedure in the library is called.

Synchronization overview

One main benefit of using threads is the ease of using synchronization facilities.

To effectively interact, threads must synchronize their activities. This includes:

- Implicit communication through the modification of shared data
- Explicit communication by informing each other of events that have occurred

More complex synchronization objects can be built using the primitive objects.

The threads library provides the following synchronization mechanisms: Although primitive, these powerful mechanisms can be used to build more complex mechanisms.

The threads library provides the following synchronization mechanisms:

- Mutexes (See [Using Mutexes](#))
- Condition variables (See [Using Condition Variables](#))
- Read-write locks (See [Using Read-Write Locks](#))
- Joins (See [Joining Threads](#))

Although primitive, these powerful mechanisms can be used to build more complex mechanisms.

Related concepts

[Creating complex synchronization objects](#)

The subroutines provided in the threads library can be used as primitives to build more complex synchronization objects.

Using mutexes

A *mutex* is a mutual exclusion lock. Only one thread can hold the lock.

Mutexes are used to protect data or other resources from concurrent access. A mutex has attributes, which specify the characteristics of the mutex.

Mutex attributes object

Like threads, mutexes are created with the help of an attributes object. The mutex attributes object is an abstract object, containing several attributes, depending on the implementation of POSIX options. It is accessed through a variable of type **pthread_mutexattr_t**. In AIX, the **pthread_mutexattr_t** data type is a pointer; on other systems, it may be a structure or another data type.

Creating and destroying the mutex attributes object

The mutex attributes object is initialized to default values by the **pthread_mutexattr_init** subroutine. The attributes are handled by subroutines. The thread attributes object is destroyed by the **pthread_mutexattr_destroy** subroutine. This subroutine may release storage dynamically allocated by the **pthread_mutexattr_init** subroutine, depending on the implementation of the threads library.

In the following example, a mutex attributes object is created and initialized with default values, then used and finally destroyed:

```
pthread_mutexattr_t attributes;
    /* the attributes object is created */
...
if (!pthread_mutexattr_init(&attributes)) {
    /* the attributes object is initialized */
    ...
    /* using the attributes object */
    ...
    pthread_mutexattr_destroy(&attributes);
    /* the attributes object is destroyed */
}
```

The same attributes object can be used to create several mutexes. It can also be modified between mutex creations. When the mutexes are created, the attributes object can be destroyed without affecting the mutexes created with it.

Mutex attributes

The following mutex attributes are defined:

Attribute	Description
Protocol	Specifies the protocol used to prevent priority inversions for a mutex. This attribute depends on either the priority inheritance or the priority protection POSIX option.
Process-shared	Specifies the process sharing of a mutex. This attribute depends on the process sharing POSIX option.

For more information on these attributes, see [Threads Library Options](#) and [Synchronization Scheduling](#).

Creating and destroying mutexes

A mutex is created by calling the **pthread_mutex_init** subroutine. You may specify a mutex attributes object. If you specify a **NULL** pointer, the mutex will have the default attributes. Thus, the following code fragment:

```
pthread_mutex_t mutex;
pthread_mutexattr_t attr;
...
pthread_mutexattr_init(&attr);
```

```
pthread_mutex_init(&mutex, &attr);
pthread_mutexattr_destroy(&attr);
```

is equivalent to the following:

```
pthread_mutex_t mutex;
...
pthread_mutex_init(&mutex, NULL);
```

The ID of the created mutex is returned to the calling thread through the *mutex* parameter. The mutex ID is an opaque object; its type is **pthread_mutex_t**. In AIX, the **pthread_mutex_t** data type is a structure; on other systems, it might be a pointer or another data type.

A mutex must be created once. However, avoid calling the **pthread_mutex_init** subroutine more than once with the same *mutex* parameter (for example, in two threads concurrently executing the same code). Ensuring the uniqueness of a mutex creation can be done in the following ways:

- Calling the **pthread_mutex_init** subroutine prior to the creation of other threads that will use this mutex; for example, in the initial thread.
- Calling the **pthread_mutex_init** subroutine within a one time initialization routine. For more information, see [One-Time Initializations](#).
- Using a static mutex initialized by the **PTHREAD_MUTEX_INITIALIZER** static initialization macro; the mutex will have default attributes.

After the mutex is no longer needed, destroy it by calling the **pthread_mutex_destroy** subroutine. This subroutine may reclaim any storage allocated by the **pthread_mutex_init** subroutine. After having destroyed a mutex, the same **pthread_mutex_t** variable can be reused to create another mutex. For example, the following code fragment is valid, although not very practical:

```
pthread_mutex_t mutex;
...
for (i = 0; i < 10; i++) {
    /* creates a mutex */
    pthread_mutex_init(&mutex, NULL);

    /* uses the mutex */

    /* destroys the mutex */
    pthread_mutex_destroy(&mutex);
}
```

Like any system resource that can be shared among threads, a mutex allocated on a thread's stack must be destroyed before the thread is terminated. The threads library maintains a linked list of mutexes. Thus, if the stack where a mutex is allocated is freed, the list will be corrupted.

Types of mutexes

The type of mutex determines how the mutex behaves when it is operated on. The following types of mutexes exist:

PTHREAD_MUTEX_DEFAULT or **PTHREAD_MUTEX_NORMAL**

Results in a deadlock if the same pthread tries to lock it a second time using the **pthread_mutex_lock** subroutine without first unlocking it. This is the default type.

PTHREAD_MUTEX_ERRORCHECK

Avoids deadlocks by returning a non-zero value if the same thread attempts to lock the same mutex more than once without first unlocking the mutex.

PTHREAD_MUTEX_RECURSIVE

Allows the same pthread to recursively lock the mutex using the **pthread_mutex_lock** subroutine without resulting in a deadlock or getting a non-zero return value from **pthread_mutex_lock**. The same pthread has to call the **pthread_mutex_unlock** subroutine the same number of times as it called **pthread_mutex_lock** subroutine in order to unlock the mutex for other pthreads to use.

When a mutex attribute is first created, it has a default type of **PTHREAD_MUTEX_NORMAL**. After creating the mutex, the type can be changed using the **pthread_mutexattr_settype** API library call.

The following is an example of creating and using a recursive mutex type:

```
pthread_mutexattr_t   attr;
pthread_mutex_t       mutex;

pthread_mutexattr_settype(&attr, PTHREAD_MUTEX_RECURSIVE);
pthread_mutex_init(&mutex, &attr);

struct {
    int a;
    int b;
    int c;
} A;

f()
{
    pthread_mutex_lock(&mutex);
    A.a++;
    g();
    A.c = 0;
    pthread_mutex_unlock(&mutex);
}

g()
{
    pthread_mutex_lock(&mutex);
    A.b += A.a;
    pthread_mutex_unlock(&mutex);
}
```

Locking and unlocking mutexes

A mutex is a simple lock, having two states: locked and unlocked. When it is created, a mutex is unlocked. The **pthread_mutex_lock** subroutine locks the specified mutex under the following conditions:

- If the mutex is unlocked, the subroutine locks it.
- If the mutex is already locked by another thread, the subroutine blocks the calling thread until the mutex is unlocked.
- If the mutex is already locked by the calling thread, the subroutine might block forever or return an error depending on the type of mutex.

The **pthread_mutex_trylock** subroutine acts like the **pthread_mutex_lock** subroutine without blocking the calling thread under the following conditions:

- If the mutex is unlocked, the subroutine locks it.
- If the mutex is already locked by any thread, the subroutine returns an error.

The thread that locked a mutex is often called the *owner* of the mutex.

The **pthread_mutex_unlock** subroutine resets the specified mutex to the unlocked state if it is owned by the calling mutex under the following conditions:

- If the mutex was already unlocked, the subroutine returns an error.
- If the mutex was owned by the calling thread, the subroutine unlocks the mutex.
- If the mutex was owned by another thread, the subroutine might return an error or unlock the mutex depending on the type of mutex. Unlocking the mutex is not recommended because mutexes are usually locked and unlocked by the same pthread.

Because locking does not provide a cancellation point, a thread blocked while waiting for a mutex cannot be canceled. Therefore, it is recommended that you use mutexes only for short periods of time, as in instances where you are protecting data from concurrent access. For more information, see [Cancellation Points](#) and [Canceling a Thread](#).

Protecting data with mutexes

Mutexes are intended to serve either as a low-level primitive from which other thread synchronization functions can be built or as a data protection lock. For more information about implementing long locks and writer-priority readers/writers locks see [“Using mutexes”](#) on page 424.

Mutex usage example

Mutexes can be used to protect data from concurrent access. For example, a database application may create several threads to handle several requests concurrently. The database itself is protected by a mutex called **db_mutex**. For example:

```
/* the initial thread */
pthread_mutex_t mutex;
int i;
...
pthread_mutex_init(&mutex, NULL); /* creates the mutex */
for (i = 0; i < num_req; i++) /* loop to create threads */
    pthread_create(th + i, NULL, rtn, &mutex);
...
pthread_mutex_destroy(&mutex); /* destroys the mutex */
...

/* the request handling thread */
...
pthread_mutex_lock(&db_mutex); /* waits for a request */
/* locks the database */
...
pthread_mutex_unlock(&db_mutex); /* handles the request */
/* unlocks the database */
...
```

The initial thread creates the mutex and all the request-handling threads. The mutex is passed to the thread using the parameter of the thread's entry point routine. In a real program, the address of the mutex may be a field of a more complex data structure passed to the created thread.

Avoiding Deadlocks

There are a number of ways that a multithreaded application can deadlock. Following are some examples:

- A mutex created with the default type, **PTHREAD_MUTEX_NORMAL**, cannot be relocked by the same pthread without resulting in a deadlock.
- An application can deadlock when locking mutexes in reverse order. For example, the following code fragment can produce a deadlock between threads A and B.

```
/* Thread A */
pthread_mutex_lock(&mutex1);
pthread_mutex_lock(&mutex2);

/* Thread B */
pthread_mutex_lock(&mutex2);
pthread_mutex_lock(&mutex1);
```

- An application can deadlock in what is called *resource* deadlock. For example:

```
struct {
    pthread_mutex_t mutex;
    char *buf;
} A;

struct {
    pthread_mutex_t mutex;
    char *buf;
} B;

struct {
    pthread_mutex_t mutex;
    char *buf;
} C;

use_all_buffers()
{
    pthread_mutex_lock(&A.mutex);
    /* use buffer A */

    pthread_mutex_lock(&B.mutex);
    /* use buffers B */

    pthread_mutex_lock(&C.mutex);
    /* use buffer C */

    /* All done */
    pthread_mutex_unlock(&C.mutex);
}
```

```

        pthread_mutex_unlock(&B.mutex);
        pthread_mutex_unlock(&A.mutex);
    }

    use_buffer_a()
    {
        pthread_mutex_lock(&A.mutex);
        /* use buffer A */
        pthread_mutex_unlock(&A.mutex);
    }

    functionB()
    {
        pthread_mutex_lock(&B.mutex);
        /* use buffer B */
        if (..some condition)
        {
            use_buffer_a();
        }
        pthread_mutex_unlock(&B.mutex);
    }

    /* Thread A */
    use_all_buffers();

    /* Thread B */
    functionB();

```

This application has two threads, thread A and thread B. Thread B starts to run first, then thread A starts shortly thereafter. If thread A executes **use_all_buffers()** and successfully locks **A.mutex**, it will then block when it tries to lock **B.mutex**, because thread B has already locked it. While thread B executes **functionB** and **some_condition** occurs while thread A is blocked, thread B will now also block trying to acquire **A.mutex**, which is already locked by thread A. This results in a deadlock.

The solution to this deadlock is for each thread to acquire all the resource locks that it needs before using the resources. If it cannot acquire the locks, it must release them and start again.

Mutexes and race conditions

Mutual exclusion locks (mutexes) can prevent data inconsistencies due to race conditions. A race condition often occurs when two or more threads must perform operations on the same memory area, but the results of computations depends on the order in which these operations are performed.

Consider, for example, a single counter, *X*, that is incremented by two threads, A and B. If *X* is originally 1, then by the time threads A and B increment the counter, *X* should be 3. Both threads are independent entities and have no synchronization between them. Although the C statement *X++* looks simple enough to be atomic, the generated assembly code may not be, as shown in the following pseudo-assembler code:

```

move    X, REG
inc     REG
move    REG, X

```

If both threads in the previous example are executed concurrently on two CPUs, or if the scheduling makes the threads alternatively execute on each instruction, the following steps may occur:

1. Thread A executes the first instruction and puts *X*, which is 1, into the thread A register. Then thread B executes and puts *X*, which is 1, into the thread B register. The following example illustrates the resulting registers and the contents of memory *X*.

```

Thread A Register = 1
Thread B Register = 1
Memory X         = 1

```

2. Thread A executes the second instruction and increments the content of its register to 2. Then thread B increments its register to 2. Nothing is moved to memory *X*, so memory *X* stays the same. The following example illustrates the resulting registers and the contents of memory *X*.

```

Thread A Register = 2
Thread B Register = 2
Memory X         = 1

```

3. Thread A moves the content of its register, which is now 2, into memory X. Then thread B moves the content of its register, which is also 2, into memory X, overwriting thread A's value. The following example illustrates the resulting registers and the contents of memory X.

```
Thread A Register = 2
Thread B Register = 2
Memory X         = 2
```

In most cases, thread A and thread B execute the three instructions one after the other, and the result would be 3, as expected. Race conditions are usually difficult to discover, because they occur intermittently.

To avoid this race condition, each thread should lock the data before accessing the counter and updating memory X. For example, if thread A takes a lock and updates the counter, it leaves memory X with a value of 2. After thread A releases the lock, thread B takes the lock and updates the counter, taking 2 as its initial value for X and increment it to 3, the expected result.

Using condition variables

Condition variables allow threads to wait until some event or condition has occurred.

A condition variable has attributes that specify the characteristics of the condition. Typically, a program uses the following objects:

- A boolean variable, indicating whether the condition is met
- A mutex to serialize the access to the boolean variable
- A condition variable to wait for the condition

Using a condition variable requires some effort from the programmer. However, condition variables allow the implementation of powerful and efficient synchronization mechanisms. For more information about implementing long locks and semaphores with condition variables, see [Creating Complex Synchronization Objects](#).

When a thread is terminated, its storage may not be reclaimed, depending on an attribute of the thread. Such threads can be joined by other threads and return information to them. A thread that wants to join another thread is blocked until the target thread terminates. This joint mechanism is a specific case of condition-variable usage, the condition is the thread termination.

Condition attributes object

Like threads and mutexes, condition variables are created with the help of an attributes object. The *condition attributes object* is an abstract object, containing at most one attribute, depending on the implementation of POSIX options. It is accessed through a variable of type **pthread_condattr_t**. In AIX, the **pthread_condattr_t** data type is a pointer; on other systems, it may be a structure or another data type.

Creating and destroying the condition attributes object

The condition attributes object is initialized to default values by the **pthread_condattr_init** subroutine. The attribute is handled by subroutines. The thread attributes object is destroyed by the **pthread_condattr_destroy** subroutine. This subroutine can release storage dynamically allocated by the **pthread_condattr_init** subroutine, depending on the implementation of the threads library.

In the following example, a condition attributes object is created and initialized with default values, then used and finally destroyed:

```
pthread_condattr_t attributes;
/* the attributes object is created */
...
if (!pthread_condattr_init(&attributes)) {
/* the attributes object is initialized */
...
/* using the attributes object */
...
pthread_condattr_destroy(&attributes);
```

```
        /* the attributes object is destroyed */  
    }
```

The same attributes object can be used to create several condition variables. It can also be modified between two condition variable creations. When the condition variables are created, the attributes object can be destroyed without affecting the condition variables created with it.

Condition attribute

The following condition attribute is supported:

Process-shared

Specifies the process sharing of a condition variable. This attribute depends on the process sharing [POSIX option](#).

Creating and destroying condition variables

A condition variable is created by calling the `pthread_cond_init` subroutine. You may specify a condition attributes object. If you specify a **NULL** pointer, the condition variable will have the default attributes. Thus, the following code fragment:

```
pthread_cond_t cond;  
pthread_condattr_t attr;  
...  
pthread_condattr_init(&attr);  
pthread_cond_init(&cond, &attr);  
pthread_condattr_destroy(&attr);
```

is equivalent to the following:

```
pthread_cond_t cond;  
...  
pthread_cond_init(&cond, NULL);
```

The ID of the created condition variable is returned to the calling thread through the *condition* parameter. The condition ID is an opaque object; its type is `pthread_cond_t`. In AIX, the `pthread_cond_t` data type is a structure; on other systems, it may be a pointer or another data type.

A condition variable must be created once. Avoid calling the `pthread_cond_init` subroutine more than once with the same *condition* parameter (for example, in two threads concurrently executing the same code). Ensuring the uniqueness of a newly created condition variable can be done in the following ways:

- Calling the `pthread_cond_init` subroutine prior to the creation of other threads that will use this variable; for example, in the initial thread.
- Calling the `pthread_cond_init` subroutine within a one-time initialization routine. For more information, see [One-Time Initializations](#).
- Using a static condition variable initialized by the `PTHREAD_COND_INITIALIZER` static initialization macro; the condition variable will have default attributes.

After the condition variable is no longer needed, destroy it by calling the `pthread_cond_destroy` subroutine. This subroutine may reclaim any storage allocated by the `pthread_cond_init` subroutine. After having destroyed a condition variable, the same `pthread_cond_t` variable can be reused to create another condition. For example, the following code fragment is valid, although not very practical:

```
pthread_cond_t cond;  
...  
for (i = 0; i < 10; i++) {  
    /* creates a condition variable */  
    pthread_cond_init(&cond, NULL);  
  
    /* uses the condition variable */  
  
    /* destroys the condition */  
    pthread_cond_destroy(&cond);  
}
```


Like any system resource that can be shared among threads, a condition variable allocated on a thread's stack must be destroyed before the thread is terminated. The threads library maintains a linked list of condition variables; thus, if the stack where a mutex is allocated is freed, the list will be corrupted.

Using condition variables

A condition variable must always be used together with a mutex. A given condition variable can have only one mutex associated with it, but a mutex can be used for more than one condition variable. It is possible to bundle into a structure the condition, the mutex, and the condition variable, as shown in the following code fragment:

```
struct condition_bundle_t {
    int          condition_predicate;
    pthread_mutex_t condition_lock;
    pthread_cond_t condition_variable;
};
```

Waiting for a condition

The mutex protecting the condition must be locked before waiting for the condition. A thread can wait for a condition to be signaled by calling the **pthread_cond_wait** or **pthread_cond_timedwait** subroutine. The subroutine atomically unlocks the mutex and blocks the calling thread until the condition is signaled. When the call returns, the mutex is locked again.

The **pthread_cond_wait** subroutine blocks the thread indefinitely. If the condition is never signaled, the thread never wakes up. Because the **pthread_cond_wait** subroutine provides a cancellation point, the only way to exit this deadlock is to cancel the blocked thread, if cancelability is enabled. For more information, see [Canceling a Thread](#).

The **pthread_cond_timedwait** subroutine blocks the thread only for a given period of time. This subroutine has an extra parameter, *timeout*, specifying an absolute date where the sleep must end. The *timeout* parameter is a pointer to a **timespec** structure. This data type is also called **timestruc_t**. It contains the following fields:

tv_sec

A long unsigned integer, specifying seconds

tv_nsec

A long integer, specifying nanoseconds

Typically, the **pthread_cond_timedwait** subroutine is used in the following manner:

```
struct timespec timeout;
...
time(&timeout.tv_sec);
timeout.tv_sec += MAXIMUM_SLEEP_DURATION;
pthread_cond_timedwait(&cond, &mutex, &timeout);
```

The *timeout* parameter specifies an absolute date. The previous code fragment shows how to specify a duration rather than an absolute date.

To use the **pthread_cond_timedwait** subroutine with an absolute date, you can use the **mktime** subroutine to calculate the value of the *tv_sec* field of the **timespec** structure. In the following example, the thread waits for the condition until 08:00 January 1, 2001, local time:

```
struct tm      date;
time_t        seconds;
struct timespec timeout;
...

date.tm_sec = 0;
date.tm_min = 0;
date.tm_hour = 8;
date.tm_mday = 1;
date.tm_mon = 0;          /* the range is 0-11 */
date.tm_year = 101;      /* 0 is 1900 */
date.tm_wday = 1;        /* this field can be omitted -
                          but it will really be a Monday! */
date.tm_yday = 0;        /* first day of the year */
date.tm_isdst = daylight;
```

```

        /* daylight is an external variable - we are assuming
           that Daylight Saving Time will still be used... */

seconds = mktime(&date);

timeout.tv_sec = (unsigned long)seconds;
timeout.tv_nsec = 0L;

pthread_cond_timedwait(&cond, &mutex, &timeout);

```

The **pthread_cond_timedwait** subroutine also provides a cancellation point, although the sleep is not indefinite. Thus, a sleeping thread can be canceled, whether or not the sleep has a timeout.

Signaling a condition

A condition can be signaled by calling either the **pthread_cond_signal** or the **pthread_cond_broadcast** subroutine.

The **pthread_cond_signal** subroutine wakes up at least one thread that is currently blocked on the specified condition. The awoken thread is chosen according to the scheduling policy; it is the thread with the most-favored scheduling priority (see [Scheduling Policy and Priority](#)). It may happen on multiprocessor systems, or some non-AIX systems, that more than one thread is awakened. Do not assume that this subroutine wakes up exactly one thread.

The **pthread_cond_broadcast** subroutine wakes up every thread that is currently blocked on the specified condition. However, a thread can start waiting on the same condition just after the call to the subroutine returns.

A call to these routines always succeeds, unless an invalid *cond* parameter is specified. This does not mean that a thread has been awakened. Furthermore, signaling a condition is not remembered by the library. For example, consider a condition C. No thread is waiting on this condition. At time *t*, thread 1 signals the condition C. The call is successful although no thread is awakened. At time *t*+1, thread 2 calls the **pthread_cond_wait** subroutine with C as *cond* parameter. Thread 2 is blocked. If no other thread signals C, thread 2 may wait until the process terminates.

You can avoid this kind of deadlock by checking the **EBUSY** error code returned by the **pthread_cond_destroy** subroutine when destroying the condition variable, as in the following code fragment:

The **pthread_yield** subroutine gives the opportunity to another thread to be scheduled; for example, one of the awoken threads. For more information about the **pthread_yield** subroutine.

The **pthread_cond_wait** and the **pthread_cond_broadcast** subroutines must not be used within a signal handler. To provide a convenient way for a thread to await a signal, the threads library provides the **sigwait** subroutine. For more information about the **sigwait** subroutine, see [Signal Management](#).

Synchronizing threads with condition variables

```

while (pthread_cond_destroy(&cond) == EBUSY) {
    pthread_cond_broadcast(&cond);
    pthread_yield();
}

```

Condition variables are used to wait until a particular condition predicate becomes true. This condition predicate is set by another thread, usually the one that signals the condition.

Condition wait semantics

A condition predicate must be protected by a mutex. When waiting for a condition, the wait subroutine (either the **pthread_cond_wait** or **pthread_cond_timedwait** subroutine) atomically unlocks the mutex and blocks the thread. When the condition is signaled, the mutex is relocked and the wait subroutine returns. It is important to note that when the subroutine returns without error, the predicate may still be false.

The reason is that more than one thread may be awoken: either a thread called the **pthread_cond_broadcast** subroutine, or an unavoidable race between two processors simultaneously

woke two threads. The first thread locking the mutex will block all other awoken threads in the wait subroutine until the mutex is unlocked by the program. Thus, the predicate may have changed when the second thread gets the mutex and returns from the wait subroutine.

In general, whenever a condition wait returns, the thread should reevaluate the predicate to determine whether it can safely proceed, should wait again, or should declare a timeout. A return from the wait subroutine does not imply that the predicate is either true or false.

It is recommended that a condition wait be enclosed in a "while loop" that checks the predicate. Basic implementation of a condition wait is shown in the following code fragment:

```
pthread_mutex_lock(&condition_lock);
while (condition_predicate == 0)
    pthread_cond_wait(&condition_variable, &condition_lock);
...
pthread_mutex_unlock(&condition_lock);
```

Timed wait semantics

When the `pthread_cond_timedwait` subroutine returns with the timeout error, the predicate may be true, due to another unavoidable race between the expiration of the timeout and the predicate state change.

Just as for non-timed wait, the thread should reevaluate the predicate when a timeout occurred to determine whether it should declare a timeout or should proceed anyway. It is recommended that you carefully check all possible cases when the `pthread_cond_timedwait` subroutine returns. The following code fragment shows how such checking could be implemented in a robust program:

```
int result = CONTINUE_LOOP;

pthread_mutex_lock(&condition_lock);
while (result == CONTINUE_LOOP) {
    switch (pthread_cond_timedwait(&condition_variable,
        &condition_lock, &timeout)) {

        case 0:
            if (condition_predicate)
                result = PROCEED;
            break;

        case ETIMEDOUT:
            result = condition_predicate ? PROCEED : TIMEOUT;
            break;

        default:
            result = ERROR;
            break;
    }
}

...
pthread_mutex_unlock(&condition_lock);
```

The **result** variable can be used to choose an action. The statements preceding the unlocking of the mutex should be done as soon as possible because a mutex should not be held for long periods of time.

Specifying an absolute date in the *timeout* parameter allows easy implementation of real-time behavior. An absolute timeout need not be recomputed if it is used multiple times in a loop, such as that enclosing a condition wait. For cases where the system clock is advanced discontinuously by an operator, using an absolute timeout ensures that the timed wait will end as soon as the system time specifies a date later than the *timeout* parameter.

Condition variables usage example

The following example provides the source code for a synchronization point routine. A *synchronization point* is a given point in a program where different threads must wait until all threads (or at least a certain number of threads) have reached that point.

A synchronization point can simply be implemented by a counter, which is protected by a lock, and a condition variable. Each thread takes the lock, increments the counter, and waits for the condition to be

signaled if the counter did not reach its maximum. Otherwise, the condition is broadcast, and all threads can proceed. The last thread that calls the routine broadcasts the condition.

```
#define SYNC_MAX_COUNT 10

void SynchronizationPoint()
{
    /* use static variables to ensure initialization */
    static mutex_t sync_lock = PTHREAD_MUTEX_INITIALIZER;
    static cond_t sync_cond = PTHREAD_COND_INITIALIZER;
    static int sync_count = 0;

    /* lock the access to the count */
    pthread_mutex_lock(&sync_lock);

    /* increment the counter */
    sync_count++;

    /* check if we should wait or not */
    if (sync_count < SYNC_MAX_COUNT)

        /* wait for the others */
        pthread_cond_wait(&sync_cond, &sync_lock);

    else

        /* broadcast that everybody reached the point */
        pthread_cond_broadcast(&sync_cond);

    /* unlocks the mutex - otherwise only one thread
       will be able to return from the routine! */
    pthread_mutex_unlock(&sync_lock);
}
```

This routine has some limitations: it can be used only once, and the number of threads that call the routine is coded by a symbolic constant. However, this example shows a basic usage of condition variables. For more complex usage examples, see [Creating Complex Synchronization Objects](#).

Related concepts

[Joining threads](#)

Joining a thread means waiting for it to terminate, which can be seen as a specific usage of condition variables.

Using read/write locks

In many situations, data is read more often than it is modified or written.

In these cases, you can allow threads to read concurrently while holding the lock and allow only one thread to hold the lock when data is modified. A multiple-reader single-writer lock (or read/write lock) does this. A read/write lock is acquired either for reading or writing, and then is released. The thread that acquires the read-write lock must be the one that releases it.

Read/write attributes object

The `pthread_rwlockattr_init` subroutine initializes a read-write lock attributes object (**attr**). The default value for all attributes is defined by the implementation. Unexpected results can occur if the `pthread_rwlockattr_init` subroutine specifies an already-initialized read/write lock attributes object.

The following examples illustrate how to call the `pthread_rwlockattr_init` subroutine with the **attr** object:

```
pthread_rwlockattr_t attr;
```

and:

```
pthread_rwlockattr_init(&attr);
```

After a read/write lock attributes object is used to initialize one or more read/write locks, any function that affects the attributes object (including destruction) does not affect any previously initialized read/write locks.

The **pthread_rwlockattr_destroy** subroutine destroys a read/write lock attributes object. Unexpected results can occur if the object is used before it is reinitialized by another call to the **pthread_rwlockattr_init** subroutine. An implementation can cause the **pthread_rwlockattr_destroy** subroutine to set the object referenced by the **attr** object to an invalid value.

Creating and Destroying Read/Write Locks

The **pthread_rwlock_init** subroutine initializes the read/write lock referenced by the **rwlock** object with the attributes referenced by the **attr** object. If the **attr** object is NULL, the default read/write lock attributes are used; the effect is the same as passing the address of a default read-write lock attributes object. Upon successful initialization, the state of the read/write lock becomes initialized and unlocked. After initialized, the lock can be used any number of times without being reinitialized. Unexpected results can occur if the call to the **pthread_rwlock_init** subroutine is called specifying an already initialized read/write lock, or if a read/write lock is used without first being initialized.

If the **pthread_rwlock_init** subroutine fails, the **rwlock** object is not initialized and the contents are undefined.

The **pthread_rwlock_destroy** subroutine destroys the read/write lock object referenced by the **rwlock** object and releases any resources used by the lock. Unexpected results can occur in any of the following situations:

- If the lock is used before it is reinitialized by another call to the **pthread_rwlock_init** subroutine.
- An implementation can cause the **pthread_rwlock_destroy** subroutine to set the object referenced by the **rwlock** object to an invalid value. Unexpected results can occur if **pthread_rwlock_destroy** is called when any thread holds the **rwlock** object.
- Attempting to destroy an uninitialized read/write lock results in unexpected results. A destroyed read/write lock object can be reinitialized by using the **pthread_rwlock_init** subroutine. Unexpected results can occur if the read/write lock object is referenced after it has been destroyed.

In cases where default read/write lock attributes are appropriate, use the **PTHREAD_RWLOCK_INITIALIZER** macro to initialize read/write locks that are statically allocated. For example:

```
pthread_rwlock_t      rwlock1 = PTHREAD_RWLOCK_INITIALIZER;
```

The effect is similar to dynamic initialization by using a call to the **pthread_rwlock_init** subroutine with the *attr* parameter specified as NULL, except that no error checks are performed. For example:

```
pthread_rwlock_init(&rwlock2, NULL);
```

The following example illustrates how to use the **pthread_rwlock_init** subroutine with the *attr* parameter initialized. For an example of how to initialize the *attr* parameter, see [Read-Write Attributes Object](#).

```
pthread_rwlock_init(&rwlock, &attr);
```

Locking a read/write lock object for reading

The **pthread_rwlock_rdlock** subroutine applies a read lock to the read/write lock referenced by the **rwlock** object. The calling thread acquires the read lock if a writer does not hold the lock and if no writers are blocked on the lock. It is unspecified whether the calling thread acquires the lock when a writer does not hold the lock and there are writers waiting for the lock. If a writer holds the lock, the calling thread will not acquire the read lock. If the read lock is not acquired, the calling thread does not return from the **pthread_rwlock_rdlock** call until it can acquire the lock. Results are undefined if the calling thread holds a write lock on the **rwlock** object at the time the call is made.

A thread can hold multiple concurrent read locks on the **rwlock** object (that is, successfully call the **pthread_rwlock_rdlock** subroutine *n* times). If so, the thread must perform matching unlocks (that is, it must call the **pthread_rwlock_unlock** subroutine *n* times).

The **pthread_rwlock_tryrdlock** subroutine applies a read lock similar to the **pthread_rwlock_rdlock** subroutine with the exception that the subroutine fails if any thread holds a write lock on the **rwlock** object or if writers are blocked on the **rwlock** object. Results are undefined if any of these functions are called with an uninitialized read/write lock.

If a signal is delivered to a thread waiting for a read/write lock for reading, upon return from the signal handler, the thread resumes waiting for the read/write lock for reading as if it was not interrupted.

Locking a read/write lock object for writing

The **pthread_rwlock_wrlock** subroutine applies a write lock to the read/write lock referenced by the **rwlock** object. The calling thread acquires the write lock if no other thread (reader or writer) holds the read/write lock on the **rwlock** object. Otherwise, the thread does not return from the **pthread_rwlock_wrlock** call until it can acquire the lock. Results are undefined if the calling thread holds the read/write lock (whether a read or write lock) at the time the call is made.

The **pthread_rwlock_trywrlock** subroutine applies a write lock similar to the **pthread_rwlock_wrlock** subroutine, with the exception that the function fails if any thread currently holds **rwlock** for reading or writing. Results are undefined if any of these functions are called with an uninitialized read/write lock.

If a signal is delivered to a thread waiting for a read/write lock for writing, upon return from the signal handler, the thread resumes waiting for the read/write lock for writing as if it was not interrupted.

Preferring a writer-thread over a reader-thread

The **pthread_rwlock_attr_setfavorwriters_np** subroutine can be used by an application to initialize the attributes of a read/write lock. You can specify the pthread library to prioritize the scheduling of the threads that requires the read/write lock in write mode. When the pthread library schedules the writer-threads (threads that write data) to get the read/write lock in write mode, the pthread library does not support recursion by threads that are holding a the read/write lock in read mode. Unexpected results can occur when the threads hold a read/write lock in the read mode more than once.

The **pthread_rwlock_attr_getfavorwriters_np** subroutine returns the current preference that is set in the read/write lock attribute structure. By default, the reader-threads are prioritized over writer-threads to get a read/write lock.

Sample read/write lock programs

The following sample programs demonstrate how to use locking subroutines. To run these programs, you need the **check.h** file and **makefile**.

check.h file:

```
#include stdio.h
#include stdio.h
#include stdio.h
#include stdio.h

/* Simple function to check the return code and exit the program
   if the function call failed
*/
static void compResults(char *string, int rc) {
    if (rc) {
        printf("Error on : %s, rc=%d",
              string, rc);
        exit(EXIT_FAILURE);
    }
    return;
}
```

Make file:

```
CC_R = xlc_r
TARGETS = test01 test02 test03
OBJS = test01.o test02.o test03.o
SRCS = $(OBJS:.o=.c)
$(TARGETS): $(OBJS)
    $(CC_R) -o $$@ $$@.o
clean:
    rm $(OBJS) $(TARGETS)
run:
    test01
    test02
    test03
```

Single-thread example

The following example uses the `pthread_rwlock_tryrdlock` subroutine with a single thread. For an example of using the `pthread_rwlock_tryrdlock` subroutine with multiple threads, see [Multiple-Thread Example](#).

```
Example: test01.c
#define _MULTI_THREADED
#include pthread.h
#include stdio.h
#include "check.h"

pthread_rwlock_t      rwlock = PTHREAD_RWLOCK_INITIALIZER;

void *rdlockThread(void *arg)
{
    int          rc;
    int          count=0;

    printf("Entered thread, getting read lock with mp wait\n");
    Retry:
    rc = pthread_rwlock_tryrdlock(&rwlock);
    if (rc == EBUSY) {
        if (count >= 10) {
            printf("Retried too many times, failure!\n");
            exit(EXIT_FAILURE);
        }
        ++count;
        printf("Could not get lock, do other work, then RETRY...\n");
        sleep(1);
        goto Retry;
    }
    compResults("pthread_rwlock_tryrdlock() 1\n", rc);

    sleep(2);

    printf("unlock the read lock\n");
    rc = pthread_rwlock_unlock(&rwlock);
    compResults("pthread_rwlock_unlock()\n", rc);

    printf("Secondary thread complete\n");
    return NULL;
}

int main(int argc, char **argv)
{
    int          rc=0;
    pthread_t    thread;

    printf("Enter test case - %s\n", argv[0]);

    printf("Main, get the write lock\n");
    rc = pthread_rwlock_wrlock(&rwlock);
    compResults("pthread_rwlock_wrlock()\n", rc);
```

```

printf("Main, create the try read lock thread\n");
rc = pthread_create(&thread, NULL, rdlockThread, NULL);
compResults("pthread_create\n", rc);

printf("Main, wait a bit holding the write lock\n");
sleep(5);

printf("Main, Now unlock the write lock\n");
rc = pthread_rwlock_unlock(&rwlock);
compResults("pthread_rwlock_unlock()\n", rc);

printf("Main, wait for the thread to end\n");
rc = pthread_join(thread, NULL);
compResults("pthread_join\n", rc);

rc = pthread_rwlock_destroy(&rwlock);
compResults("pthread_rwlock_destroy()\n", rc);
printf("Main completed\n");
return 0;
}

```

The output for this sample program will be similar to the following:

```

Enter test case - ./test01
Main, get the write lock
Main, create the try read lock thread
Main, wait a bit holding the write lock

Entered thread, getting read lock with mp wait
Could not get lock, do other work, then RETRY...
Could not get lock, do other work, then RETRY...
Could not get lock, do other work, then RETRY...
Could not get lock, do other work, then RETRY...
Could not get lock, do other work, then RETRY...
Main, Now unlock the write lock
Main, wait for the thread to end
unlock the read lock
Secondary thread complete
Main completed

```

Multiple-thread example

The following example uses the **pthread_rwlock_tryrdlock** subroutine with multiple threads. For an example of using the **pthread_rwlock_tryrdlock** subroutine with a single thread, see [Single-Thread Example](#).

```

Example: test02.c

#define _MULTI_THREADED
#include pthread.h
#include stdio.h
#include "check.h"

pthread_rwlock_t      rwlock = PTHREAD_RWLOCK_INITIALIZER;

void *wrlockThread(void *arg)
{
    int          rc;
    int          count=0;

    printf("%.8x: Entered thread, getting write lock\n",
           pthread_self());
    Retry:
    rc = pthread_rwlock_trywrlock(&rwlock);
    if (rc == EBUSY) {
        if (count >= 10) {
            printf("%.8x: Retried too many times, failure!\n",
                   pthread_self());
            exit(EXIT_FAILURE);
        }
        ++count;
        printf("%.8x: Go off an do other work, then RETRY...\n",
               pthread_self());
        sleep(1);
        goto Retry;
    }
}

```



```

}
compResults("pthread_rwlock_trywrlock() 1\n", rc);
printf("%.8x: Got the write lock\n", pthread_self());

sleep(2);

printf("%.8x: Unlock the write lock\n",
        pthread_self());
rc = pthread_rwlock_unlock(&rwlock);
compResults("pthread_rwlock_unlock()\n", rc);

printf("%.8x: Secondary thread complete\n",
        pthread_self());
return NULL;
}

int main(int argc, char **argv)
{
    int                rc=0;
    pthread_t          thread, thread2;

    printf("Enter test case - %s\n", argv[0]);

    printf("Main, get the write lock\n");
    rc = pthread_rwlock_wrlock(&rwlock);
    compResults("pthread_rwlock_wrlock()\n", rc);

    printf("Main, create the timed write lock threads\n");
    rc = pthread_create(&thread, NULL, wrlockThread, NULL);
    compResults("pthread_create\n", rc);

    rc = pthread_create(&thread2, NULL, wrlockThread, NULL);
    compResults("pthread_create\n", rc);

    printf("Main, wait a bit holding this write lock\n");
    sleep(1);

    printf("Main, Now unlock the write lock\n");
    rc = pthread_rwlock_unlock(&rwlock);
    compResults("pthread_rwlock_unlock()\n", rc);

    printf("Main, wait for the threads to end\n");
    rc = pthread_join(thread, NULL);
    compResults("pthread_join\n", rc);

    rc = pthread_join(thread2, NULL);
    compResults("pthread_join\n", rc);

    rc = pthread_rwlock_destroy(&rwlock);
    compResults("pthread_rwlock_destroy()\n", rc);
    printf("Main completed\n");
    return 0;
}

```

The output for this sample program will be similar to the following:

```

Enter test case - ./test02
Main, get the write lock
Main, create the timed write lock threads
Main, wait a bit holding this write lock
00000102: Entered thread, getting write lock
00000102: Go off an do other work, then RETRY...
00000203: Entered thread, getting write lock
00000203: Go off an do other work, then RETRY...
Main, Now unlock the write lock
Main, wait for the threads to end
00000102: Got the write lock
00000203: Go off an do other work, then RETRY...
00000203: Go off an do other work, then RETRY...
00000102: Unlock the write lock
00000102: Secondary thread complete
00000203: Got the write lock
00000203: Unlock the write lock
00000203: Secondary thread complete
Main completed

```

Read/write read-lock example

The following example uses the `pthread_rwlock_rdlock` subroutine to implement read/write read locks:

```
Example: test03.c

#define _MULTI_THREADED
#include pthread.h
#include stdio.h
#include "check.h"

pthread_rwlock_t      rwlock;

void *rdlockThread(void *arg)
{
    int rc;

    printf("Entered thread, getting read lock\n");
    rc = pthread_rwlock_rdlock(&rwlock);
    compResults("pthread_rwlock_rdlock()\n", rc);
    printf("got the rwlock read lock\n");

    sleep(5);

    printf("unlock the read lock\n");
    rc = pthread_rwlock_unlock(&rwlock);
    compResults("pthread_rwlock_unlock()\n", rc);
    printf("Secondary thread unlocked\n");
    return NULL;
}

void *wrlockThread(void *arg)
{
    int rc;

    printf("Entered thread, getting write lock\n");
    rc = pthread_rwlock_wrlock(&rwlock);
    compResults("pthread_rwlock_wrlock()\n", rc);

    printf("Got the rwlock write lock, now unlock\n");
    rc = pthread_rwlock_unlock(&rwlock);
    compResults("pthread_rwlock_unlock()\n", rc);
    printf("Secondary thread unlocked\n");
    return NULL;
}

int main(int argc, char **argv)
{
    int          rc=0;
    pthread_t    thread, thread1;

    printf("Enter test case - %s\n", argv[0]);

    printf("Main, initialize the read write lock\n");
    rc = pthread_rwlock_init(&rwlock, NULL);
    compResults("pthread_rwlock_init()\n", rc);

    printf("Main, grab a read lock\n");
    rc = pthread_rwlock_rdlock(&rwlock);
    compResults("pthread_rwlock_rdlock()\n",rc);

    printf("Main, grab the same read lock again\n");
    rc = pthread_rwlock_rdlock(&rwlock);
    compResults("pthread_rwlock_rdlock() second\n", rc);

    printf("Main, create the read lock thread\n");
    rc = pthread_create(&thread, NULL, rdlockThread, NULL);
    compResults("pthread_create\n", rc);

    printf("Main - unlock the first read lock\n");
    rc = pthread_rwlock_unlock(&rwlock);
    compResults("pthread_rwlock_unlock()\n", rc);

    printf("Main, create the write lock thread\n");
    rc = pthread_create(&thread1, NULL, wrlockThread, NULL);
    compResults("pthread_create\n", rc);

    sleep(5);
}
```

```

printf("Main - unlock the second read lock\n");
rc = pthread_rwlock_unlock(&rwlock);
compResults("pthread_rwlock_unlock()\n", rc);

printf("Main, wait for the threads\n");
rc = pthread_join(thread, NULL);
compResults("pthread_join\n", rc);

rc = pthread_join(thread1, NULL);
compResults("pthread_join\n", rc);

rc = pthread_rwlock_destroy(&rwlock);
compResults("pthread_rwlock_destroy()\n", rc);

printf("Main completed\n");
return 0;
}

```

The output for this sample program will be similar to the following:

```

$ ./test03
Enter test case - ./test03
Main, initialize the read write lock
Main, grab a read lock
Main, grab the same read lock again
Main, create the read lock thread
Main - unlock the first read lock
Main, create the write lock thread
Entered thread, getting read lock
got the rwlock read lock
Entered thread, getting write lock
Main - unlock the second read lock
Main, wait for the threads
unlock the read lock
Secondary thread unlocked
Got the rwlock write lock, now unlock
Secondary thread unlocked
Main completed

```

Joining threads

Joining a thread means waiting for it to terminate, which can be seen as a specific usage of condition variables.

Waiting for a thread

Using the **pthread_join** subroutine allows a thread to wait for another thread to terminate. More complex conditions, such as waiting for multiple threads to terminate, can be implemented by using condition variables.

Calling the pthread_join subroutine

The **pthread_join** subroutine blocks the calling thread until the specified thread terminates. The target thread (the thread whose termination is awaited) must not be detached. If the target thread is already terminated, but not detached, the **pthread_join** subroutine returns immediately. After a target thread has been joined, it is automatically detached, and its storage can be reclaimed.

The following table indicates the possible cases when a thread calls the **pthread_join** subroutine, depending on the **state** and the **detachstate** attribute of the target thread.

Target State	Undetached target	Detached target
Target is still running	The caller is blocked until the target is terminated.	The call returns immediately, indicating an error.
Target is terminated	The call returns immediately, indicating a successful completion.	

Multiple joins

Several threads can join the same target thread, if the target is not detached. The success of this operation depends on the order of the calls to the **pthread_join** subroutine and the moment when the target thread terminates.

- Any call to the **pthread_join** subroutine occurring before the target thread's termination blocks the calling thread.
- When the target thread terminates, all blocked threads are awoken, and the target thread is automatically detached.
- Any call to the **pthread_join** subroutine occurring after the target thread's termination will fail, because the thread is detached by the previous join.
- If no thread called the **pthread_join** subroutine before the target thread's termination, the first call to the **pthread_join** subroutine will return immediately, indicating a successful completion, and any further call will fail.

Join example

In the following example, the program ends after exactly five messages display in each language. This is done by blocking the initial thread until the "writer" threads exit.

```
#include <pthread.h>    /* include file for pthreads - the 1st */
#include <stdio.h>      /* include file for printf()          */

void *Thread(void *string)
{
    int i;

    /* writes five messages and exits */
    for (i=0; i<5; i++)
        printf("%s\n", (char *)string);
    pthread_exit(NULL);
}

int main()
{
    char *e_str = "Hello!";
    char *f_str = "Bonjour !";

    pthread_attr_t attr;
    pthread_t e_th;
    pthread_t f_th;

    int rc;

    /* creates the right attribute */
    pthread_attr_init(&attr);
    pthread_attr_setdetachstate(&attr,
                                PTHREAD_CREATE_UNDETACHED);

    /* creates both threads */
    rc = pthread_create(&e_th, &attr, Thread, (void *)e_str);
    if (rc)
        exit(-1);
    rc = pthread_create(&f_th, &attr, Thread, (void *)f_str);
    if (rc)
        exit(-1);
    pthread_attr_destroy(&attr);

    /* joins the threads */
    pthread_join(e_th, NULL);
    pthread_join(f_th, NULL);

    pthread_exit(NULL);
}
```

A thread cannot join itself because a deadlock would occur and it is detected by the library. However, two threads may try to join each other. They will deadlock, but this situation is not detected by the library.

Returning information from a thread

The **pthread_join** subroutine also allows a thread to return information to another thread. When a thread calls the **pthread_exit** subroutine or when it returns from its entry-point routine, it returns a pointer (see

Exiting a Thread). This pointer is stored as long as the thread is not detached, and the **pthread_join** subroutine can return it.

For example, a multithreaded **grep** command may choose the implementation in the following example. In this example, the initial thread creates one thread per file to scan, each thread having the same entry point routine. The initial thread then waits for all threads to be terminated. Each "scanning" thread stores the found lines in a dynamically allocated buffer and returns a pointer to this buffer. The initial thread prints each buffer and releases it.

```
/* "scanning" thread */
...
buffer = malloc(...);
/* finds the search pattern in the file
   and stores the lines in the buffer */
return (buffer);

/* initial thread */
...
for (/* each created thread */) {
    void *buf;
    pthread_join(thread, &buf);
    if (buf != NULL) {
        /* print all the lines in the buffer,
           preceded by the filename of the thread */
        free(buf);
    }
}
...
```

If the target thread is canceled, the **pthread_join** subroutine returns a value of -1 cast into a pointer (see [Canceling a Thread](#)). Because -1 cannot be a pointer value, getting -1 as returned pointer from a thread means that the thread was canceled.

The returned pointer can point to any kind of data. The pointer must still be valid after the thread was terminated and its storage reclaimed. Therefore, avoid returning a value, because the destructor routine is called when the thread's storage is reclaimed.

Returning a pointer to dynamically allocated storage to several threads needs special consideration. Consider the following code fragment:

```
void *returned_data;
...
pthread_join(target_thread, &returned_data);
/* retrieves information from returned_data */
free(returned_data);
```

The **returned_data** pointer is freed when it is executed by only one thread. If several threads execute the above code fragment concurrently, the **returned_data** pointer is freed several times; a situation that must be avoided. To prevent this, use a mutex-protected flag to signal that the **returned_data** pointer was freed. The following line from the previous example:

```
free(returned_data);
```

would be replaced by the following lines, where a mutex can be used for locking the access to the critical region (assuming the *flag* variable is initially 0):

```
/* lock - entering a critical region, no other thread should
   run this portion of code concurrently */
if (!flag) {
    free(returned_data);
    flag = 1;
}
/* unlock - exiting the critical region */
```

Locking access to the critical region ensures that the **returned_data** pointer is freed only once.

When returning a pointer to dynamically allocated storage to several threads all executing different code, you must ensure that exactly one thread frees the pointer.

Related concepts

Using condition variables

Condition variables allow threads to wait until some event or condition has occurred.

Thread-specific data

Many applications require that certain data be maintained on a per-thread basis across function calls.

Scheduling threads

Threads can be scheduled, and the threads library provides several facilities to handle and control the scheduling of threads.

It also provides facilities to control the scheduling of threads during synchronization operations such as locking a mutex. Each thread has its own set of scheduling parameters. These parameters can be set using the thread attributes object before the thread is created. The parameters can also be dynamically set during the thread's execution.

Controlling the scheduling of a thread can be a complicated task. Because the scheduler handles all threads system wide, the scheduling parameters of a thread interact with those of all other threads in the process and in the other processes. The following facilities are the first to be used if you want to control the scheduling of a thread.

The threads library allows the programmer to control the execution scheduling of the threads in the following ways:

- By setting scheduling attributes when creating a thread
- By dynamically changing the scheduling attributes of a created thread
- By defining the effect of a mutex on the thread's scheduling when creating a mutex (known as *synchronization scheduling*)
- By dynamically changing the scheduling of a thread during synchronization operations (known as *synchronization scheduling*)

Scheduling parameters

A thread has the following scheduling parameters:

Parameter	Description
scope	The contention scope of a thread is defined by the thread model used in the threads library.
policy	The scheduling policy of a thread defines how the scheduler treats the thread after it gains control of the CPU.
priority	The scheduling priority of a thread defines the relative importance of the work being done by each thread.

The scheduling parameters can be set before the thread's creation or during the thread's execution. In general, controlling the scheduling parameters of threads is important only for threads that are CPU-intensive. Thus, the threads library provides default values that are sufficient for most cases.

Using the inheritsched attribute

The **inheritsched** attribute of the thread attributes object specifies how the thread's scheduling attributes will be defined. The following values are valid:

Values	Description
PTHREAD_INHERIT_SCHED	Specifies that the new thread will get the scheduling attributes (schedpolicy and schedparam attributes) of its creating thread. Scheduling attributes defined in the attributes object are ignored.
PTHREAD_EXPLICIT_SCHED	Specifies that the new thread will get the scheduling attributes defined in this attributes object.

The default value of the **inheritsched** attribute is `PTHREAD_INHERIT_SCHED`. The attribute is set by calling the **pthread_attr_setinheritsched** subroutine. The current value of the attribute is returned by calling the **pthread_attr_getinheritsched** subroutine.

To set the scheduling attributes of a thread in the thread attributes object, the **inheritsched** attribute must first be set to `PTHREAD_EXPLICIT_SCHED`. Otherwise, the attributes-object scheduling attributes are ignored.

Scheduling policy and priority

The threads library provides the following scheduling policies:

Library	Description
SCHED_FIFO	First-in first-out (FIFO) scheduling. Each thread has a fixed priority; when multiple threads have the same priority level, they run to completion in FIFO order.
SCHED_RR	Round-robin (RR) scheduling. Each thread has a fixed priority; when multiple threads have the same priority level, they run for a fixed time slice in FIFO order.
SCHED_OTHER	Default AIX scheduling. Each thread has an initial priority that is dynamically modified by the scheduler, according to the thread's activity; thread execution is time-sliced. On other systems, this scheduling policy may be different.

In versions of AIX prior to 5.3, changing the priority of a thread when setting its scheduling policy to **SCHED_OTHER** is not permitted. In this case, the kernel directly manages the priority, and the only legal value that can be passed to the **pthread_setschedparam** subroutine is the **DEFAULT_PRIO** value. The **DEFAULT_PRIO** value is defined in **pthread.h** file as 1, and any other passed values are ignored.

Beginning with AIX 5.3, you can change the priority of a thread when you set its scheduling policy to **SCHED_OTHER**. The legal values that can be passed to the **pthread_setschedparam** subroutine are from 40 to 80, however, only privileged users can set a priority greater than 60. A priority in the range of 1 to 39 provides the same priority as that of 40, and a priority in the range of 81 to 127 provides the same priority as that of 80.

Note: In AIX, the kernel inverts the priority levels. For the AIX kernel, the priority is in the range from 0 to 127, where 0 is the most favored priority and 127 the least favored priority. Commands, such as the **ps** command, report the kernel priority.

The threads library handles the priority through a **sched_param** structure, defined in the **sys/sched.h** header file. This structure contains the following fields:

Fields	Description
<code>sched_priority</code>	Specifies the priority.
<code>sched_policy</code>	This field is ignored by the threads library. Do not use.

Setting the scheduling policy and priority at creation time

The scheduling policy can be set when creating a thread by setting the **schedpolicy** attribute of the thread attributes object. The **pthread_attr_setschedpolicy** subroutine sets the scheduling policy to one of the previously defined scheduling policies. The current value of the **schedpolicy** attribute of a thread attributes object can be obtained by using the **pthread_attr_getschedpolicy** subroutine.

The scheduling priority can be set at creation time of a thread by setting the **schedparam** attribute of the thread attributes object. The **pthread_attr_setschedparam** subroutine sets the value of the **schedparam** attribute, copying the value of the specified structure. The **pthread_attr_getschedparam** subroutine gets the **schedparam** attribute.

In the following code fragment, a thread is created with the round-robin scheduling policy, using a priority level of 3:

```
sched_param schedparam;  
schedparam.sched_priority = 3;
```

```
pthread_attr_init(&attr);
pthread_attr_setinheritsched(&attr, PTHREAD_EXPLICIT_SCHED);
pthread_attr_setschedpolicy(&attr, SCHED_RR);
pthread_attr_setschedparam(&attr, &schedparam);

pthread_create(&thread, &attr, &start_routine, &args);
pthread_attr_destroy(&attr);
```

For more information about the **inheritsched** attribute, see [Using the inheritsched Attribute](#).

Setting the scheduling attributes at execution time

The **pthread_getschedparam** subroutine returns the **schedpolicy** and **schedparam** attributes of a thread. These attributes can be set by calling the **pthread_setschedparam** subroutine. If the target thread is currently running on a processor, the new scheduling policy and priority will be implemented the next time the thread is scheduled. If the target thread is not running, it can be scheduled immediately at the end of the subroutine call.

For example, consider a thread T that is currently running with round-robin policy at the moment the **schedpolicy** attribute of T is changed to FIFO. T will run until the end of its time slice, at which time its scheduling attributes are then re-evaluated. If no threads have higher priority, T will be rescheduled, even before other threads having the same priority. Consider a second example where a low-priority thread is not running. If this thread's priority is raised by another thread calling the **pthread_setschedparam** subroutine, the target thread will be scheduled immediately if it is the highest priority runnable thread.

Note: Both subroutines use a *policy* parameter and a **sched_param** structure. Although this structure contains a `sched_policy` field, programs should not use it. The subroutines use the *policy* parameter to pass the scheduling policy, and the subroutines then ignore the `sched_policy` field.

Scheduling-policy considerations

Applications should use the default scheduling policy, unless a specific application requires the use of a fixed-priority scheduling policy. Consider the following points about using the nondefault policies:

- Using the round-robin policy ensures that all threads having the same priority level will be scheduled equally, regardless of their activity. This can be useful in programs where threads must read sensors or write actuators.
- Using the FIFO policy should be done with great care. A thread running with FIFO policy runs to completion, unless it is blocked by some calls, such as performing input and output operations. A high-priority FIFO thread may not be preempted and can affect the global performance of the system. For example, threads doing intensive calculations, such as inverting a large matrix, should never run with FIFO policy.

The setting of scheduling policy and priority is also influenced by the contention scope of threads. Using the FIFO or the round-robin policy may not always be allowed.

sched_yield subroutine

The **sched_yield** subroutine is the equivalent for threads of the **yield** subroutine. The **sched_yield** subroutine forces the calling thread to relinquish the use of its processor and gives other threads an opportunity to be scheduled. The next scheduled thread may belong to the same process as the calling thread or to another process. Do not use the **yield** subroutine in a multithreaded program.

The interface **pthread_yield** subroutine is not available in Single UNIX Specification, Version 2.

Related concepts

[Synchronization scheduling](#)

Programmers can control the execution scheduling of threads when there are constraints, especially time constraints, that require certain threads to be executed faster than other ones.

[List of scheduling subroutines](#)

This section lists scheduling subroutines.

[Developing multithreaded programs](#)

Developing multithreaded programs is similar to developing programs with multiple processes. Developing programs also consists of compiling and debugging the code.

List of scheduling subroutines

This section lists scheduling subroutines.

Subroutine	Description
<u>pthread_attr_getschedparam</u>	Returns the value of the schedparam attribute of a thread attributes object.
<u>pthread_attr_setschedparam</u>	Sets the value of the schedparam attribute of a thread attributes object.
<u>pthread_getschedparam</u>	Returns the value of the schedpolicy and schedparam attributes of a thread.
<u>sched_yield</u>	Forces the calling thread to relinquish use of its processor.

Related concepts

[Scheduling threads](#)

Threads can be scheduled, and the threads library provides several facilities to handle and control the scheduling of threads.

Contention scope and concurrency level

The *contention scope* of a user thread defines how it is mapped to a kernel thread

. The threads library defines the following contention scopes:

PTHREAD_SCOPE_PROCESS

Process contention scope, sometimes called *local contention scope*. Specifies that the thread will be scheduled against all other local contention scope threads in the process. A process-contention-scope user thread is a user thread that shares a kernel thread with other process-contention-scope user threads in the process. All user threads in an M:1 thread model have process contention scope.

PTHREAD_SCOPE_SYSTEM

System contention scope, sometimes called *global contention scope*. Specifies that the thread will be scheduled against all other threads in the system and is directly mapped to one kernel thread. All user threads in a 1:1 thread model have system contention scope.

In an M:N thread model, user threads can have either system or process contention scope. Therefore, an M:N thread model is often referred as a *mixed-scope* model.

The *concurrency level* is a property of M:N threads libraries. It defines the number of virtual processors used to run the process-contention scope user threads. This number cannot exceed the number of process-contention-scope user threads and is usually dynamically set by the threads library. The system also sets a limit to the number of available kernel threads.

Setting the contention scope

The contention scope can only be set before a thread is created by setting the contention-scope attribute of a thread attributes object. The **pthread_attr_setscope** subroutine sets the value of the attribute; the **pthread_attr_getscope** returns it.

The contention scope is only meaningful in a mixed-scope M:N library implementation. A **TestImplementation** routine could be written as follows:

```
int TestImplementation()
{
    pthread_attr_t a;
    int result;

    pthread_attr_init(&a);
    switch (pthread_attr_setscope(&a, PTHREAD_SCOPE_PROCESS))
    {
```

```

        case 0:      result = LIB_MN; break;
        case ENOTSUP: result = LIB_11; break;
        case ENOSYS: result = NO_PRIO_OPTION; break;
        default:    result = ERROR; break;
    }

    pthread_attr_destroy(&a);
    return result;
}

```

Impacts of contention scope on scheduling

The contention scope of a thread influences its scheduling. Each contention-scope thread is bound to one kernel thread. Thus, changing the scheduling policy and priority of a global user thread results in changing the scheduling policy and priority of the underlying kernel thread.

In AIX, only kernel threads with root authority can use a fixed-priority scheduling policy (FIFO or round-robin). The following code will always return the **EPERM** error code if the calling thread has system contention scope but does not have root authority. This code would not fail, if the calling thread had process contention scope.

```

    schedparam.sched_priority = 3;
    pthread_setschedparam(pthread_self(), SCHED_FIFO, schedparam);

```

Note: Root authority is not required to control the scheduling parameters of user threads having process contention scope.

Local user threads can set any scheduling policy and priority, within the valid range of values. However, two threads having the same scheduling policy and priority but having different contention scope will not be scheduled in the same way. Threads having process contention scope are executed by kernel threads whose scheduling parameters are set by the library.

Related concepts

[Understanding threads and processes](#)

A *thread* is an independent flow of control that operates within the same address space as other independent flows of controls within a process.

Synchronization scheduling

Programmers can control the execution scheduling of threads when there are constraints, especially time constraints, that require certain threads to be executed faster than other ones.

Synchronization objects, such as mutexes, may block even high-priority threads. In some cases, undesirable behavior, known as *priority inversion*, may occur. The threads library provides the *mutex protocols* facility to avoid priority inversions.

Synchronization scheduling defines how the execution scheduling, especially the priority, of a thread is modified by holding a mutex. This allows custom-defined behavior and avoids priority inversions. It is useful when using complex locking schemes. Some implementations of the threads library do not provide synchronization scheduling.

Priority inversion

Priority inversion occurs when a low-priority thread holds a mutex, blocking a high-priority thread. Due to its low priority, the mutex owner may hold the mutex for an unbounded duration. As a result, it becomes impossible to guarantee thread deadlines.

The following example illustrates a typical priority inversion. In this example, the case of a uniprocessor system is considered. Priority inversions also occur on multiprocessor systems in a similar way.

In our example, a mutex *M* is used to protect some common data. Thread *A* has a priority level of 100 and is scheduled very often. Thread *B* has a priority level of 20 and is a background thread. Other threads in the process have priority levels near 60. A code fragment from thread *A* is as follows:

```
pthread_mutex_lock(&M);           /* 1 */
...
pthread_mutex_unlock(&M);
```

A code fragment from thread *B* is as follows:

```
pthread_mutex_lock(&M);           /* 2 */
...
fprintf(...);                    /* 3 */
...
pthread_mutex_unlock(&M);
```

Consider the following execution chronology. Thread *B* is scheduled and executes line 2. While executing line 3, thread *B* is preempted by thread *A*. Thread *A* executes line 1 and is blocked, because the mutex *M* is held by thread *B*. Thus, other threads in the process are scheduled. Because thread *B* has a very low priority, it may not be rescheduled for a long period, blocking thread *A*, although thread *A* has a very high priority.

Mutex protocols

To avoid priority inversions, the following mutex protocols are provided by the threads library:

Priority inheritance protocol

Sometimes called *basic priority inheritance protocol*. In the priority inheritance protocol, the mutex holder inherits the priority of the highest-priority blocked thread. When a thread tries to lock a mutex using this protocol and is blocked, the mutex owner temporarily receives the blocked thread's priority, if that priority is higher than the owner's. It recovers its original priority when it unlocks the mutex.

Priority protection protocol

Sometimes called *priority ceiling protocol emulation*. In the priority protection protocol, each mutex has a *priority ceiling*. It is a priority level within the valid range of priorities. When a thread owns a mutex, it temporarily receives the mutex priority ceiling, if the ceiling is higher than its own priority. It recovers its original priority when it unlocks the mutex. The priority ceiling should have the value of the highest priority of all threads that may lock the mutex. Otherwise, priority inversions or even deadlocks may occur, and the protocol would be inefficient.

Both protocols increase the priority of a thread holding a specific mutex, so that deadlines can be guaranteed. Furthermore, when correctly used, mutex protocols can prevent mutual deadlocks. Mutex protocols are individually assigned to mutexes.

Choosing a mutex protocol

The choice of a mutex protocol is made by setting attributes when creating a mutex. The mutex protocol is controlled through the protocol attribute. This attribute can be set in the mutex attributes object by using the `pthread_mutexattr_getprotocol` and `pthread_mutexattr_setprotocol` subroutines. The protocol attribute can have one of the following values:

Value	Description
<code>PTHREAD_PRIO_DEFAULT</code>	No value
<code>PTHREAD_PRIO_NONE</code>	Denotes no protocol.
<code>PTHREAD_PRIO_INHERIT</code>	Denotes the priority inheritance protocol.
<code>PTHREAD_PRIO_PROTECT</code>	Denotes the priority protection protocol.

Note: The behavior of `PTHREAD_PRIO_DEFAULT` is the same as the `PTHREAD_PRIO_INHERIT` attribute. With reference to the mutex locking, the threads acting with the default attribute will temporarily boost the priority of a mutex holder when a user is locked and has a higher priority than the owner. Therefore, there are only three behaviors that are possible, although there are four values for the possible priority in the attribute structure.

The priority protection protocol uses one additional attribute: the prioceiling attribute. This attribute contains the priority ceiling of the mutex. The prioceiling attribute can be controlled in the mutex attributes object, by using the **pthread_mutexattr_getprioceiling** and **pthread_mutexattr_setprioceiling** subroutines.

The prioceiling attribute of a mutex can also be dynamically controlled by using the **pthread_mutex_getprioceiling** and **pthread_mutex_setprioceiling** subroutines. When dynamically changing the priority ceiling of a mutex, the mutex is locked by the library; it should not be held by the thread calling the **pthread_mutex_setprioceiling** subroutine to avoid a deadlock. Dynamically setting the priority ceiling of a mutex can be useful when increasing the priority of a thread.

The implementation of mutex protocols is optional. Each protocol is a POSIX option.

Inheritance or protection

Both protocols are similar and result in promoting the priority of the thread holding the mutex. If both protocols are available, programmers must choose a protocol. The choice depends on whether the priorities of the threads that will lock the mutex are available to the programmer who is creating the mutex. Typically, mutexes defined by a library and used by application threads will use the inheritance protocol, whereas mutexes created within the application program will use the protection protocol.

In performance-critical programs, performance considerations may also influence the choice. In most implementations, especially in AIX, changing the priority of a thread results in making a system call. Therefore, the two mutex protocols differ in the amount of system calls they generate, as follows:

- Using the inheritance protocol, a system call is made each time a thread is blocked when trying to lock the mutex.
- Using the protection protocol, one system call is always made each time the mutex is locked by a thread.

In most performance-critical programs, the inheritance protocol should be chosen, because mutexes are low contention objects. Mutexes are not held for long periods of time; thus, it is not likely that threads are blocked when trying to lock them.

Related concepts

Scheduling threads

Threads can be scheduled, and the threads library provides several facilities to handle and control the scheduling of threads.

List of synchronization subroutines

This section lists synchronization subroutines.

pthread_mutex_destroy

Deletes a mutex.

pthread_mutex_init

Initializes a mutex and sets its attributes.

PTHREAD_MUTEX_INITIALIZER

Initializes a static mutex with default attributes.

pthread_mutex_lock or pthread_mutex_trylock

Locks a mutex.

pthread_mutex_unlock

Unlocks a mutex.

pthread_mutexattr_destroy

Deletes a mutex attributes object.

pthread_mutexattr_init

Creates a mutex attributes object and initializes it with default values.

pthread_cond_destroy

Deletes a condition variable.

pthread_cond_init

Initializes a condition variable and sets its attributes.

PTHREAD_COND_INITIALIZER

Initializes a static condition variable with default attributes.

pthread_cond_signal or pthread_cond_broadcast

Unblocks one or more threads blocked on a condition.

pthread_cond_wait or pthread_cond_timedwait

Blocks the calling thread on a condition.

pthread_condattr_destroy

Deletes a condition attributes object.

pthread_condattr_init

Creates a condition attributes object and initializes it with default values.

One-time initializations

Some C libraries are designed for dynamic initialization, in which the global initialization for the library is performed when the first procedure in the library is called.

In a single-threaded program, this is usually implemented using a static variable whose value is checked on entry to each routine, as in the following code fragment:

```
static int isInitialized = 0;
extern void Initialize();

int function()
{
    if (isInitialized == 0) {
        Initialize();
        isInitialized = 1;
    }
    ...
}
```

For dynamic library initialization in a multithreaded program, a simple initialization flag is not sufficient. This flag must be protected against modification by multiple threads simultaneously calling a library function. Protecting the flag requires the use of a mutex; however, mutexes must be initialized before they are used. Ensuring that the mutex is only initialized once requires a recursive solution to this problem.

To keep the same structure in a multithreaded program, use the **pthread_once** subroutine. Otherwise, library initialization must be accomplished by an explicit call to a library exported initialization function prior to any use of the library. The **pthread_once** subroutine also provides an alternative for initializing mutexes and condition variables.

One-time initialization object

The uniqueness of the initialization is ensured by the one-time initialization object. It is a variable having the **pthread_once_t** data type. In AIX and most other implementations of the threads library, the **pthread_once_t** data type is a structure.

A one-time initialization object is typically a global variable. It must be initialized with the **PTHREAD_ONCE_INIT** macro, as in the following example:

```
static pthread_once_t once_block = PTHREAD_ONCE_INIT;
```

The initialization can also be done in the initial thread or in any other thread. Several one-time initialization objects can be used in the same program. The only requirement is that the one-time initialization object be initialized with the macro.

One-time initialization routine

The **pthread_once** subroutine calls the specified initialization routine associated with the specified one-time initialization object if it is the first time it is called; otherwise, it does nothing. The same initialization

routine must always be used with the same one-time initialization object. The initialization routine must have the following prototype:

```
void init_routine();
```

The **pthread_once** subroutine does not provide a cancelation point. However, the initialization routine may provide cancelation points, and, if cancelability is enabled, the first thread calling the **pthread_once** subroutine may be canceled during the execution of the initialization routine. In this case, the routine is not considered as executed, and the next call to the **pthread_once** subroutine would result in recalling the initialization routine.

It is recommended to use cleanup handlers in one-time initialization routines, especially when performing non-idempotent operations, such as opening a file, locking a mutex, or allocating memory.

One-time initialization routines can be used for initializing mutexes or condition variables or to perform dynamic initialization. In a multithreaded library, the code fragment shown above (`void init_routine();`) would be written as follows:

```
static pthread_once_t once_block = PTHREAD_ONCE_INIT;
extern void Initialize();

int function()
{
    pthread_once(&once_block, Initialize);
    ...
}
```

Related concepts

[Terminating threads](#)

A thread automatically terminates when it returns from its entry-point routine.

[Thread-specific data](#)

Many applications require that certain data be maintained on a per-thread basis across function calls.

Related reference

[Writing reentrant and threadsafe code](#)

In single-threaded processes, only one flow of control exists. The code executed by these processes thus need not be reentrant or threadsafe. In multithreaded programs, the same functions and the same resources may be accessed concurrently by several flows of control.

Thread-specific data

Many applications require that certain data be maintained on a per-thread basis across function calls.

For example, a multithreaded **grep** command using one thread for each file must have thread-specific file handlers and list of found strings. The thread-specific data interface is provided by the threads library to meet these needs.

Thread-specific data may be viewed as a two-dimensional array of values, with keys serving as the row index and thread IDs as the column index. A thread-specific data *key* is an opaque object, of the **pthread_key_t** data type. The same key can be used by all threads in a process. Although all threads use the same key, they set and access different thread-specific data values associated with that key. Thread-specific data are void pointers, which allows referencing any kind of data, such as dynamically allocated strings or structures.

In the following figure, thread T2 has a thread-specific data value of 12 associated with the key K3. Thread T4 has the value of 2 associated with the same key.

Keys	T1 Thread	T2 Thread	T3 Thread	T4 Thread
K1	6	56	4	1
K2	87	21	0	9
K3	23	12	61	2

Keys	T1 Thread	T2 Thread	T3 Thread	T4 Thread
K4	11	76	47	88

Creating and destroying keys

Thread-specific data keys must be created before being used. Their values can be automatically destroyed when the corresponding threads terminate. A key can also be destroyed upon request to reclaim its storage.

Key creation

A thread-specific data key is created by calling the `pthread_key_create` subroutine. This subroutine returns a key. The thread-specific data is set to a value of **NULL** for all threads, including threads not yet created.

For example, consider two threads *A* and *B*. Thread *A* performs the following operations in chronological order:

1. Create a thread-specific data key *K*.

Threads *A* and *B* can use the key *K*. The value for both threads is **NULL**.

2. Create a thread *C*.

Thread *C* can also use the key *K*. The value for thread *C* is **NULL**.

The number of thread-specific data keys is limited to 450 per process. This number can be retrieved by the `PTHREAD_KEYS_MAX` symbolic constant.

The `pthread_key_create` subroutine must be called only once. Otherwise, two different keys are created. For example, consider the following code fragment:

```

/* a global variable */
static pthread_key_t theKey;

/* thread A */
...
pthread_key_create(&theKey, NULL); /* call 1 */
...

/* thread B */
...
pthread_key_create(&theKey, NULL); /* call 2 */
...

```

In our example, threads *A* and *B* run concurrently, but call 1 happens before call 2. Call 1 will create a key *K1* and store it in the `theKey` variable. Call 2 will create another key *K2*, and store it also in the `theKey` variable, thus overriding *K1*. As a result, thread *A* will use *K2*, assuming it is *K1*. This situation should be avoided for the following reasons:

- Key *K1* is lost, thus its storage will never be reclaimed until the process terminates. Because the number of keys is limited, you may not have enough keys.
- If thread *A* stores a thread-specific data using the `theKey` variable before call 2, the data will be bound to key *K1*. After call 2, the `theKey` variable contains *K2*; if thread *A* then tries to fetch its thread-specific data, it would always get **NULL**.

Ensuring that keys are created uniquely can be done in the following ways:

- Using the one-time initialization facility.
- Creating the key before the threads that will use it. This is often possible, for example, when using a pool of threads with thread-specific data to perform similar operations. This pool of threads is usually created by one thread, the initial (or another "driver") thread.

It is the programmer's responsibility to ensure the uniqueness of key creation. The threads library provides no way to check if a key has been created more than once.

Destructor routine

A destructor routine may be associated with each thread-specific data key. Whenever a thread is terminated, if there is non-**NULL**, thread-specific data for this thread bound to any key, the destructor routine associated with that key is called. This allows dynamically allocated thread-specific data to be automatically freed when the thread is terminated. The destructor routine has one parameter, the value of the thread-specific data.

For example, a thread-specific data key may be used for dynamically allocated buffers. A destructor routine should be provided to ensure that when the thread terminates the buffer is freed, the **free** subroutine can be used as follows:

```
pthread_key_create(&key, free);
```

More complex destructors may be used. If a multithreaded **grep** command, using a thread per file to scan, has thread-specific data to store a structure containing a work buffer and the thread's file descriptor, the destructor routine may be as follows:

```
typedef struct {
    FILE *stream;
    char *buffer;
} data_t;
...

void destructor(void *data)
{
    fclose(((data_t *)data)->stream);
    free(((data_t *)data)->buffer);
    free(data);
    *data = NULL;
}
```

Destructor calls can be repeated up to four times.

Key destruction

A thread-specific data key can be destroyed by calling the **pthread_key_delete** subroutine. The **pthread_key_delete** subroutine does not actually call the destructor routine for each thread having data. After a data key is destroyed, it can be reused by another call to the **pthread_key_create** subroutine. Thus, the **pthread_key_delete** subroutine is useful especially when using many data keys. For example, in the following code fragment, the loop would never end:

```
/* bad example - do not write such code! */
pthread_key_t key;

while (pthread_key_create(&key, NULL))
    pthread_key_delete(key);
```

Using thread-specific data

Thread-specific data is accessed using the **pthread_getspecific** and **pthread_setspecific** subroutines. The **pthread_getspecific** subroutine reads the value bound to the specified key and is specific to the calling thread; the **pthread_setspecific** subroutine sets the value.

Setting successive values

The value bound to a specific key should be a pointer, which can point to any kind of data. Thread-specific data is typically used for dynamically allocated storage, as in the following code fragment:

```
private_data = malloc(...);
pthread_setspecific(key, private_data);
```

When setting a value, the previous value is lost. For example, in the following code fragment, the value of the **old** pointer is lost, and the storage it pointed to may not be recoverable:

```
pthread_setspecific(key, old);
...
pthread_setspecific(key, new);
```


It is the programmer's responsibility to retrieve the old thread-specific data value to reclaim storage before setting the new value. For example, it is possible to implement a **swap_specific** routine in the following manner:

```
int swap_specific(pthread_key_t key, void **old_pt, void *new)
{
    *old_pt = pthread_getspecific(key);
    if (*old_pt == NULL)
        return -1;
    else
        return pthread_setspecific(key, new);
}
```

Such a routine does not exist in the threads library because it is not always necessary to retrieve the previous value of thread-specific data. Such a case occurs, for example, when thread-specific data are pointers to specific locations in a memory pool allocated by the initial thread.

Using destructor routines

When using dynamically allocated thread-specific data, the programmer must provide a destructor routine when calling the **pthread_key_create** subroutine. The programmer must also ensure that, when releasing the storage allocated for thread-specific data, the pointer is set to **NULL**. Otherwise, the destructor routine might be called with an illegal parameter. For example:

```
pthread_key_create(&key, free);
...

...
private_data = malloc(...);
pthread_setspecific(key, private_data);
...

/* bad example! */
...
pthread_getspecific(key, &data);
free(data);
...
```

When the thread terminates, the destructor routine is called for its thread-specific data. Because the value is a pointer to already released memory, an error can occur. To correct this, the following code fragment should be substituted:

```
/* better example! */
...
pthread_getspecific(key, &data);
free(data);
pthread_setspecific(key, NULL);
...
```

When the thread terminates, the destructor routine is not called, because there is no thread-specific data.

Using non-pointer values

Although it is possible to store values that are not pointers, it is not recommended for the following reasons:

- Casting a pointer into a scalar type may not be portable.
- The **NULL** pointer value is implementation-dependent; several systems assign the **NULL** pointer a non-zero value.

If you are sure that your program will never be ported to another system, you may use integer values for thread-specific data.

Related concepts

[Joining threads](#)

Joining a thread means waiting for it to terminate, which can be seen as a specific usage of condition variables.

[One-time initializations](#)

Some C libraries are designed for dynamic initialization, in which the global initialization for the library is performed when the first procedure in the library is called.

Creating complex synchronization objects

The subroutines provided in the threads library can be used as primitives to build more complex synchronization objects.

List of threads-processes interactions subroutines

This section lists threads-processes interactions subroutines.

Creating complex synchronization objects

The subroutines provided in the threads library can be used as primitives to build more complex synchronization objects.

Long Locks

The mutexes provided by the threads library are low-contention objects and should not be held for a very long time. Long locks are implemented with mutexes and condition variables, so that a long lock can be held for a long time without affecting the performance of the program. Long locks should not be used if cancelability is enabled.

A long lock has the **long_lock_t** data type. It must be initialized by the **long_lock_init** routine. The **long_lock**, **long_trylock**, and **long_unlock** subroutine performs similar operations to the **pthread_mutex_lock**, **pthread_mutex_trylock**, and **pthread_mutex_unlock** subroutine.

The following example shows a typical use of condition variables. In this example, the lock owner is not checked. As a result, any thread can unlock any lock. Error handling and cancelation handling are not performed.

```
typedef struct {
    pthread_mutex_t lock;
    pthread_cond_t cond;
    int free;
    int wanted;
} long_lock_t;

void long_lock_init(long_lock_t *ll)
{
    pthread_mutex_init(&ll->lock, NULL);
    pthread_cond_init(&ll->cond);
    ll->free = 1;
    ll->wanted = 0;
}

void long_lock_destroy(long_lock_t *ll)
{
    pthread_mutex_destroy(&ll->lock);
    pthread_cond_destroy(&ll->cond);
}

void long_lock(long_lock_t *ll)
{
    pthread_mutex_lock(&ll->lock);
    ll->wanted++;
    while(!ll->free)
        pthread_cond_wait(&ll->cond);
    ll->wanted--;
    ll->free = 0;
    pthread_mutex_unlock(&ll->lock);
}

int long_trylock(long_lock_t *ll)
{
    int got_the_lock;

    pthread_mutex_lock(&ll->lock);
    got_the_lock = ll->free;
    if (got_the_lock)
        ll->free = 0;
    pthread_mutex_unlock(&ll->lock);
    return got_the_lock;
}
```

```

void long_unlock(long_lock_t *ll)
{
    pthread_mutex_lock(&ll->lock);
    ll->free = 1;
    if (ll->wanted)
        pthread_cond_signal(&ll->cond);
    pthread_mutex_unlock(&ll->lock);
}

```

Semaphores

Traditional semaphores in UNIX systems are interprocess-synchronization facilities. For specific usage, you can implement interthread semaphores.

A semaphore has the **sema_t** data type. It must be initialized by the **sema_init** routine and destroyed with the **sema_destroy** routine. The semaphore wait and semaphore post operations are respectively performed by the **sema_p** and **sema_v** routines.

In the following basic implementation, error handling is not performed, but cancelations are properly handled with cleanup handlers whenever required:

```

typedef struct {
    pthread_mutex_t lock;
    pthread_cond_t cond;
    int count;
} sema_t;

void sema_init(sema_t *sem)
{
    pthread_mutex_init(&sem->lock, NULL);
    pthread_cond_init(&sem->cond, NULL);
    sem->count = 1;
}

void sema_destroy(sema_t *sem)
{
    pthread_mutex_destroy(&sem->lock);
    pthread_cond_destroy(&sem->cond);
}

void p_operation_cleanup(void *arg)
{
    sema_t *sem;

    sem = (sema_t *)arg;
    pthread_mutex_unlock(&sem->lock);
}

void sema_p(sema_t *sem)
{
    pthread_mutex_lock(&sem->lock);
    pthread_cleanup_push(p_operation_cleanup, sem);
    while (sem->count <= 0)
        pthread_cond_wait(&sem->cond, &sem->lock);
    sem->count--;
    /*
     * Note that the pthread_cleanup_pop subroutine will
     * execute the p_operation_cleanup routine
     */
    pthread_cleanup_pop(1);
}

```

```

void sema_v(sema_t *sem)
{
    pthread_mutex_lock(&sem->lock);
    if (sem->count <= 0)
        pthread_cond_signal(&sem->cond);
    sem->count++;
    pthread_mutex_unlock(&sem->lock);
}

```

The counter specifies the number of users that are allowed to use the semaphore. It is never strictly negative; thus, it does not specify the number of waiting users, as for traditional semaphores. This implementation provides a typical solution to the multiple wakeup problem on the **pthread_cond_wait**

subroutine. The semaphore wait operation is cancelable, because the **pthread_cond_wait** subroutine provides a cancellation point.

Write-Priority Read/Write Locks

A write-priority read/write lock provides multiple threads with simultaneous read-only access to a protected resource, and a single thread with write access to the resource while excluding reads. When a writer releases a lock, other waiting writers will get the lock before any waiting reader. Write-priority read/write locks are usually used to protect resources that are more often read than written.

A write-priority read/write lock has the **rwlock_t** data type. It must be initialized by the **rwlock_init** routine. The **rwlock_lock_read** routine locks the lock for a reader (multiple readers are allowed), the **rwlock_unlock_read** routine unlocks it. The **rwlock_lock_write** routine locks the lock for a writer, the **rwlock_unlock_write** routine unlocks it. The proper unlocking routine (for the reader or for the writer) must be called.

In the following example, the lock owner is not checked. As a result, any thread can unlock any lock. Routines, such as the **pthread_mutex_trylock** subroutine, are missing and error handling is not performed, but cancellations are properly handled with cleanup handlers whenever required.

```
typedef struct {
    pthread_mutex_t lock;
    pthread_cond_t rcond;
    pthread_cond_t wcond;
    int lock_count; /* < 0 .. held by writer          */
                  /* > 0 .. held by lock_count readers */
                  /* = 0 .. held by nobody          */
    int waiting_writers; /* count of waiting writers */
} rwlock_t;

void rwlock_init(rwlock_t *rwl)
{
    pthread_mutex_init(&rwl->lock, NULL);
    pthread_cond_init(&rwl->wcond, NULL);
    pthread_cond_init(&rwl->rcond, NULL);
    rwl->lock_count = 0;
    rwl->waiting_writers = 0;
}

void waiting_reader_cleanup(void *arg)
{
    rwlock_t *rwl;

    rwl = (rwlock_t *)arg;
    pthread_mutex_unlock(&rwl->lock);
}

void rwlock_lock_read(rwlock_t *rwl)
{
    pthread_mutex_lock(&rwl->lock);
    pthread_cleanup_push(waiting_reader_cleanup, rwl);
    while ((rwl->lock_count < 0) && (rwl->waiting_writers))
        pthread_cond_wait(&rwl->rcond, &rwl->lock);
    rwl->lock_count++;
    /*
     * Note that the pthread_cleanup_pop subroutine will
     * execute the waiting_reader_cleanup routine
     */
    pthread_cleanup_pop(1);
}

void rwlock_unlock_read(rwlock_t *rwl)
{
    pthread_mutex_lock(&rwl->lock);
    rwl->lock_count--;
    if (!rwl->lock_count)
        pthread_cond_signal(&rwl->wcond);
    pthread_mutex_unlock(&rwl->lock);
}

void waiting_writer_cleanup(void *arg)
{
    rwlock_t *rwl;

    rwl = (rwlock_t *)arg;
    rwl->waiting_writers--;
```

```

        if ((!rwl->waiting_writers) && (rwl->lock_count >= 0))
            /*
             * This only happens if we have been canceled
             */
            pthread_cond_broadcast(&rwl->wcond);
            pthread_mutex_unlock(&rwl->lock);
    }

void rwlock_lock_write(rwlock_t *rwl)
{
    pthread_mutex_lock(&rwl->lock);
    rwl->waiting_writers++;
    pthread_cleanup_push(waiting_writer_cleanup, rwl);
    while (rwl->lock_count)
        pthread_cond_wait(&rwl->wcond, &rwl->lock);
    rwl->lock_count = -1;
    /*
     * Note that the pthread_cleanup_pop subroutine will
     * execute the waiting_writer_cleanup routine
     */
    pthread_cleanup_pop(1);
}

void rwlock_unlock_write(rwlock_t *rwl)
{
    pthread_mutex_lock(&rwl->lock);
    rwl->lock_count = 0;
    if (!rwl->waiting_writers)
        pthread_cond_broadcast(&rwl->rcond);
    else
        pthread_cond_signal(&rwl->wcond);
    pthread_mutex_unlock(&rwl->lock);
}

```

Readers are counted only. When the count reaches zero, a waiting writer may take the lock. Only one writer can hold the lock. When the lock is released by a writer, another writer is awakened, if there is one. Otherwise, all waiting readers are awakened.

The locking routines are cancelable, because they call the [pthread_cond_wait](#) subroutine. Cleanup handlers are therefore registered before calling the subroutine.

Related concepts

[Thread-specific data](#)

Many applications require that certain data be maintained on a per-thread basis across function calls.

[Synchronization overview](#)

One main benefit of using threads is the ease of using synchronization facilities.

[List of threads-processes interactions subroutines](#)

This section lists threads-processes interactions subroutines.

Signal management

Signals in multithreaded processes are an extension of signals in traditional single-threaded programs.

Signal management in multithreaded processes is shared by the process and thread levels, and consists of the following:

- Per-process signal handlers
- Per-thread signal masks
- Single delivery of each signal

Signal handlers and signal masks

Signal handlers are maintained at process level. It is strongly recommended to use the [sigwait](#) subroutine when waiting for signals. The [sigaction](#) subroutine is not recommended because the list of signal handlers is maintained at process level and any thread within the process might change it. If two threads set a signal handler on the same signal, the last thread that called the [sigaction](#) subroutine overrides the setting of the previous thread call; and in most cases, the order in which threads are scheduled cannot be predicted.

Signal masks are maintained at thread level. Each thread can have its own set of signals that will be blocked from delivery. The **sigthreadmask** subroutine must be used to get and set the calling thread's signal mask. The **sigprocmask** subroutine must not be used in multithreaded programs; because unexpected behavior might result.

The **pthread_sigmask** subroutine is very similar to the **sigprocmask** subroutine. The parameters and usage of both subroutines are identical. When porting existing code to support the threads library, you can replace the **sigprocmask** subroutine with the **pthread_sigmask** subroutine.

Signal generation

Signals generated by some action attributable to a particular thread, such as a hardware fault, are sent to the thread that caused the signal to be generated. Signals generated in association with a process ID, a process group ID, or an asynchronous event (such as terminal activity) are sent to the process.

- The **pthread_kill** subroutine sends a signal to a thread. Because thread IDs identify threads within a process, this subroutine can only send signals to threads within the same process.
- The **kill** subroutine (and thus the **kill** command) sends a signal to a process. A thread can send a **Signal** signal to its process by executing the following call:

```
kill(getpid(), Signal);
```

- The **raise** subroutine cannot be used to send a signal to the calling thread's process. The **raise** subroutine sends a signal to the calling thread, as in the following call:

```
pthread_kill(pthread_self(), Signal);
```

This ensures that the signal is sent to the caller of the **raise** subroutine. Thus, library routines written for single-threaded programs can easily be ported to a multithreaded system, because the **raise** subroutine is usually intended to send the signal to the caller.

- The **alarm** subroutine requests that a signal be sent later to the process, and alarm states are maintained at process level. Thus, the last thread that called the **alarm** subroutine overrides the settings of other threads in the process. In a multithreaded program, the **SIGALRM** signal is not necessarily delivered to the thread that called the **alarm** subroutine. The calling thread might even be terminated; and therefore, it cannot receive the signal.

Handling signals

Signal handlers are called within the thread to which the signal is delivered. The following limitations to signal handlers are introduced by the threads library:

- Signal handlers might call the **longjmp** or **siglongjmp** subroutine only if the corresponding call to the **setjmp** or **sigsetjmp** subroutine is performed in the same thread.

Usually, a program that wants to wait for a signal installs a signal handler that calls the **longjmp** subroutine to continue execution at the point where the corresponding **setjmp** subroutine is called. This cannot be done in a multithreaded program, because the signal might be delivered to a thread other than the one that called the **setjmp** subroutine, thus causing the handler to be executed by the wrong thread.

Note: Using **longjmp** from a signal handler can result in undefined behavior.

- No **pthread** routines can be called from a signal handler. Calling a **pthread** routine from a signal handler can lead to an application deadlock.

To allow a thread to wait for asynchronously generated signals, the threads library provides the **sigwait** subroutine. The **sigwait** subroutine blocks the calling thread until one of the awaited signals is sent to the process or to the thread. There must not be a signal handler installed on the awaited signal using the **sigwait** subroutine.

Typically, programs might create a dedicated thread to wait for asynchronously generated signals. Such a thread loops on a **sigwait** subroutine call and handles the signals. It is recommended that such a thread block all the signals. The following code fragment gives an example of such a signal-waiter thread:

```
#include <pthread.h>
#include <signal.h>

static pthread_mutex_t mutex;
sigset_t set;
static int sig_cond = 0;

void *run_me(void *id)
{
    int sig;
    int err;
    sigset_t sigs;
    sigset_t oldSigSet;
    sigfillset(&sigs);
    sigthreadmask(SIG_BLOCK, &sigs, &oldSigSet);

    err = sigwait(&set, &sig);

    if(err)
    {
        /* do error code */
    }
    else
    {
        printf("SIGINT caught\n");
        pthread_mutex_lock(&mutex);
        sig_cond = 1;
        pthread_mutex_unlock(&mutex);
    }

    return;
}

main()
{
    pthread_t tid;

    sigemptyset(&set);
    sigaddset(&set, SIGINT);
    pthread_sigmask(SIG_BLOCK, &set, 0);

    pthread_mutex_init(&mutex, NULL);

    pthread_create(&tid, NULL, run_me, (void *)1);

    while(1)
    {
        sleep(1);
        /* or so something here */

        pthread_mutex_lock(&mutex);
        if(sig_cond)
        {
            /* do exit stuff */
            return;
        }
        pthread_mutex_unlock(&mutex);
    }
}
```

If more than one thread called the **sigwait** subroutine, exactly one call returns when a matching signal is sent. Which thread is awakened cannot be predicted. If a thread is going to do **sigwait** as well as handling of some other signals for which it is not doing **sigwait**, the user-defined signal handlers need to block the sigwaiter signals for the proper handling. Note that the **sigwait** subroutine provides a cancellation point.

Because a dedicated thread is not a real signal handler, it might signal a condition to any other thread. It is possible to implement a **sigwait_multiple** routine that would awaken all threads waiting for a specific signal. Each caller of the **sigwait_multiple** routine registers a set of signals. The caller then waits on a condition variable. A single thread calls the **sigwait** subroutine on the union of all registered signals. When the call to the **sigwait** subroutine returns, the appropriate state is set and condition variables are

broadcasted. New callers to the **sigwait_multiple** subroutine cause the pending **sigwait** subroutine call to be canceled and reissued to update the set of signals being waited for.

Signal delivery

A signal is delivered to a thread, unless its action is set to ignore. The following rules govern signal delivery in a multithreaded process:

- A signal whose action is set to terminate, stop, or continue the target thread or process respectively terminates, stops, or continues the entire process (and thus all of its threads). Single-threaded programs can thus be rewritten as multithreaded programs without changing their externally visible signal behavior.

For example, consider a multithreaded user command, such as the **grep** command. A user can start the command in his favorite shell and then decide to stop it by sending a signal with the **kill** command. The signal should stop the entire process running the **grep** command.

- Signals generated for a specific thread, using the **pthread_kill** or the **raise** subroutines, are delivered to that thread. If the thread has blocked the signal from delivery, the signal is set pending on the thread until the signal is unblocked from delivery. If the thread is terminated before the signal delivery, the signal will be ignored.
- Signals generated for a process, using the **kill** subroutine for example, are delivered to exactly one thread in the process. If one or more threads called the **sigwait** subroutine, the signal is delivered to exactly one of these threads. Otherwise, the signal is delivered to exactly one thread that did not block the signal from delivery. If no thread matches these conditions, the signal is set pending on the process until a thread calls the **sigwait** subroutine specifying this signal or a thread unblocks the signal from delivery.

If the action associated with a pending signal (on a thread or on a process) is set to ignore, the signal is ignored.

Related concepts

Process duplication and termination

Because all processes have at least one thread, creating (that is, duplicating) and terminating a process implies the creation and the termination of threads.

List of threads-processes interactions subroutines

This section lists threads-processes interactions subroutines.

Subroutine	Description
<u>alarm</u>	Causes a signal to be sent to the calling process after a specified timeout.
<u>kill</u> or <u>killpg</u>	Sends a signal to a process or a group of processes.
<u>pthread_atfork</u>	Registers fork cleanup handlers.
<u>pthread_kill</u>	Sends a signal to the specified thread.
<u>pthread_sigmask</u>	Sets the signal mask of a thread.
<u>raise</u>	Sends a signal to the executing thread.
<u>sigaction</u> , <u>sigvec</u> , or <u>signal</u>	Specifies the action to take upon delivery of a signal.
<u>sigsuspend</u> or <u>sigpause</u>	Atomically changes the set of blocked signals and waits for a signal.
<u>sigthreadmask</u>	Sets the signal mask of a thread.
<u>sigwait</u>	Blocks the calling thread until a specified signal is received.

Related concepts

Thread-specific data

Many applications require that certain data be maintained on a per-thread basis across function calls.

Creating complex synchronization objects

The subroutines provided in the threads library can be used as primitives to build more complex synchronization objects.

Process duplication and termination

Because all processes have at least one thread, creating (that is, duplicating) and terminating a process implies the creation and the termination of threads.

Process duplication and termination

Because all processes have at least one thread, creating (that is, duplicating) and terminating a process implies the creation and the termination of threads.

This section describes the interactions between threads and processes when duplicating and terminating a process.

Forking

Programmers call the **fork** subroutine in the following cases:

- To create a new flow of control within the same program. AIX creates a new process.
- To create a new process running a different program. In this case, the call to the **fork** subroutine is soon followed by a call to one of the **exec** subroutines.

In a multithreaded program, the first use of the **fork** subroutine, creating new flows of control, is provided by the **pthread_create** subroutine. The **fork** subroutine should thus be used only to run new programs.

The **fork** subroutine duplicates the parent process, but duplicates only the calling thread; the child process is a single-threaded process. The calling thread of the parent process becomes the initial thread of the child process; it may not be the initial thread of the parent process. Thus, if the initial thread of the child process returns from its entry-point routine, the child process terminates.

When duplicating the parent process, the **fork** subroutine also duplicates all the synchronization variables, including their state. Thus, for example, mutexes may be held by threads that no longer exist in the child process and any associated resource may be inconsistent.

It is strongly recommended that the **fork** subroutine be used only to run new programs, and to call one of the **exec** subroutines as soon as possible after the call to the **fork** subroutine in the child process.

Fork handlers

The preceding forking rule does not address the needs of multithreaded libraries. Application programs may not be aware that a multithreaded library is in use and will call any number of library routines between the **fork** and the **exec** subroutines, just as they always have. Indeed, they may be old single-threaded programs and cannot, therefore, be expected to obey new restrictions imposed by the threads library.

On the other hand, multithreaded libraries need a way to protect their internal state during a fork in case a routine is called later in the child process. The problem arises especially in multithreaded input/output libraries, which are almost sure to be invoked between the **fork** and the **exec** subroutines to affect input/output redirection.

The **pthread_atfork** subroutine provides a way for multithreaded libraries to protect themselves from innocent application programs that call the **fork** subroutine. It also provides multithreaded application programs with a standard mechanism for protecting themselves from calls to the **fork** subroutine in a library routine or the application itself.

The **pthread_atfork** subroutine registers fork handlers to be called before and after the call to the **fork** subroutine. The fork handlers are executed in the thread that called the **fork** subroutine. The following fork handlers exist:

Subroutine Description

Prepare	The prepare fork handler is called just before the processing of the fork subroutine begins.
Parent	The parent fork handler is called just after the processing of the fork subroutine is completed in the parent process.
Child	The child fork handler is called just after the processing of the fork subroutine is completed in the child process.

Process termination

The prepare fork handlers are called in last-in first-out (LIFO) order, whereas the parent and child fork handlers are called in first-in first-out (FIFO) order. This allows programs to preserve any desired locking order.

When a process terminates, by calling the **exit**, **atexit**, or **_exit** subroutine either explicitly or implicitly, all threads within the process are terminated. Neither the cleanup handlers nor the thread-specific data destructors are called.

Note: The `unatexit` subroutine unregisters functions that were previously registered by the `atexit` subroutine. If the referenced function is found, it is removed from the list of functions that are called at normal program termination.

The reason for this behavior is that there is no state to leave clean and no thread-specific storage to reclaim, because the whole process terminates, including all the threads, and all the process storage is reclaimed, including all thread-specific storage.

Related concepts

[Understanding threads and processes](#)

A *thread* is an independent flow of control that operates within the same address space as other independent flows of controls within a process.

[Signal management](#)

Signals in multithreaded processes are an extension of signals in traditional single-threaded programs.

[List of threads-processes interactions subroutines](#)

This section lists threads-processes interactions subroutines.

Threads library options

This section describes special attributes of threads, mutexes, and condition variables.

The POSIX standard for the threads library specifies the implementation of some parts as optional. All subroutines defined by the threads library API are always available. Depending on the available options, some subroutines may not be implemented. Unimplemented subroutines can be called by applications, but they always return the **ENOSYS** error code.

Stack attributes

A stack is allocated for each thread. Stack management is implementation-dependent. Thus, the following information applies only to AIX, although similar features may exist on other systems.

The stack is dynamically allocated when the thread is created. Using advanced thread attributes, it is possible for the user to control the stack size and address of the stack. The following information does not apply to the initial thread, which is created by the system.

Stack size

The stack size option enables the control of the **stacksize** attribute of a thread attributes object. This attribute specifies the minimum stack size to be used for the created thread.

The **stacksize** attribute is defined in AIX. The following attribute and subroutines are available when the option is implemented:

- The **stacksize** attribute of the thread attributes object
- The **pthread_attr_getstacksize** returns the value of the attribute
- and **pthread_attr_setstacksize** subroutines sets the value

The default value of the **stacksize** attribute is 96 KB. The minimum value of the **stacksize** attribute is 16 KB. If the assigned value is less than the minimum value, the minimum value is allocated.

In the AIX implementation of the threads library, a chunk of data, called *user thread area*, is allocated for each created thread. The area is divided into the following sections:

- A *red zone*, which is both read-protected and write-protected for stack overflow detection. There is no red zone in programs that use large pages.
- A default stack.
- A pthread structure.
- A thread structure.
- A thread attribute structure.

Note: The user thread area described here has no relationship to the **uthread** structure used in the AIX kernel. The user thread area is accessed only in user mode and is exclusively handled by the threads library, whereas the **uthread** structure only exists within the kernel environment.

Stack address POSIX option

The stack address option enables the control of the **stackaddr** attribute of a thread attributes object. This attribute specifies the location of storage to be used for the created thread's stack.

The following attribute and subroutines are available when the option is implemented:

- The **stackaddr** attribute of the thread attributes object specifies the address of the stack that will be allocated for a thread.
- The **pthread_attr_getstackaddr** subroutine returns the value of the attribute.
- and **pthread_attr_setstackaddr** subroutine sets the value.

If no stack address is specified, the stack is allocated by the system at an arbitrary address. If you must have the stack at a known location, you can use the **stackaddr** attribute. For example, if you need a very large stack, you can set its address to an unused segment, guaranteeing that the allocation will succeed.

If a stack address is specified when calling the **pthread_create** subroutine, the system attempts to allocate the stack at the given address. If it fails, the **pthread_create** subroutine returns **EINVAL**.

Because the **pthread_attr_setstackaddr** subroutine does not actually allocate the stack, it only returns an error if the specified stack address exceeds the addressing space.

Priority scheduling POSIX option

The priority scheduling option enables the control of execution scheduling at thread level. When this option is disabled, all threads within a process share the scheduling properties of the process. When this option is enabled, each thread has its own scheduling properties. For local contention scope threads, the scheduling properties are handled at process level by a library scheduler, while for global contention scope threads, the scheduling properties are handled at system level by the kernel scheduler.

The following attributes and subroutines are available when the option is implemented:

- The **inheritsched** attribute of the thread attributes object
- The **schedparam** attribute of the thread attributes object and the thread
- The **schedpolicy** attribute of the thread attributes objects and the thread
- The **contention-scope** attribute of the thread attributes objects and the thread
- The **pthread_attr_getschedparam** and **pthread_attr_setschedparam** subroutines
- The **pthread_getschedparam** subroutine

Checking the availability of an option

Options can be checked at compile time or at run time. Portable programs should check the availability of options before using them, so that they need not be rewritten when ported to other systems.

Compile-time checking

When an option is not available, you can stop the compilation, as in the following example:

```
#ifndef _POSIX_THREAD_ATTR_STACKSIZE
#error "The stack size POSIX option is required"
#endif
```

The **pthread.h** header file also defines the following symbols that can be used by other header files or by programs:

_POSIX_REENTRANT_FUNCTIONS

Denotes that reentrant functions are required

_POSIX_THREADS

Denotes the implementation of the threads library

Run-time checking

The **sysconf** subroutine can be used to get the availability of options on the system where the program is executed. This is useful when porting programs between systems that have a binary compatibility, such as two versions of AIX.

The following list indicates the symbols that are associated with each option and that must be used for the *Name* parameter of the **sysconf** subroutine. The symbolic constants are defined in the **unistd.h** header file.

Stack address

_SC_THREAD_ATTR_STACKADDR

Stack size

_SC_THREAD_ATTR_STACKSIZE

Priority scheduling

_SC_THREAD_PRIORITY_SCHEDULING

Priority inheritance

_SC_THREAD_PRIO_INHERIT

Priority protection

_SC_THREAD_PRIO_PROTECT

Process sharing

_SC_THREAD_PROCESS_SHARED

To check the general options, use the **sysconf** subroutine with the following *Name* parameter values:

_SC_REENTRANT_FUNCTIONS

Denotes that reentrant functions are required.

_SC_THREADS

Denotes the implementation of the threads library.

Process sharing

AIX and most UNIX systems allow several processes to share a common data space, known as *shared memory*. The process-sharing attributes for condition variables and mutexes are meant to allow these objects to be allocated in shared memory to support synchronization among threads belonging to different processes. However, because there is no industry-standard interface for shared memory management, the process-sharing POSIX option is not implemented in the AIX threads library.

Threads data types

The following data types are defined for the threads library. The definition of these data types can vary between systems:

pthread_t

Identifies a thread

pthread_attr_t

Identifies a thread attributes object

pthread_cond_t

Identifies a condition variable

pthread_condattr_t

Identifies a condition attributes object

pthread_key_t

Identifies a thread-specific data key

pthread_mutex_t

Identifies a mutex

pthread_mutexattr_t

Identifies a mutex attributes object

pthread_once_t

Identifies a one-time initialization object

Limits and Default Values

The threads library has some implementation-dependent limits and default values. These limits and default values can be retrieved by symbolic constants to enhance the portability of programs:

- The maximum number of threads per process is 512. The maximum number of threads can be retrieved at compilation time using the **PTHREAD_THREADS_MAX** symbolic constant defined in the **pthread.h** header file. If an application is compiled with the **-D_LARGE_THREADS** flag, the maximum number of threads per process is 32767.
- The minimum stack size for a thread is 8 K. The default stack size is 96 KB. This number can be retrieved at compilation time using the **PTHREAD_STACK_MIN** symbolic constant defined in the **pthread.h** header file.

Note: The maximum stack size is 256 MB, the size of a segment. This limit is indicated by the **PTHREAD_STACK_MAX** symbolic constant in the **pthread.h** header file.

- The maximum number of thread-specific data keys is limited to 508. This number can be retrieved at compilation time using the **PTHREAD_KEYS_MAX** symbolic constant defined in the **pthread.h** header file.

Default attribute values

The default values for the thread attributes object are defined in the **pthread.h** header file by the following symbolic constants:

- The default value for the **DEFAULT_DETACHSTATE** symbolic constant is **PTHREAD_CREATE_DETACHED**, which specifies the default value for the **detachstate** attribute.
- The default value for the **DEFAULT_JOINABLE** symbolic constant is **PTHREAD_CREATE_JOINABLE**, which specifies the default value for the joinable state.
- The default value for the **DEFAULT_INHERIT** symbolic constant is **PTHREAD_INHERIT_SCHED**, which specifies the default value for the **inheritsched** attribute.
- The default value for the **DEFAULT_PRIO** symbolic constant is 1, which specifies the default value for the **sched_prio** field of the **schedparam** attribute.
- The default value for the **DEFAULT_SCHED** symbolic constant is **SCHED_OTHER**, which specifies the default value for the **schedpolicy** attribute of a thread attributes object.
- The default value for the **DEFAULT_SCOPE** symbolic constant is **PTHREAD_SCOPE_LOCAL**, which specifies the default value for the **contention-scope** attribute.

Related concepts

[Threadsafe and threaded libraries in AIX](#)

This section describes the thread libraries in the AIX.

List of threads advanced-feature subroutines

This section lists threads advanced-features subroutines.

Subroutine	Description
<u>pthread_attr_getstackaddr</u>	Returns the value of the stackaddr attribute of a thread attributes object.
<u>pthread_attr_getstacksize</u>	Returns the value of the stacksize attribute of a thread attributes object.
<u>pthread_attr_setstackaddr</u>	Sets the value of the stackaddr attribute of a thread attributes object.
<u>pthread_attr_setstacksize</u>	Sets the value of the stacksize attribute of a thread attributes object.
<u>pthread_condattr_getpshared</u>	Returns the value of the process-shared attribute of a condition attributes object.
<u>pthread_condattr_setpshared</u>	Sets the value of the process-shared attribute of a condition attributes object.
<u>pthread_getspecific</u>	Returns the thread-specific data associated with the specified key.
<u>pthread_key_create</u>	Creates a thread-specific data key.
<u>pthread_key_delete</u>	Deletes a thread-specific data key.
<u>pthread_mutexattr_getpshared</u>	Returns the value of the process-shared attribute of a mutex attributes object.
<u>pthread_mutexattr_setpshared</u>	Sets the value of the process-shared attribute of a mutex attributes object.
<u>pthread_once</u>	Executes a routine exactly once in a process.
<u>PTHREAD_ONCE_INIT</u>	Initializes a one-time synchronization control structure.
<u>pthread_setspecific</u>	Sets the thread-specific data associated with the specified key.

Supported interfaces

On AIX systems, the **_POSIX_THREADS**, **_POSIX_THREAD_ATTR_STACKADDR**, **_POSIX_THREAD_ATTR_STACKSIZE** and **_POSIX_THREAD_PROCESS_SHARED** symbols are always defined.

Therefore, the following threads interfaces are supported.

POSIX interfaces

The following is a list of POSIX interfaces:

- pthread_atfork
- pthread_attr_destroy
- pthread_attr_getdetachstate
- pthread_attr_getschedparam
- pthread_attr_getstacksize
- pthread_attr_getstackaddr

- pthread_attr_init
- pthread_attr_setdetachstate
- pthread_attr_setschedparam
- pthread_attr_setstackaddr
- pthread_attr_setstacksize
- pthread_cancel
- pthread_cleanup_pop
- pthread_cleanup_push
- pthread_detach
- pthread_equal
- pthread_exit
- pthread_getspecific
- pthread_join
- pthread_key_create
- pthread_key_delete
- pthread_kill
- pthread_mutex_destroy
- pthread_mutex_init
- pthread_mutex_lock
- pthread_mutex_trylock
- pthread_mutex_unlock
- pthread_mutexattr_destroy
- pthread_mutexattr_getpshared
- pthread_mutexattr_init
- pthread_mutexattr_setpshared
- pthread_once
- pthread_self
- pthread_setcancelstate
- pthread_setcanceltype
- pthread_setspecific
- pthread_sigmask
- pthread_testcancel
- pthread_cond_broadcast
- pthread_cond_destroy
- pthread_cond_init
- pthread_cond_signal
- pthread_cond_timedwait
- pthread_cond_wait
- pthread_condattr_destroy
- pthread_condattr_getpshared
- pthread_condattr_init
- pthread_condattr_setpshared
- pthread_create

- sigwait

Single UNIX Specification, Version 2 Interfaces

The following is a list of Single UNIX Specification, Version 2 interfaces:

- pthread_attr_getguardsize
- pthread_attr_setguardsize
- pthread_getconcurrency
- pthread_mutexattr_gettype
- pthread_mutexattr_settype
- pthread_rwlock_destroy
- pthread_rwlock_init
- pthread_rwlock_rdlock
- pthread_rwlock_tryrdlock
- pthread_rwlock_trywrlock
- pthread_rwlock_unlock
- pthread_rwlock_wrlock
- pthread_rwlockattr_destroy
- pthread_rwlockattr_getpshared
- pthread_rwlockattr_init
- pthread_rwlockattr_setpshared
- pthread_setconcurrency

On AIX systems, **_POSIX_THREAD_SAFE_FUNCTIONS** symbol is always defined. Therefore, the following interfaces are always supported:

- asctime_r
- ctime_r
- flockfile
- ftrylockfile
- funlockfile
- getc_unlocked
- getchar_unlocked
- getgrgid_r
- getgrnam_r
- getpwnam_r
- getpwuid_r
- gmtime_r
- localtime_r
- putc_unlocked
- putchar_unlocked
- rand_r
- readdir_r
- strtok_r

AIX does not support the following interfaces; the symbols are provided but they always return an error and set the errno to ENOSYS:

- pthread_mutex_getprioceiling

- pthread_mutex_setprioceiling
- pthread_mutexattr_getprioceiling
- pthread_mutexattr_getprotocol
- pthread_mutexattr_setprioceiling
- pthread_mutexattr_setprotocol

Non-threadsafe interfaces

libc.a library (standard functions):

- advance
- asctime
- brk
- catgets
- chroot
- compile
- ctime
- cuserid
- dbm_clearerr
- dbm_close
- dbm_delete
- dbm_error
- dbm_fetch
- dbm_firstkey
- dbm_nextkey
- dbm_open
- dbm_store
- dirname
- drand48
- ecvt
- encrypt
- endgrent
- endpwent
- endutxent
- fcvt
- gamma
- gcvt
- getc_unlocked
- getchar_unlocked
- getdate
- getdtablesize
- getgrent
- getgrgid
- getgrnam
- getlogin
- getopt

- getpagesize
- getpass
- getpwent
- getpwnam
- getpwuid
- getutxent
- getutxid
- getutxline
- getw
- getw
- gmtime
- l64a
- lgamma
- localtime
- lrand48
- mrand48
- nl_langinfo
- ptsname
- putc_unlocked
- putchar_unlocked
- pututxline
- putw
- rand
- random
- readdir
- re_comp
- re_exec
- regcmp
- regex
- sbrk
- setgrent
- setkey
- setpwent
- setutxent
- sigstack
- srand48
- srandom
- step
- strerror
- strtok
- ttyname
- ttyslot
- wait3

The following AIX interfaces are not threadsafe.

libc.a Library (AIX-specific functions):

- endfsent
- endttyent
- endutent
- getfsent
- getfsfile
- getfsspec
- getfstype
- getttyent
- getttynam
- getutent
- getutid
- getutline
- pututline
- setfsent
- setttyent
- setutent
- utmpname

libbsd.a library:

- timezone

libm.a and **libmsaa.a** libraries:

- gamma
- lgamma

None of the functions in the following libraries are threadsafe:

- libPW.a
- libblas.a
- libcur.a
- libcurses.a
- libplot.a
- libprint.a

The **ctermid** and **tmpnam** interfaces are not threadsafe if they are passed a NULL argument.

In a multi-threaded program it is not recommended to execute **setlocale()** subroutine simultaneously from multiple threads if one of the threads calls **setlocale()** subroutine from within a module-initialization routine.

Note: Certain subroutines may be implemented as macros on some systems. Avoid using the address of threads subroutines.

Writing reentrant and threadsafe code

In single-threaded processes, only one flow of control exists. The code executed by these processes thus need not be reentrant or threadsafe. In multithreaded programs, the same functions and the same resources may be accessed concurrently by several flows of control.

To protect resource integrity, code written for multithreaded programs must be reentrant and threadsafe.

Reentrance and thread safety are both related to the way that functions handle resources. Reentrance and thread safety are separate concepts: a function can be either reentrant, threadsafe, both, or neither.

This section provides information about writing reentrant and threadsafe programs. It does not cover the topic of writing thread-efficient programs. Thread-efficient programs are efficiently parallelized programs. You must consider thread efficiency during the design of the program. Existing single-threaded programs can be made thread-efficient, but this requires that they be completely redesigned and rewritten.

Reentrance

A reentrant function does not hold static data over successive calls, nor does it return a pointer to static data. All data is provided by the caller of the function. A reentrant function must not call non-reentrant functions.

A non-reentrant function can often, but not always, be identified by its external interface and its usage. For example, the **strtok** subroutine is not reentrant, because it holds the string to be broken into tokens. The **ctime** subroutine is also not reentrant; it returns a pointer to static data that is overwritten by each call.

Thread safety

A threadsafe function protects shared resources from concurrent access by locks. Thread safety concerns only the implementation of a function and does not affect its external interface.

In C language, local variables are dynamically allocated on the stack. Therefore, any function that does not use static data or other shared resources is trivially threadsafe, as in the following example:

```
/* threadsafe function */
int diff(int x, int y)
{
    int delta;

    delta = y - x;
    if (delta < 0)
        delta = -delta;

    return delta;
}
```

The use of global data is thread-unsafe. Global data should be maintained per thread or encapsulated, so that its access can be serialized. A thread may read an error code corresponding to an error caused by another thread. In AIX, each thread has its own **errno** value.

Making a function reentrant

In most cases, non-reentrant functions must be replaced by functions with a modified interface to be reentrant. Non-reentrant functions cannot be used by multiple threads. Furthermore, it may be impossible to make a non-reentrant function threadsafe.

Returning data

Many non-reentrant functions return a pointer to static data. This can be avoided in the following ways:

- Returning dynamically allocated data. In this case, it will be the caller's responsibility to free the storage. The benefit is that the interface does not need to be modified. However, backward compatibility is not ensured; existing single-threaded programs using the modified functions without changes would not free the storage, leading to memory leaks.
- Using caller-provided storage. This method is recommended, although the interface must be modified.

For example, a **strtoupper** function, converting a string to uppercase, could be implemented as in the following code fragment:

```
/* non-reentrant function */
char *strtoupper(char *string)
```

```

{
    static char buffer[MAX_STRING_SIZE];
    int index;

    for (index = 0; string[index]; index++)
        buffer[index] = toupper(string[index]);
    buffer[index] = 0

    return buffer;
}

```

This function is not reentrant (nor threadsafe). To make the function reentrant by returning dynamically allocated data, the function would be similar to the following code fragment:

```

/* reentrant function (a poor solution) */
char *strttoupper(char *string)
{
    char *buffer;
    int index;

    /* error-checking should be performed! */
    buffer = malloc(MAX_STRING_SIZE);

    for (index = 0; string[index]; index++)
        buffer[index] = toupper(string[index]);
    buffer[index] = 0

    return buffer;
}

```

A better solution consists of modifying the interface. The caller must provide the storage for both input and output strings, as in the following code fragment:

```

/* reentrant function (a better solution) */
char *strttoupper_r(char *in_str, char *out_str)
{
    int index;

    for (index = 0; in_str[index]; index++)
        out_str[index] = toupper(in_str[index]);
    out_str[index] = 0

    return out_str;
}

```

The non-reentrant standard C library subroutines were made reentrant using caller-provided storage.

Keeping data over successive calls

No data should be kept over successive calls, because different threads may successively call the function. If a function must maintain some data over successive calls, such as a working buffer or a pointer, the caller should provide this data.

Consider the following example. A function returns the successive lowercase characters of a string. The string is provided only on the first call, as with the **strtok** subroutine. The function returns 0 when it reaches the end of the string. The function could be implemented as in the following code fragment:

```

/* non-reentrant function */
char lowercase_c(char *string)
{
    static char *buffer;
    static int index;
    char c = 0;

    /* stores the string on first call */
    if (string != NULL) {
        buffer = string;
        index = 0;
    }

    /* searches a lowercase character */
    for (; c = buffer[index]; index++) {
        if (islower(c)) {
            index++;
        }
    }
}

```

```

        }
        }
        return c;
}

```

h

This function is not reentrant. To make it reentrant, the static data, the **index** variable, must be maintained by the caller. The reentrant version of the function could be implemented as in the following code fragment:

```

/* reentrant function */
char reentrant_lowercase_c(char *string, int *p_index)
{
    char c = 0;

    /* no initialization - the caller should have done it */

    /* searches a lowercase character */
    for (; c = string[*p_index]; (*p_index)++) {
        if (islower(c)) {
            (*p_index)++;
            break;
        }
    }
    return c;
}

```

The interface of the function changed and so did its usage. The caller must provide the string on each call and must initialize the index to 0 before the first call, as in the following code fragment:

```

char *my_string;
char my_char;
int my_index;
...
my_index = 0;
while (my_char = reentrant_lowercase_c(my_string, &my_index)) {
    ...
}

```

Making a function threadsafe

In multithreaded programs, all functions called by multiple threads must be threadsafe. However, a workaround exists for using thread-unsafe subroutines in multithreaded programs. Non-reentrant functions usually are thread-unsafe, but making them reentrant often makes them threadsafe, too.

Locking shared resources

Functions that use static data or any other shared resources, such as files or terminals, must serialize the access to these resources by locks in order to be threadsafe. For example, the following function is thread-unsafe:

```

/* thread-unsafe function */
int increment_counter()
{
    static int counter = 0;

    counter++;
    return counter;
}

```

To be threadsafe, the static variable **counter** must be protected by a static lock, as in the following example:

```

/* pseudo-code threadsafe function */
int increment_counter();
{
    static int counter = 0;
    static lock_type counter_lock = LOCK_INITIALIZER;
}

```

```

pthread_mutex_lock(counter_lock);
counter++;
pthread_mutex_unlock(counter_lock);
return counter;
}

```

In a multithreaded application program using the threads library, mutexes should be used for serializing shared resources. Independent libraries may need to work outside the context of threads and, thus, use other kinds of locks.

Workarounds for thread-unsafe functions

It is possible to use a workaround to use thread-unsafe functions called by multiple threads. This can be useful, especially when using a thread-unsafe library in a multithreaded program, for testing or while waiting for a threadsafe version of the library to be available. The workaround leads to some overhead, because it consists of serializing the entire function or even a group of functions. The following are possible workarounds:

- Use a global lock for the library, and lock it each time you use the library (calling a library routine or using a library global variable). This solution can create performance bottlenecks because only one thread can access any part of the library at any given time. The solution in the following pseudocode is acceptable only if the library is seldom accessed, or as an initial, quickly implemented workaround.

```

/* this is pseudo code! */

lock(library_lock);
library_call();
unlock(library_lock);

lock(library_lock);
x = library_var;
unlock(library_lock);

```

- Use a lock for each library component (routine or global variable) or group of components. This solution is somewhat more complicated to implement than the previous example, but it can improve performance. Because this workaround should only be used in application programs and not in libraries, mutexes can be used for locking the library.

```

/* this is pseudo-code! */

lock(library_moduleA_lock);
library_moduleA_call();
unlock(library_moduleA_lock);

lock(library_moduleB_lock);
x = library_moduleB_var;
unlock(library_moduleB_lock);

```

Reentrant and threadsafe libraries

Reentrant and threadsafe libraries are useful in a wide range of parallel (and asynchronous) programming environments, not just within threads. It is a good programming practice to always use and write reentrant and threadsafe functions.

Using libraries

Several libraries shipped with the AIX Base Operating System are threadsafe. In the current version of AIX, the following libraries are threadsafe:

- Standard C library (**libc.a**)
- Berkeley compatibility library (**libbsd.a**)

Some of the standard C subroutines are non-reentrant, such as the **ctime** and **strtok** subroutines. The reentrant version of the subroutines have the name of the original subroutine with a suffix **_r** (underscore followed by the letter *r*).

When writing multithreaded programs, use the reentrant versions of subroutines instead of the original version. For example, the following code fragment:

```
token[0] = strtok(string, separators);
i = 0;
do {
    i++;
    token[i] = strtok(NULL, separators);
} while (token[i] != NULL);
```

should be replaced in a multithreaded program by the following code fragment:

```
char *pointer;
...
token[0] = strtok_r(string, separators, &pointer);
i = 0;
do {
    i++;
    token[i] = strtok_r(NULL, separators, &pointer);
} while (token[i] != NULL);
```

Thread-unsafe libraries may be used by only one thread in a program. Ensure the uniqueness of the thread using the library; otherwise, the program will have unexpected behavior, or may even stop.

Converting libraries

Consider the following when converting an existing library to a reentrant and threadsafe library. This information applies only to C language libraries.

- Identify exported global variables. Those variables are usually defined in a header file with the **export** keyword. Exported global variables should be encapsulated. The variable should be made private (defined with the **static** keyword in the library source code), and access (read and write) subroutines should be created.
- Identify static variables and other shared resources. Static variables are usually defined with the **static** keyword. Locks should be associated with any shared resource. The granularity of the locking, thus choosing the number of locks, impacts the performance of the library. To initialize the locks, the one-time initialization facility may be used.
- Identify non-reentrant functions and make them reentrant. For more information, see [Making a Function Reentrant](#).
- Identify thread-unsafe functions and make them threadsafe. For more information, see [Making a Function threadsafe](#).

Related concepts

[One-time initializations](#)

Some C libraries are designed for dynamic initialization, in which the global initialization for the library is performed when the first procedure in the library is called.

Related information

[admin](#)

[cdc](#)

[delta](#)

[get](#)

[prs](#)

[sccsdiff](#)

[scsfile](#)

Developing multithreaded programs

Developing multithreaded programs is similar to developing programs with multiple processes. Developing programs also consists of compiling and debugging the code.

Compiling a multithreaded program

This section explains how to generate a multithreaded program. It describes the following:

- The required header file
- Invoking the compiler, which is used to generate multithreaded programs.

Header file

All subroutine prototypes, macros, and other definitions for using the threads library are in the **pthread.h** header file, which is located in the **/usr/include** directory. The **pthread.h** header file must be included in each source file using the threads library.

The **pthread.h** header includes the **unistd.h** header, which provides the following global definitions:

_POSIX_REENTRANT_FUNCTIONS

Specifies that all functions should be reentrant. Several header files use this symbol to define supplementary reentrant subroutines, such as the **localtime_r** subroutine.

_POSIX_THREADS

Denotes the POSIX threads API. This symbol is used to check if the POSIX threads API is available. Macros or subroutines may be defined in different ways, depending on whether the POSIX or some other threads API is used.

The **pthread.h** file also includes **errno.h**, in which the **errno** global variable is redefined to be thread-specific. The **errno** identifier is, therefore, no longer an **l-value** in a multithreaded program.

Invoking the compiler

When compiling a multithreaded program, invoke the C compiler using one of the following commands:

xlc_r

Invokes the compiler with default language level of **ansi**

cc_r

Invokes the compiler with default language level of **extended**

These commands ensure that the adequate options and libraries are used to be compliant with the Single UNIX Specification, Version 2. The POSIX Threads Specification 1003.1c is a subset of the Single UNIX Specification, Version 2.

The following libraries are automatically linked with your program when using the **xlc_r** and **cc_r** commands:

libpthreads.a

Threads library

libc.a

Standard C library

For example, the following command compiles the **foo.c** multithreaded C source file and produces the **foo** executable file:

```
cc_r -o foo foo.c
```

Invoking the compiler for draft 7 of POSIX 1003.1c

AIX provides source code compatibility for Draft 7 applications. It is recommended that developers port their threaded application to the latest standard.

When compiling a multithreaded program for Draft 7 support of threads, invoke the C compiler using one of the following commands:

xlc_r7

Invokes the compiler with default language level of **ansi**

cc_r7

Invokes the compiler with default language level of **extended**

The following libraries are automatically linked with your program when using the **xlc_r7** and **cc_r7** commands:

libpthreads_compat.a

Draft 7 Compatibility Threads library

libpthreads.a

Threads library

libc.a

Standard C library

To achieve source code compatibility, use the compiler directive **_AIX_PTHREADS_D7**. It is also necessary to link the libraries in the following order: **libpthreads_compat.a**, **libpthreads.a**, and **libc.a**. Most users do not need to know this information, because the commands provide the necessary options. These options are provided for those who do not have the latest AIX compiler.

Porting draft 7 applications to the &Symbol.unixspec;

Differences exist between Draft 7 and the final standard include:

- Minor **errno** differences. The most prevalent is the use of **ESRCH** to denote the specified pthread could not be found. Draft 7 frequently returned **EINVAL** for this failure.
- The default state when a pthread is created is *joinable*. This is a significant change because it can result in a memory leak if ignored.
- The default pthread scheduling parameter is *scope*.
- The **pthread_yield** subroutine has been replaced by the **sched_yield** subroutine.
- The various scheduling policies associated with the mutex locks are slightly different.

Memory requirements of a multithreaded program

AIX supports up to 32768 threads in a single process. Each individual pthread requires some amount of process address space, so the actual maximum number of pthreads that a process can have depends on the memory model and the use of process address space, for other purposes. The amount of memory that a pthread needs includes the stack size and the guard region size, plus some amount for internal use. The user can control the size of the stack with the **pthread_attr_setstacksize** subroutine and the size of the guard region with the **pthread_attr_setguardsize** subroutine.

Note: The soft limit on stack size imposed by the command `ulimit -s` applies only to the stack of the main thread of the application.

The following table indicates the maximum number of pthreads that could be created in a 32-bit process using a simple program which does nothing other than create pthreads in a loop using the NULL pthread attribute. In a real program, the actual numbers depend on other memory usage in the program. For a 64-bit process, the **ulimit** subroutine controls how many threads can be created. Therefore, the big data model is not necessary and in fact, can decrease the maximum number of threads.

Data Model	-bmaxdata	Maximum Pthreads
Small Data	n/a	1084
Big Data	0x10000000	2169
Big Data	0x20000000	4340
Big Data	0x30000000	6510
Big Data	0x40000000	8681
Big Data	0x50000000	10852
Big Data	0x60000000	13022
Big Data	0x70000000	15193
Big Data	0x80000000	17364

The **NUM_SPAREVP** environment variable can be set to control the number of spare virtual processors that are maintained by the library. It is not necessary to modify this variable. In some circumstances, applications that use only a few megabytes of memory can reduce memory overhead by setting the **NUM_SPAREVP** environment variable to a lower value. Typical settings include the number of CPUs on the system or the peak number of process threads. Setting this variable does not affect process performance. The default setting is 256.

Note: The **NUM_SPAREVP** environment variable is available only in AIX 5.1.

Example of a multithreaded program

The following short multithreaded program displays "Hello!" in both English and French for five seconds. Compile with **cc_r** or **xlC_r**. F

```
#include <pthread.h>    /* include file for pthreads - the 1st */
#include <stdio.h>      /* include file for printf()          */
#include <unistd.h>     /* include file for sleep()           */

void *Thread(void *string)
{
    while (1)
        printf("%s\n", (char *)string);
    pthread_exit(NULL);
}

int main()
{
    char *e_str = "Hello!";
    char *f_str = "Bonjour !";

    pthread_t e_th;
    pthread_t f_th;

    int rc;

    rc = pthread_create(&e_th, NULL, Thread, (void *)e_str);
    if (rc)
        exit(-1);
    rc = pthread_create(&f_th, NULL, Thread, (void *)f_str);
    if (rc)
        exit(-1);
    sleep(5);

    /* usually the exit subroutine should not be used
       see below to get more information */
    exit(0);
}
```

The initial thread (executing the **main** routine) creates two threads. Both threads have the same entry-point routine (the **Thread** routine), but a different parameter. The parameter is a pointer to the string that will be displayed.

Debugging a multithreaded program

The following tools are available to debug multithreaded programs:

- Application programmers can use the **dbx** command to perform debugging. Several subcommands are available for displaying thread-related objects, including **attribute**, **condition**, **mutex**, and **thread**.
- Kernel programmers can use the kernel debug program to perform debugging on kernel extensions and device drivers. The kernel debug program provides limited access to user threads, and primarily handles kernel threads. Several subcommands support multiple kernel threads and processors, including:
 - The **cpu** subcommand, which changes the current processor
 - The **ppd** subcommand, which displays per-processor data structures
 - The **thread** subcommand, which displays thread table entries
 - The **uthread** subcommand, which displays the **uthread** structure of a thread

For more information on the kernel debug program, see the *Kernel Extensions and Device Support Programming Concepts*.

Core File requirements of a multithreaded program

By default, processes do not generate a full core file. If an application must debug data in shared memory regions, particularly thread stacks, it is necessary to generate a full core dump. To generate full core file information, run the following command as root user:

```
chdev -l sys0 -a fullcore=true
```

Each individual pthread adds to the size of the generated core file. The amount of core file space that a pthread needs includes the stack size, which the user can control with the **pthread_attr_setstacksize** subroutine. For pthreads created with the NULL pthread attribute, each pthread in a 32-bit process adds 128 KB to the size of the core file, and each pthread in a 64-bit process adds 256 KB to the size of the core file.

Related concepts

[Threadsafe and threaded libraries in AIX](#)

This section describes the thread libraries in the AIX.

[Creating threads](#)

Thread creation differs from process creation in that no parent-child relation exists between threads.

[Scheduling threads](#)

Threads can be scheduled, and the threads library provides several facilities to handle and control the scheduling of threads.

[Developing multithreaded programs](#)

Developing multithreaded programs is similar to developing programs with multiple processes.

Developing programs also consists of compiling and debugging the code.

Developing multithreaded programs to examine and modify pthread library objects

The pthread debug library (**libpthreaddebug.a**) provides a set of functions that enable application developers to examine and modify pthread library objects.

This library can be used for both 32-bit applications and 64-bit applications. This library is threadsafe. The pthread debug library contains a 32-bit shared object and a 64-bit shared object.

The pthread debug library provides applications with access to the pthread library information. This includes information on pthreads, pthread attributes, mutexes, mutex attributes, condition variables, condition variable attributes, read/write locks, read/write lock attributes, and information about the state of the pthread library.

Note: All data (addresses, registers) returned by this library is in 64-bit format both for 64-bit and 32-bit applications. It is the application's responsibility to convert these values into 32-bit format for 32-bit applications. When debugging a 32-bit application, the top half of addresses and registers is ignored.

The pthread debug library does not report information on mutexes, mutex attributes, condition variables, condition variable attributes, read/write locks, and read/write lock attributes that have the pshared value of **PTHREAD_PROCESS_SHARED**.

Initialization

The application must initialize a pthread debug library session for each pthreaded process. The **pthread_session_init** function must be called from each pthreaded process after the process has been loaded. The pthread debug library supports one session for a single process. The application must assign a unique user identifier and pass it to the **pthread_session_init** function, which in turn assigns a unique session identifier that must be passed as the first parameter to all other **pthread debug library** functions, except **pthread_session_pthreaded** function, in return. Whenever the pthread debug library invokes a call back function, it will pass the unique application assigned user identifier back to the application. The **pthread_session_init** function checks the list of call back functions provided by the application, and initializes the session's data structures. Also, this function sets the session flags. An

application must pass the **PTHDB_FLAG_SUSPEND** flag to the **pthdb_session_init** Function. See the **pthdb_session_setflags** function for a full list of flags.

Call back functions

The pthread debug library uses the call back functions to obtain and write data, as well as to give storage management to the application. Required call back functions for an application are as follows:

read_data

Retrieves **pthread library** object information

alloc

Allocates memory in the pthread debug library

realloc

Reallocates memory in the pthread debug library

dealloc

Frees allocated memory in the pthread debug library

Optional call back functions for an application are as follows:

read_regs

Necessary only for the **pthdb_pthread_context** and **pthdb_pthread_setcontext** subroutines

write_data

Necessary only for the **pthdb_pthread_setcontext** subroutine

write_regs

Necessary only for **pthdb_pthread_setcontext** subroutine

Update function

Each time the application is stopped, after the session has been initialized, it is necessary to call the **pthdb_session_update** function. This function sets or resets the lists of pthreads, pthread attributes, mutexes, mutex attributes, condition variables, condition variable attributes, read/write locks, read/write lock attributes, pthread specific keys, and active keys. It uses call back functions to manage memory for the lists.

Context functions

The **pthdb_pthread_context** function obtains the context information, and the **pthdb_pthread_setcontext** function sets the context. The **pthdb_pthread_context** function obtains the context information of a pthread from either the kernel or the pthread data structure in the application's address space. If the pthread is not associated with a kernel thread, the context information saved by the pthread library is obtained. If a pthread is associated with a kernel thread, the information is obtained from the application using the call back functions. The application must determine if the kernel thread is in kernel mode or user mode and then to provide the correct information for that mode.

When a pthread with kernel thread is in kernel mode, you cannot get the full user mode context because the kernel does not save it in one place. The **getthrds** function can be used to obtain part of this information, because it always saves the user mode stack. The application can discover this by checking the **thrdsinfo64.ti_scount** structure. If this is non-zero, the user mode stack is available in the **thrdsinfo64.ti_ustk** structure. From the user mode stack, it is possible to determine the instruction address register (IAR) and the call back frames, but not the other register values. The **thrdsinfo64** structure is defined in **procinfo.h** file.

List functions

The pthread debug library maintains lists for pthreads, pthread attributes, mutexes, mutex attributes, condition variables, condition variables attributes, read/write locks, read/write lock attributes, pthread specific keys and active keys, each represented by a type-specific handle. The **pthdb_object** functions return the next handle in the appropriate list, where *object* is one of the following: **pthread**, **attr**, **mutex**,

mutexattr, cond, condattr, rwlock, rwlockattr or **key**. If the list is empty or the end of the list is reached, **PTHDB_INVALID_OBJECT** is reported, where *OBJECT* is one of the following: **PTHREAD, ATTR, MUTEX, MUTEXATTR, COND, CONDATTR, RWLOCK, RWLOCKATTR** or **KEY**.

Field functions

Detailed information about an object can be obtained by using the appropriate object member function, **pthdb_object_field**, where *object* is one of the following: **pthread, attr, mutex, mutexattr, cond, condattr, rwlock, rwlockattr** or **key** and where *field* is the name of a field of the detailed information for the object.

Customizing the session

The **pthdb_session_setflags** function allows the application to change the flags that customize the session. These flags control the number of registers that are read or written during context operations.

The **pthdb_session_flags** function obtains the current flags for the session.

Terminating the session

At the end of the session, the session data structures must be deallocated, and the session data must be deleted. This is accomplished by calling the **pthdb_session_destroy** function, which uses a call back function to deallocate the memory. All of the memory allocated by the **pthdb_session_init**, and **pthdb_session_update** functions will be deallocated.

The following example shows how an application can connect to the pthread debug library:

```
/* includes */
#include <thread.h>
#include <ys/pthdebug.h>
...
int my_read_data(pthdb_user_t user, pthdb_symbol_t symbols[],int count)
{
    int rc;

    rc=memcpy(buf,(void *)addr,len);
    if (rc==NULL) {
        fprintf(stderr,&odq;Error message\n&cdq;);
        return(1);
    }
    return(0);
}
int my_alloc(pthdb_user_t user, size_t len, void **bufp)
{
    *bufp=malloc(len);
    if(!*bufp) {
        fprintf(stderr,&odq;Error message\n&cdq;);
        return(1);
    }
    return(0);
}
int my_realloc(pthdb_user_t user, void *buf, size_t len, void **bufp)
{
    *bufp=realloc(buf,len);
    if(!*bufp) {
        fprintf(stderr,"Error message\n");
        return(1);
    }
    return(0);
}
int my_dealloc(pthdb_user_t user,void *buf)
{
    free(buf);
    return(0);
}
status()
{
    pthdb_callbacks_t callbacks =
```

```

        { NULL,
          my_read_data,
          NULL,
          NULL,
          NULL,
          my_alloc,
          my_realloc,
          my_dealloc,
          NULL
        };

...

rc=pthread_suspend_others_np();
if (rc!=0)
    deal with error

if (not initialized)
    rc=pthdb_session_init(user,exec_mode,PTHDB_SUSPEND|PTHDB_REGS,callbacks,
                        &session);
    if (rc!=PTHDB_SUCCESS)
        deal with error

rc=pthdb_session_update(session);
if (rc!=PTHDB_SUCCESS)
    deal with error

    retrieve pthread object information using the object list functions and
    the object field functions

...

rc=pthread_continue_others_np();
if (rc!=0)
    deal with error
}

...

main()
{
    ...
}

```

Developing multithreaded program debuggers

The pthread debug library (**libpthdebug.a**) provides a set of functions that allows developers to provide debug capabilities for applications that use the pthread library.

The pthread debug library is used to debug both 32-bit and 64-bit pthreaded applications. This library is used to debug targeted debug processes only. It can also be used to examine pthread information of its own application. This library can be used by a multithreaded debugger to debug a multithreaded application. Multithreaded debuggers are supported in the **libpthread.a** library, which is threadsafe. The pthread debug library contains a 32-bit shared object and a 64-bit shared object.

Debuggers using the ptrace facility must link to the 32-bit version of the library, because the ptrace facility is not supported in 64-bit mode. Debuggers using the /proc facility can link to either the 32-bit version or the 64-bit version of this library.

The pthread debug library provides debuggers with access to pthread library information. This includes information on pthreads, pthread attributes, mutexes, mutex attributes, condition variables, condition variable attributes, read/write locks, read/write lock attributes, and information about the state of the pthread library. This library also provides help with controlling the execution of pthreads.

Note: All data (addresses, registers) returned by this library is in 64-bit format both for 64-bit and 32-bit applications. It is the debugger's responsibility to convert these values into 32-bit format for 32-bit applications. When debugging a 32-bit application, the top half of addresses and registers is ignored.

The pthread debug library does not report mutexes, mutex attributes, condition variables, condition variable attributes, read/write locks, and read/write lock attributes that have the pshared value of PTHREAD_PROCESS_SHARED.

Initialization

The debugger must initialize a pthread debug library session for each debug process. This cannot be done until the pthread library has been initialized in the debug process. The **`pthdb_session_pthreaded`** function has been provided to tell the debugger when the pthread library has been initialized in the debug process. Each time the **`pthdb_session_pthreaded`** function is called, it checks to see if the pthread library has been initialized. If initialized, it returns `PTHDB_SUCCESS`. Otherwise, it returns `PTHDB_NOT_PTHREADED`. In both cases, it returns a function name that can be used to set a breakpoint for immediate notification that the pthread library has been initialized. Therefore, the **`pthdb_session_pthreaded`** function provides the following methods for determining when the pthread library has been initialized:

- The debugger calls the function each time the debug process stops, to see if the program that is being debugged is pthreaded.
- The debugger calls the function once and if the program that is being debugged is not pthreaded, sets a breakpoint to notify the debugger when the debug process is pthreaded.

After the debug process is pthreaded, the debugger must call the **`pthdb_session_init`** function, to initialize a session for the debug process. The pthread debug library supports one session for a single debug process. The debugger must assign a unique user identifier and pass it to **`pthdb_session_init`** which in turn will assign a unique session identifier which must be passed as the first parameter to all other pthread debug library functions, except **`pthdb_session_pthreaded`**, in return. Whenever the pthread debug library invokes a call back function, it will pass the unique debugger assigned user identifier back to the debugger. The **`pthdb_session_init`** function checks the list of call back functions provided by the debugger, and initializes the session's data structures. Also, this function sets the session flags.

Call back functions

The pthread debug library uses call back functions to do the following:

- Obtain addresses and data
- Write data
- Give storage management to the debugger
- Aid debugging of the pthread debug library

Update function

Each time the debugger is stopped, after the session has been initialized, it is necessary to call the **`pthdb_session_update`** function. This function sets or resets the lists of pthreads, pthread attributes, mutexes, mutex attributes, condition variables, condition variable attributes, read/write locks, read/write lock attributes, pthread specific keys, and active keys. It uses call back functions to manage memory for the lists.

Hold and unhold functions

Debuggers must support hold and unhold of threads for the following reasons:

- To allow a user to single step a single thread, it must be possible to hold one or more of the other threads.
- For users to continue through a subset of available threads, it must be possible to hold threads not in the set.

The following list of functions perform hold and unhold tasks:

- The **`pthdb_pthread_hold`** function sets the *hold state* of a pthread to hold.
- The **`pthdb_pthread_unhold`** function sets the *hold state* of a pthread to unhold.

Note: The **`pthdb_pthread_hold`** and **`pthdb_pthread_unhold`** functions must always be used, whether or not a pthread has a kernel thread.

- The **`pthdb_pthread_holdstate`** function returns the *hold state* of the pthread.

- The **pthdb_session_committed** function reports the function name of the function that is called after all of the hold and unhold changes are committed. A breakpoint can be placed at this function to notify the debugger when the hold and unhold changes have been committed.
- The **pthdb_session_stop_tid** function informs the **pthread debug library**, which informs the **pthread library** the thread ID (TID) of the thread that stopped the debugger.
- The **pthdb_session_commit_tid** function returns the list of kernel threads, one kernel thread at a time, that must be continued to commit the hold and unhold changes. This function must be called repeatedly, until PTHDB_INVALID_TID is reported. If the list of kernel threads is empty, it is not necessary to continue any threads for the commit operation.

The debugger can determine when all of the hold and unhold changes have been committed in the following ways:

- Before the commit operation (continuing all of the tids returned by the **pthdb_session_commit_tid** function) is started, the debugger can call the **pthdb_session_committed** function to get the function name and set a breakpoint. (This method can be done once for the life of the process.)
- Before the commit operation is started, the debugger calls the **pthdb_session_stop_tid** function with the TID of the thread that stopped the debugger. When the commit operation is complete, the pthread library ensures that the same TID is stopped as before the commit operation.

To hold or unhold pthreads, use the follow the following procedure, before continuing a group of pthreads or single-stepping a single pthread:

1. Use the **pthdb_thread_hold** and **pthdb_thread_unhold** functions to set up which pthreads will be held and which will be unheld.
2. Select the method that will determine when all of the hold and unhold changes have been committed.
3. Use the **pthdb_session_commit_tid** function to determine the list of TIDs that must be continued to commit the hold and unhold changes.
4. Continue the TIDs in the previous step and the thread that stopped the debugger.

The **pthdb_session_continue_tid** function allows the debugger to obtain the list of kernel threads that must be continued before it proceeds with single-stepping a single pthread or continuing a group of pthreads. This function must be called repeatedly, until PTHDB_INVALID_TID is reported. If the list of kernel threads is not empty, the debugger must continue these kernel threads along with the others that it is explicitly interested in. The debugger is responsible for parking the stop thread and continuing the stop thread. The stop thread is the thread that caused the debugger to be entered.

Context functions

The **pthdb_thread_context** function obtains the context information and the **pthdb_thread_setcontext** function sets the context. The **pthdb_thread_context** function obtains the context information of a pthread from either the kernel or the pthread data structure in the debug process's address space. If the pthread is not associated with a kernel thread, the context information saved by the pthread library is obtained. If a pthread is associated with a kernel thread, the information is obtained from the debugger using call backs. It is the debugger's responsibility to determine if the kernel thread is in kernel mode or user mode and then to provide the correct information for that mode.

When a pthread with kernel thread is in kernel mode, you cannot get the full user mode context because the kernel does not save it in one place. The **getthrds** function can be used to obtain part of this information, because it always saves the user mode stack. The debugger can discover this by checking the **thrdsinfo64.ti_scount** structure. If this is non-zero, the user mode stack is available in the **thrdsinfo64.ti_ustk** structure. From user mode stack, it is possible to determine the instruction address register (IAR) and the call back frames, but not the other register values. The **thrdsinfo64** structure is defined in **procinfo.h** file.

List functions

The pthread debug library maintains lists for pthreads, pthread attributes, mutexes, mutex attributes, condition variables, condition variables attributes, read/write locks, read/write lock attributes, pthread specific keys and active keys, each represented by a type-specific handle. The **pthdb_object** functions return the next handle in the appropriate list, where *object* is one of the following: **pthread**, **attr**, **mutex**,

mutexattr, cond, condattr, rwlock, rwlockattr or **key**. If the list is empty or the end of the list is reached, **PTHDB_INVALID_object** is reported, where *object* is one of the following: **PTHREAD, ATTR, MUTEX, MUTEXATTR, COND, CONDATTR, RWLOCK, RWLOCKATTR** or **KEY**.

Field Functions

Detailed information about an object can be obtained by using the appropriate object member function, **pthdb_object_field**, where *object* is one of the following: **pthread, attr, mutex, mutexattr, cond, condattr, rwlock, rwlockattr** or **key** and where *field* is the name of a field of the detailed information for the object.

Customizing the session

The **pthdb_session_setflags** function allows the debugger to change the flags that customize the session. These flags control the number of registers that are read or written to during context operations, and to control the printing of debug information.

The **pthdb_session_flags** function obtains the current flags for the session.

Terminating the session

At the end of the debug session, the session data structures must be deallocated, and the session data must be deleted. This is accomplished by calling the **pthdb_session_destroy** function, which uses a call back function to deallocate the memory. All of the memory allocated by the **pthdb_session_init** and **pthdb_session_update** functions will be deallocated.

Example of hold/unhold functions

The following pseudocode example shows how the debugger uses the hold/unhold code:

```
/* includes */
#include <sys/ptdebug.h>

main()
{
    tid_t stop_tid; /* thread which stopped the process */
    pthdb_user_t user = <unique debugger value>;
    pthdb_session_t session; /* <unique library value> */
    pthdb_callbacks_t callbacks = <callback functions>;
    char *pthreaded_symbol=NULL;
    char *committed_symbol;
    int pthreaded = 0;
    int pthdb_init = 0;
    char *committed_symbol;

    /* fork/exec or attach to the program that is being debugged */

    /* the program that is being debugged uses ptrace()/ptracex() with PT_TRACE_ME */

    while (/* waiting on an event */)
    {
        /* debugger waits on the program that is being debugged */

        if (pthreaded_symbol==NULL) {
            rc = pthdb_session_pthreaded(user, &callbacks, pthreaded_symbol);
            if (rc == PTHDB_NOT_PTHREADED)
            {
                /* set breakpoint at pthreaded_symbol */
            }
            else
                pthreaded=1;
        }
        if (pthreaded == 1 && pthdb_init == 0) {
            rc = pthdb_session_init(user, &session, PEM_32BIT, flags, &callbacks);
            if (rc)
                /* handle error and exit */
                pthdb_init=1;
        }

        rc = pthdb_session_update(session)
        if ( rc != PTHDB_SUCCESS)
            /* handle error and exit */

        while (/* accepting debugger commands */)
        {
            switch (/* debugger command */)
            {
                ...
                case DB_HOLD:
                    /* regardless of pthread with or without kernel thread */
                    rc = pthdb_pthread_hold(session, pthread);
            }
        }
    }
}
```

```

        if (rc)
            /* handle error and exit */
        case DB_UNHOLD:
            /* regardless of pthread with or without kernel thread */
            rc = pthread_unhold(session, pthread);
            if (rc)
                /* handle error and exit */
        case DB_CONTINUE:
            /* unless we have never held threads for the life */
            /* of the process */
            if (pthreaded)
            {
                /* debugger must handle list of any size */
                struct pthread commit_tids;
                int commit_count = 0;
                /* debugger must handle list of any size */
                struct pthread continue_tids;
                int continue_count = 0;

            rc = pthread_session_committed(session, committed_symbol);
            if (rc != PTHDB_SUCCESS)
                /* handle error */
                /* set break point at committed_symbol */

                /* gather any tids necessary to commit hold/unhold */
                /* operations */
                do
                {
                    rc = pthread_session_commit_tid(session,
                                                    &commit_tids.th[commit_count++]);
                    if (rc != PTHDB_SUCCESS)
                        /* handle error and exit */
                } while (commit_tids.th[commit_count - 1] != PTHDB_INVALID_TID);

                /* set up thread which stopped the process to be */
                /* parked using the stop_park function*/

            if (commit_count > 0) {
                rc = ptrace(PTT_CONTINUE, stop_tid, stop_park, 0,
                            &commit_tids);
                if (rc)
                    /* handle error and exit */

                /* wait on process to stop */
            }

            /* gather any tids necessary to continue */
            /* interesting threads */
            do
            {
                rc = pthread_session_continue_tid(session,
                                                    &continue_tids.th[continue_count++]);
                if (rc != PTHDB_SUCCESS)
                    /* handle error and exit */
            } while (continue_tids.th[continue_count - 1] != PTHDB_INVALID_TID);

            /* add interesting threads to continue_tids */

            /* set up thread which stopped the process to be parked */
            /* unless it is an interesting thread */

            rc = ptrace(PTT_CONTINUE, stop_tid, stop_park, 0,
                        &continue_tids);
            if (rc)
                /* handle error and exit */
            }
        case DB_EXIT:
            rc = pthread_session_destroy(session);
            /* other clean up code */
            exit(0);
            ...
        }
    }
}
exit(0);
}

```

Benefits of threads

Multithreaded programs can improve performance compared to traditional parallel programs that use multiple processes. Furthermore, improved performance can be obtained on multiprocessor systems using threads.

Managing threads

Managing threads; that is, creating threads and controlling their execution, requires fewer system resources than managing processes. Creating a thread, for example, only requires the allocation of the

thread's private data area, usually 64 KB, and two system calls. Creating a process is far more expensive, because the entire parent process addressing space is duplicated.

The threads library API is also easier to use than the library for managing processes. Thread creation requires only the **pthread_create** subroutine.

Inter-thread communications

Inter-thread communication is far more efficient and easier to use than inter-process communication. Because all threads within a process share the same address space, they need not use shared memory. Protect shared data from concurrent access by using mutexes or other synchronization tools.

Synchronization facilities provided by the threads library ease implementation of flexible and powerful synchronization tools. These tools can replace traditional inter-process communication facilities, such as message queues. Pipes can be used as an inter-thread communication path.

Multiprocessor systems

On a multiprocessor system, multiple threads can concurrently run on multiple CPUs. Therefore, multithreaded programs can run much faster than on a uniprocessor system. They can also be faster than a program using multiple processes, because threads require fewer resources and generate less overhead. For example, switching threads in the same process can be faster, especially in the M:N library model where context switches can often be avoided. Finally, a major advantage of using threads is that a single multithreaded program will work on a uniprocessor system, but can naturally take advantage of a multiprocessor system, without recompiling.

Limitations

Multithreaded programming is useful for implementing parallelized algorithms using several independent entities. However, there are some cases where multiple processes should be used instead of multiple threads.

Many operating system identifiers, resources, states, or limitations are defined at the process level and, thus, are shared by all threads in a process. For example, user and group IDs and their associated permissions are handled at process level. Programs that need to assign different user IDs to their programming entities need to use multiple processes, instead of a single multithreaded process. Other examples include file-system attributes, such as the current working directory, and the state and maximum number of open files. Multithreaded programs may not be appropriate if these attributes are better handled independently. For example, a multi-processed program can let each process open a large number of files without interference from other processes.

Related concepts

[Threadsafe and threaded libraries in AIX](#)

This section describes the thread libraries in the AIX.

lex and yacc program information

For a program to receive input, either interactively or in a batch environment, you must provide another program or a routine to receive the input. Complicated input requires additional code to break the input into pieces that mean something to the program.

You can use the **lex** and **yacc** commands to develop this type of input program.

The **lex** command generates a lexical analyzer program that analyzes input and breaks it into tokens, such as numbers, letters, or operators. The tokens are defined by grammar rules set up in the **lex** specification file. The **yacc** command generates a parser program that analyzes input using the tokens identified by the lexical analyzer (generated by the **lex** command and stored in the **lex** specification file) and performs specified actions, such as flagging improper syntax. Together these commands generate a lexical analyzer and parser program for interpreting input and output handling.

Related information

[printf](#)

[ed](#)
[ex](#)
[sed](#)
[yacc](#)

Generating a lexical analyzer with the `lex` command

The `lex` command helps write a C language program that can receive and translate character-stream input into program actions.

To use the `lex` command, you must supply or write a specification file that contains:

Extended regular expressions

Character patterns that the generated lexical analyzer recognizes.

Action statements

C language program fragments that define how the generated lexical analyzer reacts to extended regular expressions it recognizes.

The `lex` command generates a C language program that can analyze an input stream using information in the specification file. The `lex` command then stores the output program in a `lex.yy.c` file. If the output program recognizes a simple, one-word input structure, you can compile the `lex.yy.c` output file with the following command to produce an executable lexical analyzer:

```
cc lex.yy.c -ll
```

However, if the lexical analyzer must recognize more complex syntax, you can create a parser program to use with the output file to ensure proper handling of any input.

You can move a `lex.yy.c` output file to another system if it has a C compiler that supports the `lex` library functions.

The compiled lexical analyzer performs the following functions:

- Reads an input stream of characters.
- Copies the input stream to an output stream.
- Breaks the input stream into smaller strings that match the extended regular expressions in the `lex` specification file.
- Executes an action for each extended regular expression that it recognizes. These actions are C language program fragments in the `lex` specification file. Each action fragment can call actions or subroutines outside of itself.

The lexical analyzer generated by the `lex` command uses an analysis method called a *deterministic finite-state automaton*. This method provides for a limited number of conditions in which the lexical analyzer can exist, along with the rules that determine the state of the lexical analyzer.

The automaton allows the generated lexical analyzer to look ahead more than one or two characters in an input stream. For example, suppose you define two rules in the `lex` specification file: one looks for the string `ab` and the other looks for the string `abcdefg`. If the lexical analyzer receives an input string of `abcdefh`, it reads characters to the end of the input string before determining that it does not match the string `abcdefg`. The lexical analyzer then returns to the rule that looks for the string `ab`, decides that it matches part of the input, and begins trying to find another match using the remaining input `cdefh`.

Compiling the lexical analyzer

To compile a `lex` program, do the following:

1. Use the `lex` program to change the specification file into a C language program. The resulting program is in the `lex.yy.c` file.
2. Use the `cc` command with the `-ll` flag to compile and link the program with a library of `lex` subroutines. The resulting executable program is in the `a.out` file.

For example, if the **lex** specification file is called **lextest**, enter the following commands:

```
lex lextest
cc lex.yy.c -ll
```

Related concepts

[Tools and utilities](#)

This section provides an overview of the tools and utilities that you can use to develop C compiled language programs.

[Creating a parser with the yacc program](#)

The **yacc** program creates parsers that define and enforce structure for character input to a computer program.

Using the lex program with the yacc program

You can also use the **lex** program with a parser generator, such as the **yacc** command. The **yacc** command generates a program, called a *parser*, that analyzes the construction of more than one-word input.

This parser program operates well with the lexical analyzers that the **lex** command generates. The parsers recognize many types of grammar with no regard to context. These parsers need a preprocessor to recognize input tokens such as the preprocessor that the **lex** command produces.

The **lex** program recognizes only extended regular expressions and formats them into character packages called *tokens*, as specified by the input file. When using the **lex** program to make a lexical analyzer for a parser, the lexical analyzer (created from the **lex** command) partitions the input stream. The parser (from the **yacc** command) assigns structure to the resulting pieces. You can also use other programs along with the programs generated by either the **lex** or **yacc** commands.

A token is the smallest independent unit of meaning as defined by either the parser or the lexical analyzer. A token can contain data, a language keyword, an identifier, or other parts of a language syntax.

The **yacc** program looks for a lexical analyzer subroutine named **yylex**, which is generated by the **lex** command. Normally, the default main program in the **lex** library calls the **yylex** subroutine. However, if the **yacc** command is installed and its main program is used, the **yacc** program calls the **yylex** subroutine. In this case, where the appropriate token value is returned, each **lex** program rule should end with the following:

```
return(token);
```

The **yacc** command assigns an integer value to each token defined in the **yacc** grammar file through a **#define** preprocessor statement. The lexical analyzer must have access to these macros to return the tokens to the parser. Use the **yacc -d** option to create a **y.tab.h** file, and include the **y.tab.h** file in the **lex** specification file by adding the following lines to the definition section of the **lex** specification file:

```
%{
#include "y.tab.h"
%}
```

Alternately, you can include the **lex.yy.c** file in the **yacc** output file by adding the following line after the second **%%** (percent sign, percent sign) delimiter in the **yacc** grammar file:

```
#include "lex.yy.c"
```

The **yacc** library should be loaded before the **lex** library to obtain a main program that invokes the **yacc** parser. You can generate **lex** and **yacc** programs in either order.

Extended regular expressions in the `lex` command

Specifying extended regular expressions in a `lex` specification file is similar to methods used in the `sed` or `ed` commands.

An extended regular expression specifies a set of strings to be matched. The expression contains both text characters and operator characters. Text characters match the corresponding characters in the strings being compared. Operator characters specify repetitions, choices, and other features.

Numbers and letters of the alphabet are considered text characters. For example, the extended regular expression `integer` matches the string `integer`, and the expression `a57D` looks for the string `a57D`.

Operators

The following list describes how operators are used to specify extended regular expressions:

Character

Matches the character *Character*.

Example: `a` matches the literal character `a`; `b` matches the literal character `b`, and `c` matches the literal character `c`.

"String"

Matches the string enclosed within quotes, even if the string includes an operator.

Example: To prevent the `lex` command from interpreting `$` (dollar sign) as an operator, enclose the symbol in quotes.

\Character or \Digits

Escape character. When preceding a character class operator used in a string, the `\` character indicates that the operator symbol represents a literal character rather than an operator. Valid escape sequences include:

\a

Alert

\b

Backspace

\f

Form-feed

\n

New line character (Do not use the actual new line character in an expression.)

\r

Return

\t

Tab

\v

Vertical tab

Backslash

\Digits

The character whose encoding is represented by the one-digit, two-digit, or three-digit octal integer specified by the *Digits* string.

\xDigits

The character whose encoding is represented by the sequence of hexadecimal characters specified by the *Digits* string.

When the `\` character precedes a character that is not in the preceding list of escape sequences, the `lex` command interprets the character literally.

Example: `\c` is interpreted as the `c` character unchanged, and `[\^abc]` represents the class of characters that includes the characters `^abc`.

Note: Never use `\0` or `\x0` in the **lex** command.

[List]

Matches any one character in the enclosed range (`[x-y]`) or the enclosed list (`[xyz]`) based on the locale in which the **lex** command is invoked. All operator symbols, with the exception of the following, lose their special meaning within a bracket expression: - (dash), ^ (caret), and \ (backslash).

Example: `[abc-f]` matches a, b, c, d, e, or f in the `en_US` locale.

[:Class:]

Matches any of the characters belonging to the character class specified between the `[:]` delimiters as defined in the `LC_TYPE` category in the current locale. The following character class names are supported in all locales:

alnum	cntrl	lower	space
alpha	digit	print	upper
blank	graph	punct	xdigit

The **lex** command also recognizes user-defined character class names. The `[:]` operator is valid only in a `[]` expression.

Example: `[[:alpha:]]` matches any character in the **alpha** character class in the current locale, but `[:alpha:]` matches only the characters `:a,l,p,` and `h`.

[. CollatingSymbol .]

Matches the collating symbol specified within the `[.]` delimiters as a single character. The `[.]` operator is valid only in a `[]` expression. The collating symbol must be a valid collating symbol for the current locale.

Example: `[[. ch .]]` matches c and h together while `[ch]` matches c or h.

[=CollatingElement=]

Matches the collating element specified within the `[=]` delimiters and all collating elements belonging to its equivalence class. The `[=]` operator is valid only in a `[]` expression.

Example: If w and v belong to the same equivalence class, `[[=w=]]` is the same as `[wv]` and matches w or v. If w does not belong to an equivalence class, then `[[=w=]]` matches w only.

[^Character]

Matches any character except the one following the ^ (caret) symbol. The resultant character class consists solely of single-byte characters. The character following the ^ symbol can be a multibyte character. However, for this operator to match multibyte characters, you must set **%h** and **%m** to greater than zero in the definitions section.

Example: `[^c]` matches any character except c.

CollatingElement-CollatingElement

In a character class, indicates a range of characters within the collating sequence defined for the current locale. Ranges must be in ascending order. The ending range point must collate equal to or higher than the starting range point. Because the range is based on the collating sequence of the current locale, a given range may match different characters, depending on the locale in which the **lex** command was invoked.

Expression?

Matches either zero or one occurrence of the expression immediately preceding the ? operator.

Example: `ab?c` matches either `ac` or `abc`.

Period character (.)

Matches any character except the new line character. In order for the period character (`.`) to match multi-byte characters, **%z** must be set to greater than 0 in the definitions section of the **lex** specification file. If **%z** is not set, the period character (`.`) matches single-byte characters only.

Expression*

Matches zero or more occurrences of the expression immediately preceding the * operator. For example, a* is any number of consecutive a characters, including zero. The usefulness of matching zero occurrences is more obvious in complicated expressions.

Example: The expression, [A-Za-z][A-Za-z0-9]* indicates all alphanumeric strings with a leading alphabetic character, including strings that are only one alphabetic character. You can use this expression for recognizing identifiers in computer languages.

Expression+

Matches one or more occurrences of the pattern immediately preceding the + operator.

Example: a+ matches one or more instances of a. Also, [a-z]+ matches all strings of lowercase letters.

Expression | Expression

Indicates a match for the expression that precedes or follows the | (pipe) operator.

Example: ab|cd matches either ab or cd.

(Expression)

Matches the expression in the parentheses. The () (parentheses) operator is used for grouping and causes the expression within parentheses to be read into the **yytext** array. A group in parentheses can be used in place of any single character in any other pattern.

Example: (ab|cd+)?(ef)* matches such strings as abefef, efefef, cdef, or cddd; but not abc, abcd, or abcdef.

^Expression

Indicates a match only when *Expression* is at the beginning of the line and the ^ (caret) operator is the first character in an expression.

Example: ^h matches an h at the beginning of a line.

Expression\$

Indicates a match only when *Expression* is at the end of the line and the \$ (dollar sign) operator is the last character in an expression.

Example: h\$ matches an h at the end of a line.

Expression1 / Expression2

Indicates a match only if *Expression2* immediately follows *Expression1*. The / (slash) operator reads only the first expression into the **yytext** array.

Example: ab/cd matches the string ab, but only if followed by cd, and then reads ab into the **yytext** array.

Note: Only one / trailing context operator can be used in a single extended regular expression. The ^ (caret) and \$ (dollar sign) operators cannot be used in the same expression with the / operator as they indicate special cases of trailing context.

{DefinedName}

Matches the name as you defined it in the definitions section.

Example: If you defined D to be numerical digits, {D} matches all numerical digits.

{Number1, Number2}

Matches *Number1* to *Number2* occurrences of the pattern immediately preceding it. The expressions {*Number*} and {*Number*, } are also allowed and match exactly *Number* occurrences of the pattern preceding the expression.

Example: xyz{2, 4} matches either xyzxyz, xyzxyzxyz, or xyzxyzxyzxyz. This differs from the +, * and ? operators in that these operators match only the character immediately preceding them. To match only the character preceding the interval expression, use the grouping operator. For example, xy(z{2, 4}) matches xyzz, xyzzz or xyzzzz.

<StartCondition>

Executes the associated action only if the lexical analyzer is in the indicated start condition

Example: If being at the beginning of a line is start condition ONE, then the ^ (caret) operator equals the expression <ONE>.

To use the operator characters as text characters, use one of the escape sequences: " " (double quotation marks) or \ (backslash). The " " operator indicates that what is enclosed is text. Thus, the following example matches the string xyz++:

```
xyz"++"
```

A portion of a string can be quoted. Quoting an ordinary text character has no effect. For example, the following expression is equivalent to the previous example:

```
"xyz++"
```

To ensure that text is interpreted as text, quote all characters that are not letters or numbers.

Another way to convert an operator character to a text character is to put a \ (backslash) character before the operator character. For example, the following expression is equivalent to the preceding examples:

```
xyz\+\+
```

Related concepts

lex program start conditions

A rule may be associated with any start condition.

Passing code to the generated lex program

The **lex** command passes C code, unchanged, to the lexical analyzer in the following circumstances:

- Lines beginning with a blank or tab in the definitions section, or at the start of the rules section before the first rule, are copied into the lexical analyzer. If the entry is in the definitions section, it is copied to the external declaration area of the **lex.yy.c** file. If the entry is at the start of the rules section, the entry is copied to the local declaration area of the **yylex** subroutine in the **lex.yy.c** file.
- Lines that lie between delimiter lines containing only %{ (percent sign, left brace) and %} (percent sign, right brace) either in the definitions section or at the start of the rules section are copied into the lexical analyzer in the same way as lines beginning with a blank or tab.
- Any lines occurring after the second %% (percent sign, percent sign) delimiter are copied to the lexical analyzer without format restrictions.

Defining lex substitution strings

You can define string macros that the **lex** program expands when it generates the lexical analyzer.

Define them before the first %% delimiter in the **lex** specification file. Any line in this section that begins in column 1 and that does not lie between %{ and %} defines a **lex** substitution string. Substitution string definitions have the following general format:

```
name                translation
```

where name and translation are separated by at least one blank or tab, and the specified name begins with a letter. When the **lex** program finds the string defined by name enclosed in {} (braces) in the rules part of the specification file, it changes that name to the string defined in translation and deletes the braces.

For example, to define the names D and E, put the following definitions before the first %% delimiter in the specification file:

```
D      [0-9]
E      [DEde] [-+]{D}+
```

Then, use these names in the rules section of the specification file to make the rules shorter:

```
{D}+           printf("integer");
{D}+"."{E}?   |
{D}*"."{D}+({E})? |
{D}+{E}      printf("real");
```

You can also include the following items in the definitions section:

- Character set table
- List of start conditions
- Changes to size of arrays to accommodate larger source programs

lex library

The **lex** library contains the following subroutines:

Subroutine	Description
main()	Invokes the lexical analyzer by calling the yylex subroutine.
yywrap()	Returns the value 1 when the end of input occurs.
yyomore()	Appends the next matched string to the current value of the yytext array rather than replacing the contents of the yytext array.
yyless(int n)	Retains <i>n</i> initial characters in the yytext array and returns the remaining characters to the input stream.
yyreject()	Allows the lexical analyzer to match multiple rules for the same input string. (The yyreject subroutine is called when the special action REJECT is used.)

Some of the **lex** subroutines can be substituted by user-supplied routines. For example, the **lex** command supports user-supplied versions of the **main** and **yywrap** subroutines. The library versions of these routines, provided as a base, are as follows:

main subroutine

```
#include <stdio.h>
#include <locale.h>
main() {
    setlocale(LC_ALL, "");
    yylex();
    exit(0);
}
```

yywrap subroutine

```
yywrap() {
    return(1);
}
```

The **yyomore**, **yyless**, and **yyreject** subroutines are available only through the **lex** library. However, these subroutines are required only when used in **lex** command actions.

Actions taken by the lexical analyzer

When the lexical analyzer matches one of the extended regular expressions in the rules section of the specification file, it executes the *action* that corresponds to the extended regular expression. Without sufficient rules to match all strings in the input stream, the lexical analyzer copies the input to standard output. Therefore, do not create a rule that only copies the input to the output. The default output can help find gaps in the rules.

When using the **lex** command to process input for a parser that the **yacc** command produces, provide rules to match all input strings. Those rules must generate output that the **yacc** command can interpret.

Null action

To ignore the input associated with an extended regular expression, use a ; (C language null statement) as an action. The following example ignores the three spacing characters (blank, tab, and new line):

```
[ \t\n] ;
```

Same as next action

To avoid repeatedly writing the same action, use the | (pipe symbol). This character indicates that the action for this rule is the same as the action for the next rule. For instance, the previous example that ignores blank, tab, and new line characters can also be written as follows:

```
" " |
"\t" |
"\n" ;
```

The quotation marks that surround \n and \t are not required.

Printing a matched string

To determine what text matched an expression in the rules section of the specification file, you can include a C language **printf** subroutine call as one of the actions for that expression. When the lexical analyzer finds a match in the input stream, the program puts the matched string into the external character (**char**) and wide character (**wchar_t**) arrays, called **yytext** and **yywtext**, respectively. For example, you can use the following rule to print the matched string:

```
[a-z]+      printf("%s",yytext);
```

The C language **printf** subroutine accepts a format argument and data to be printed. In this example, the arguments to the **printf** subroutine have the following meanings:

%s

A symbol that converts the data to type string before printing

%S

A symbol that converts the data to wide character string (**wchar_t**) before printing

yytext

The name of the array containing the data to be printed

yywtext

The name of the array containing the multibyte type (**wchar_t**) data to be printed

The **lex** command defines **ECHO**; as a special action to print the contents of **yytext**. For example, the following two rules are equivalent:

```
[a-z]+      ECHO;
[a-z]+      printf("%s",yytext);
```

You can change the representation of **yytext** by using either **%array** or **%pointer** in the definitions section of the **lex** specification file, as follows:

%array

Defines **yytext** as a null-terminated character array. This is the default action.

%pointer

Defines **yytext** as a pointer to a null-terminated character string.

Finding the length of a matched string

To find the number of characters that the lexical analyzer matched for a particular extended regular expression, use the **yytext** or the **yywtext** external variables.

yytext

Tracks the number of bytes that are matched.

yywle^{ng}

Tracks the number of wide characters in the matched string. Multibyte characters have a length greater than 1.

To count both the number of words and the number of characters in words in the input, use the following action:

```
[a-zA-Z]+      {words++;chars += yywleng;} 
```

This action totals the number of characters in the words matched and puts that number in `chars`.

The following expression finds the last character in the string matched:

```
yytext[yywleng-1]
```

Matching strings within strings

The **lex** command partitions the input stream and does not search for all possible matches of each expression. Each character is accounted for only once. To override this choice and search for items that may overlap or include each other, use the **REJECT** action. For example, to count all instances of `she` and `he`, including the instances of `he` that are included in `she`, use the following action:

```
she      {s++; REJECT;}
he       {h++;}
\n       |
.        ;
```

After counting the occurrences of `she`, the **lex** command rejects the input string and then counts the occurrences of `he`. Because `he` does not include `she`, a **REJECT** action is not necessary on `he`.

Adding results to the yytext array

Typically, the next string from the input stream overwrites the current entry in the **yytext** array. If you use the **yymore** subroutine, the next string from the input stream is added to the end of the current entry in the **yytext** array.

For example, the following lexical analyzer looks for strings:

```
%s instring
%%
<INITIAL>\"    { /* start of string */
    BEGIN instring;
    yymore();
}
<instring>\"    { /* end of string */
    printf("matched %s\n", yytext);
    BEGIN INITIAL;
}
<instring>.     {
    yymore();
}
<instring>\n    {
    printf("Error, new line in string\n");
    BEGIN INITIAL;
}
```

Even though a string may be recognized by matching several rules, repeated calls to the **yymore** subroutine ensure that the **yytext** array will contain the entire string.

Returning characters to the input stream

To return characters to the input stream, use the following call:

```
yyless(n)
```

where `n` is the number of characters of the current string to keep. Characters in the string beyond this number are returned to the input stream. The **yyless** subroutine provides the same type of look-ahead function that the `/` (slash) operator uses, but it allows more control over its usage.

Use the **yyless** subroutine to process text more than once. For example, when parsing a C language program, an expression such as `x=-a` is difficult to understand. Does it mean *x is equal to minus a*, or is it an older representation of `x -= a`, which means *decrease x by the value of a*? To treat this expression as *x is equal to minus a*, but print a warning message, use a rule such as the following:

```
--[a-zA-Z]      {
                  printf("Operator (=-) ambiguous\n");
                  yyless(yylen-1);
                  ... action for = ...
                }
```

Input/Output subroutines

The **lex** program allows a program to use the following input/output (I/O) subroutines:

input()

Returns the next input character

output(c)

Writes the character `c` on the output

unput(c)

Pushes the character `c` back onto the input stream to be read later by the **input** subroutine

winput()

Returns the next multibyte input character

woutput(C)

Writes the multibyte character `C` back onto the output stream

wunput(C)

Pushes the multibyte character `C` back onto the input stream to be read by the **winput** subroutine

The **lex** program provides these subroutines as macro definitions. The subroutines are coded in the **lex.yy.c** file. You can override them and provide other versions.

The **winput**, **wunput**, and **woutput** macros are defined to use the **yywinput**, **yywunput**, and **yywoutput** subroutines. For compatibility, the **yy** subroutines subsequently use the **input**, **unput**, and **output** subroutine to read, replace, and write the necessary number of bytes in a complete multibyte character.

These subroutines define the relationship between external files and internal characters. If you change the subroutines, change them all in the same way. These subroutines should follow these rules:

- All subroutines must use the same character set.
- The **input** subroutine must return a value of 0 to indicate end of file.
- Do not change the relationship of the **unput** subroutine to the **input** subroutine or the look-ahead functions will not work.

The **lex.yy.c** file allows the lexical analyzer to back up a maximum of 200 characters.

To read a file containing nulls, create a different version of the **input** subroutine. In the normal version of the **input** subroutine, the returned value of 0 (from the null characters) indicates the end of file and ends the input.

Character set

The lexical analyzers that the **lex** command generates process character I/O through the **input**, **output**, and **unput** subroutines. Therefore, to return values in the **yytext** subroutine, the **lex** command uses the character representation that these subroutines use. Internally, however, the **lex** command represents each character with a small integer. When using the standard library, this integer is the value of the bit pattern the computer uses to represent the character. Normally, the letter *a* is represented in the same form as the character constant *a*. If you change this interpretation with different I/O subroutines, put a translation table in the definitions section of the specification file. The translation table begins and ends with lines that contain only the following entries:

```
%T
```

The translation table contains additional lines that indicate the value associated with each character. For example:

```
%T
{integer}      {character string}
{integer}      {character string}
{integer}      {character string}
%T
```

End-of-file processing

When the lexical analyzer reaches the end of a file, it calls the **yywrap** library subroutine, which returns a value of 1 to indicate to the lexical analyzer that it should continue with normal wrap-up at the end of input.

However, if the lexical analyzer receives input from more than one source, change the **yywrap** subroutine. The new function must get the new input and return a value of 0 to the lexical analyzer. A return value of 0 indicates that the program should continue processing.

You can also include code to print summary reports and tables when the lexical analyzer ends in a new version of the **yywrap** subroutine. The **yywrap** subroutine is the only way to force the **yylex** subroutine to recognize the end of input.

lex program start conditions

A rule may be associated with any start condition.

However, the **lex** program recognizes the rule only when in that associated start condition. You can change the current start condition at any time.

Define start conditions in the *definitions* section of the specification file by using a line in the following form:

```
%Start name1 name2
```

where *name1* and *name2* define names that represent conditions. There is no limit to the number of conditions, and they can appear in any order. You can also shorten the word *Start* to *s* or *S*.

When using a start condition in the rules section of the specification file, enclose the name of the start condition in <> (less than, greater than) symbols at the beginning of the rule. The following example defines a rule, *expression*, that the **lex** program recognizes only when the **lex** program is in start condition *name1*:

```
<name1> expression
```

To put the **lex** program in a particular start condition, execute the action statement in the action part of a rule; for instance, *BEGIN* in the following line:

```
BEGIN name1;
```

This statement changes the start condition to *name1*.

To resume the normal state, enter:

```
BEGIN 0;
```

or

```
BEGIN INITIAL;
```

where *INITIAL* is defined to be 0 by the **lex** program. *BEGIN 0;* resets the **lex** program to its initial condition.

The **lex** program also supports exclusive start conditions specified with **%x** (percent sign, lowercase x) or **%X** (percent sign, uppercase X) operator followed by a list of exclusive start names in the same format as

regular start conditions. Exclusive start conditions differ from regular start conditions in that rules that do not begin with a start condition are not active when the lexical analyzer is in an exclusive start state. For example:

```
%s      one
%x      two
%%
abc     {printf("matched ");ECHO;BEGIN one;}
<one>def  printf("matched ");ECHO;BEGIN two;}
<two>ghi  {printf("matched ");ECHO;BEGIN INITIAL;}
```

In start state one in the preceding example, both `abc` and `def` can be matched. In start state two, only `ghi` can be matched.

Related concepts

[Extended regular expressions in the `lex` command](#)

Specifying extended regular expressions in a **lex** specification file is similar to methods used in the **sed** or **ed** commands.

Creating a parser with the `yacc` program

The **yacc** program creates parsers that define and enforce structure for character input to a computer program.

To use this program, you must supply the following inputs:

grammar file

A source file that contains the specifications for the language to recognize. This file also contains the **main**, **yyerror**, and **yylex** subroutines. You must supply these subroutines.

main

A C language subroutine that, as a minimum, contains a call to the **yyparse** subroutine generated by the **yacc** program. A limited form of this subroutine is available in the **yacc** library.

yyerror

A C language subroutine to handle errors that can occur during parser operation. A limited form of this subroutine is available in the **yacc** library.

yylex

A C language subroutine to perform lexical analysis on the input stream and pass tokens to the parser. You can use the **lex** command to generate this lexical analyzer subroutine.

When the **yacc** command gets a specification, it generates a file of C language functions called **y.tab.c**. When compiled using the **cc** command, these functions form the **yyparse** subroutine and return an integer. When called, the **yyparse** subroutine calls the **yylex** subroutine to get input tokens. The **yylex** subroutine continues providing input until either the parser detects an error or the **yylex** subroutine returns an end-marker token to indicate the end of operation. If an error occurs and the **yyparse** subroutine cannot recover, it returns a value of 1 to the **main** subroutine. If it finds the end-marker token, the **yyparse** subroutine returns a value of 0 to the **main** subroutine.

Related concepts

[Generating a lexical analyzer with the `lex` command](#)

The **lex** command helps write a C language program that can receive and translate character-stream input into program actions.

The `yacc` grammar file

To use the **yacc** command to generate a parser, provide it with a grammar file that describes the input data stream and what the parser is to do with the data.

The grammar file includes rules describing the input structure, code to be invoked when these rules are recognized, and a subroutine to do the basic input.

The **yacc** command uses the information in the grammar file to generate a parser that controls the input process. This parser calls an input subroutine (the lexical analyzer) to pick up the basic items (called

tokens) from the input stream. A token is a symbol or name that tells the parser which pattern is being sent to it by the input subroutine. A nonterminal symbol is the structure that the parser recognizes. The parser organizes these tokens according to the structure rules in the grammar file. The structure rules are called *grammar rules*. When the parser recognizes one of these rules, it executes the user code supplied for that rule. The user code is called an *action*. Actions return values and use the values returned by other actions.

Use the C programming language to write the action code and other subroutines. The **yacc** command uses many of the C language syntax conventions for the grammar file.

main and yyerror subroutines

You must provide the **main** and **yyerror** subroutines for the parser. To ease the initial effort of using the **yacc** command, the **yacc** library contains simple versions of the **main** and **yyerror** subroutines. Include these subroutines by using the *-ly* argument to the **ld** command (or to the **cc** command). The source code for the **main** library program is as follows:

```
#include <locale.h>
main()
{
    setlocale(LC_ALL, "");
    yyparse();
}
```

The source code for the **yyerror** library program is as follows:

```
#include <stdio.h>
yyerror(s)
    char *s;
{
    fprintf( stderr, "%s\n" ,s);
}
```

The argument to the **yyerror** subroutine is a string containing an error message, usually the string syntax **error**.

Because these programs are limited, provide more function in these subroutines. For example, keep track of the input line number and print it along with the message when a syntax error is detected. You may also want to use the value in the external integer variable *yychar*. This variable contains the look-ahead token number at the time the error was detected.

yylex Subroutine

The input subroutine that you supply to the grammar file must be able to do the following:

- Read the input stream.
- Recognize basic patterns in the input stream.
- Pass the patterns to the parser, along with tokens that define the pattern to the parser.

For example, the input subroutine separates an input stream into the tokens of **WORD**, **NUMBER**, and **PUNCTUATION**, and it receives the following input:

```
I have 9 turkeys.
```

The program could choose to pass the following strings and tokens to the parser:

String	Token
I	WORD
have	WORD
9	NUMBER
turkeys	WORD
.	PUNCTUATION

The parser must contain definitions for the tokens passed to it by the input subroutine. Using the **-d** option for the **yacc** command, it generates a list of tokens in a file called **y.tab.h**. This list is a set of **#define** statements that allow the lexical analyzer (**yylex**) to use the same tokens as the parser.

Note: To avoid conflict with the parser, do not use subroutine names that begin with the letters **yy**.

You can use the **lex** command to generate the input subroutine, or you can write the routine in the C language.

Using the yacc grammar file

A **yacc** grammar file consists of the following sections:

- Declarations
- Rules
- Programs

Two adjacent **%%** (percent sign, percent sign) separate each section of the grammar file. To make the file easier to read, put the **%%** on a line by themselves. A complete grammar file looks like the following:

```
declarations
%%
rules
%%
programs
```

The declarations section may be empty. If you omit the programs section, omit the second set of **%%**. Therefore, the smallest **yacc** grammar file is as follows:

```
%%
rules
```

The **yacc** command ignores blanks, tabs, and new line characters in the grammar file. Therefore, use these characters to make the grammar file easier to read. Do not, however, use blanks, tabs or new line characters in names or reserved symbols.

Using comments

To explain what the program is doing, put comments in the grammar file. You can put comments anywhere in the grammar file that you can put a name. However, to make the file easier to read, put the comments on lines by themselves at the beginning of functional blocks of rules. A comment in a **yacc** grammar file looks the same as a comment in a C language program. The comment is enclosed between **/*** (backslash, asterisk) and ***/** (asterisk, backslash). For example:

```
/* This is a comment on a line by itself. */
```

Using literal strings

A literal string is one or more characters enclosed in **' '** (single quotes). As in the C language, the **** (backslash) is an escape character within literals, and all the C language escape codes are recognized. Thus, the **yacc** command accepts the symbols in the following table:

Symbol	Definition
'\a'	Alert
'\b'	Backspace
'\f'	Form-feed
'\n'	New line
'\r'	Return
'\t'	Tab
'\v'	Vertical tab

declaring a name in the declarations section, you can use that name only as a nonterminal symbol. The **#include** statements are identical to C language syntax and perform the same function.

The **yacc** program has a set of keywords that define processing conditions for the generated parser. Each of the keywords begin with a % (percent sign), which is followed by a token or nonterminal name. These keywords are as follows:

Keyword	Description
%left	Identifies tokens that are left-associative with other tokens.
%nonassoc	Identifies tokens that are not associative with other tokens.
%right	Identifies tokens that are right-associative with other tokens.
%start	Identifies a nonterminal name for the start symbol.
%token	Identifies the token names that the yacc command accepts. Declares all token names in the declarations section.
%type	Identifies the type of nonterminals. Type-checking is performed when this construct is present.
%union	Identifies the yacc value stack as the union of the various type of values desired. By default, the values returned are integers. The effect of this construct is to provide the declaration of YYSTYPE directly from the input.

```
%{  
Code  
%}
```

Copies the specified *Code* into the code file. This construct can be used to add C language declarations and definitions to the declarations section.

Note: The %{ (percent sign, left bracket) and %} (percent sign, right bracket) symbols must appear on lines by themselves.

The **%token**, **%left**, **%right**, and **%nonassoc** keywords optionally support the name of a C union member (as defined by **%union**) called a *<Tag>* (literal angle brackets surrounding a union member name). The **%type** keyword requires a *<Tag>*. The use of *<Tag>* specifies that the tokens named on the line are to be of the same C type as the union member referenced by *<Tag>*. For example, the following declaration declares the *Name* parameter to be a token:

```
%token [<Tag>] Name [Number] [Name [Number]]...
```

If *<Tag>* is present, the C type for all tokens on this line are declared to be of the type referenced by *<Tag>*. If a positive integer, *Number*, follows the *Name* parameter, that value is assigned to the token.

All of the tokens on the same line have the same precedence level and associativity. The lines appear in the file in order of increasing precedence or binding strength. For example, the following describes the precedence and associativity of the four arithmetic operators:

```
%left '+' '-'  
%left '*' '/'
```

The + (plus sign) and - (minus sign) are left associative and have lower precedence than * (asterisk) and / (slash), which are also left associative.

Defining global variables

To define variables to be used by some or all actions, as well as by the lexical analyzer, enclose the declarations for those variables between %{ (percent sign, left bracket) and %} (percent sign, right bracket) symbols. Declarations enclosed in these symbols are called *global variables*. For example, to

make the **var** variable available to all parts of the complete program, use the following entry in the declarations section of the grammar file:

```
%{  
int var = 0;  
%}
```

Start conditions

The parser recognizes a special symbol called the *start* symbol. The start symbol is the name of the rule in the rules section of the grammar file that describes the most general structure of the language to be parsed. Because it is the most general structure, the parser starts in its top-down analysis of the input stream at this point. Declare the start symbol in the declarations section using the **%start** keyword. If you do not declare the name of the start symbol, the parser uses the name of the first grammar rule in the grammar file.

For example, when parsing a C language function, the most general structure for the parser to recognize is as follows:

```
main()  
{  
    code_segment  
}
```

The start symbol points to the rule that describes this structure. All remaining rules in the file describe ways to identify lower-level structures within the function.

Token numbers

Token numbers are nonnegative integers that represent the names of tokens. If the lexical analyzer passes the token number to the parser, instead of the actual token name, both programs must agree on the numbers assigned to the tokens.

You can assign numbers to the tokens used in the **yacc** grammar file. If you do not assign numbers to the tokens, the **yacc** grammar file assigns numbers using the following rules:

- A literal character is the numerical value of the character in the ASCII character set.
- Other names are assigned token numbers starting at 257.

Note: Do not assign a token number of 0. This number is assigned to the endmarker token. You cannot redefine it.

To assign a number to a token (including literals) in the declarations section of the grammar file, put a positive integer (not 0) immediately following the token name in the **%token** line. This integer is the token number of the name or literal. Each token number must be unique. All lexical analyzers used with the **yacc** command must return a 0 or a negative value for a token when they reach the end of their input.

yacc rules

The rules section of the grammar file contains one or more grammar rules. Each rule describes a structure and gives it a name.

A grammar rule has the following form:

```
A : BODY;
```

where A is a nonterminal name, and BODY is a sequence of 0 or more names, literals, and semantic actions that can optionally be followed by precedence rules. Only the names and literals are required to form the grammar. Semantic actions and precedence rules are optional. The colon and the semicolon are required **yacc** punctuation.

Semantic actions allow you to associate actions to be performed each time that a rule is recognized in the input process. An action can be an arbitrary C statement, and as such, perform input or output, call subprograms, or alter external variables. Actions can also refer to the actions of the parser; for example, shift and reduce.

Precedence rules are defined by the **%prec** keyword and change the precedence level associated with a particular grammar rule. The reserved symbol **%prec** can appear immediately after the body of the grammar rule and can be followed by a token name or a literal. The construct causes the precedence of the grammar rule to become that of the token name or literal.

Repeating nonterminal names

If several grammar rules have the same nonterminal name, use the | (pipe symbol) to avoid rewriting the left side. In addition, use the ; (semicolon) only at the end of all rules joined by pipe symbols. For example, the following grammar rules:

```
A : B C D ;
A : E F ;
A : G ;
```

can be given to the **yacc** command by using the pipe symbol as follows:

```
A : B C D
  | E F
  | G
  ;
```

Using recursion in a grammar file

Recursion is the process of using a function to define itself. In language definitions, these rules normally take the following form:

```
rule : EndCase
     | rule EndCase
```

Therefore, the simplest case of the `rule` is the *EndCase*, but `rule` can also consist of more than one occurrence of *EndCase*. The entry in the second line that uses `rule` in the definition of `rule` is the recursion. The parser cycles through the input until the stream is reduced to the final *EndCase*.

When using recursion in a rule, always put the call to the name of the rule as the leftmost entry in the rule (as it is in the preceding example). If the call to the name of the rule occurs later in the line, such as in the following example, the parser may run out of internal stack space and stop.

```
rule : EndCase
     | EndCase rule
```

The following example defines the `line` rule as one or more combinations of a string followed by a newline character (`\n`):

```
lines : line
      | lines line
      ;
line  : string '\n'
      ;
```

Empty string

To indicate a nonterminal symbol that matches the empty string, use a ; (semicolon) by itself in the body of the rule. To define a symbol empty that matches the empty string, use a rule similar to the following rule:

```
empty : ;
      | x;
```

OR

```
empty :
      | x
      ;
```

End-of-input marker

When the lexical analyzer reaches the end of the input stream, it sends an end-of-input marker to the parser. This marker is a special token called *endmarker*, which has a token value of 0. When the parser receives an end-of-input marker, it checks to see that it has assigned all input to defined grammar rules and that the processed input forms a complete unit (as defined in the **yacc** grammar file). If the input is a complete unit, the parser stops. If the input is not a complete unit, the parser signals an error and stops.

The lexical analyzer must send the end-of-input marker at the appropriate time, such as the end of a file, or the end of a record.

yacc actions

With each grammar rule, you can specify actions to be performed each time the parser recognizes the rule in the input stream. An action is a C language statement that does input and output, calls subprograms, and alters external vectors and variables.

Actions return values and obtain the values returned by previous actions. The lexical analyzer can also return values for tokens.

Specify an action in the grammar file with one or more statements enclosed in `{ }` (braces). The following examples are grammar rules with actions:

```
A : '('B')'
  {
    hello(1, "abc" );
  }
```

AND

```
XXX : YYY ZZZ
     {
       printf("a message\n");
       flag = 25;
     }
```

Passing Values between Actions

To get values generated by other actions, an action can use the **yacc** parameter keywords that begin with a dollar sign (`$1`, `$2`, ...). These keywords refer to the values returned by the components of the right side of a rule, reading from left to right. For example, if the rule is:

```
A : B C D ;
```

then `$1` has the value returned by the rule that recognized B, `$2` has the value returned by the rule that recognized C, and `$3` the value returned by the rule that recognized D.

To return a value, the action sets the pseudo-variable `$$` to some value. For example, the following action returns a value of 1:

```
{ $$ = 1; }
```

By default, the value of a rule is the value of the first element in it (`$1`). Therefore, you do not need to provide an action for rules that have the following form:

```
A : B ;
```

The following additional **yacc** parameter keywords beginning with a `$` (dollar sign) allow for type-checking:

- `$<Tag>$`
- `$<Tag>Number`

`$<Tag>Number` imposes on the reference the type of the union member referenced by `<Tag>`. This adds `.tag` to the reference so that the union member identified by `Tag` is accessed. This construct is equivalent to specifying `$$. Tag` or `$1 . Tag`. You can use this construct when you use actions in the middle of rules where the return type cannot be specified through a `%type` declaration. If a `%type` has

been declared for a nonterminal name, do not use the `<Tag>` construct; the union reference will be done automatically.

Putting actions in the middle of rules

To get control of the parsing process before a rule is completed, write an action in the middle of a rule. If this rule returns a value through the `$` keywords, actions that follow this rule can use that value. This rule can also use values returned by actions that precede it. Therefore, the following rule sets `x` to 1 and `y` to the value returned by `C`. The value of rule `A` is the value returned by `B`, following the default rule.

```
A : B
    {
        $$ =1;
    }
    C
    {
        x = $2;
        y = $3;
    }
    ;
```

Internally, the **yacc** command creates a new nonterminal symbol name for the action that occurs in the middle. It also creates a new rule matching this name to the empty string. Therefore, the **yacc** command treats the preceding program as if it were written in the following form:

```
$ACT : /* empty */
    {
        $$ = 1;
    }
    ;
A : B $ACT C
    {
        x = $2;
        y = $3;
    }
    ;
```

where `$ACT` is an empty action.

yacc program error handling

When the parser reads an input stream, that input stream might not match the rules in the grammar file.

The parser detects the problem as early as possible. If there is an error-handling subroutine in the grammar file, the parser can allow for entering the data again, ignoring the bad data, or initiating a cleanup and recovery action. When the parser finds an error, for example, it may need to reclaim parse tree storage, delete or alter symbol table entries, and set switches to avoid generating further output.

When an error occurs, the parser stops unless you provide error-handling subroutines. To continue processing the input to find more errors, restart the parser at a point in the input stream where the parser can try to recognize more input. One way to restart the parser when an error occurs is to discard some of the tokens following the error. Then try to restart the parser at that point in the input stream.

The **yacc** command uses a special token name, **error**, for error handling. Put this token in the rules file at places that an input error might occur so that you can provide a recovery subroutine. If an input error occurs in this position, the parser executes the action for the **error** token, rather than the normal action.

The following macros can be placed in **yacc** actions to assist in error handling:

Macros	Description
YYERROR	Causes the parser to initiate error handling
YYABORT	Causes the parser to return with a value of 1
YYACCEPT	Causes the parser to return with a value of 0
YYRECOVERING()	Returns a value of 1 if a syntax error has been detected and the parser has not yet fully recovered

To prevent a single error from producing many error messages, the parser remains in error state until it processes three tokens following an error. If another error occurs while the parser is in the error state, the parser discards the input token and does not produce a message.

For example, a rule of the following form:

```
stat : error ';' ;
```

tells the parser that when there is an error, it should ignore the token and all following tokens until it finds the next semicolon. All tokens after the error and before the next semicolon are discarded. After finding the semicolon, the parser reduces this rule and performs any cleanup action associated with it.

Providing for error correction

You can also allow the person entering the input stream in an interactive environment to correct any input errors by entering a line in the data stream again. The following example shows one way to do this.

```
input : error '\n'
      {
        printf(" Reenter last line: " );
      }
      input
      {
        $$ = $4;
      }
      ;
```

However, in this example, the parser stays in the error state for three input tokens following the error. If the corrected line contains an error in the first three tokens, the parser deletes the tokens and does not produce a message. To allow for this condition, use the following **yacc** statement:

```
yyerror;
```

When the parser finds this statement, it leaves the error state and begins processing normally. The error-recovery example then becomes:

```
input : error '\n'
      {
        yyerror;
        printf(" Reenter last line: " );
      }
      input
      {
        $$ = $4;
      }
      ;
```

Clearing the look-ahead token

The *look-ahead token* is the next token that the parser examines. When an error occurs, the look-ahead token becomes the token at which the error was detected. However, if the error recovery action includes code to find the correct place to start processing again, that code must also change the look-ahead token. To clear the look-ahead token, include the following statement in the error-recovery action:

```
yyclearin ;
```

Parser operation generated by the yacc command

The **yacc** command converts the grammar file to a C language program.

That program, when compiled and executed, parses the input according to the grammar specification provided.

The parser is a finite state machine with a stack. The parser can read and remember the look-ahead token. The current state is always the state at the top of the stack. The states of the finite state machine are represented by small integers. Initially, the machine is in state 0, the stack contains only 0, and no look-ahead token has been read.

The machine can perform one of the following actions:

Action	Description
shift <i>State</i>	The parser pushes the current state onto the stack, makes <i>State</i> the current state, and clears the look-ahead token.
reduce <i>Rule</i>	When the parser finds a string defined by <i>Rule</i> (a rule number) in the input stream, the parser replaces that string with <i>Rule</i> in the output stream.
accept	The parser looks at all input, matches it to the grammar specification, and recognizes the input as satisfying the highest-level structure (defined by the start symbol). This action appears only when the look-ahead token is the endmarker and indicates that the parser has successfully done its job.
error	The parser cannot continue processing the input stream and still successfully match it with any rule defined in the grammar specification. The input tokens that the parser looked at, together with the look-ahead token, cannot be followed by anything that would result in valid input. The parser reports an error and attempts to recover the situation and resume parsing.

The parser performs the following actions during one process step:

1. Based on its current state, the parser decides whether it needs a look-ahead token to determine the action to be taken. If the parser needs a look-ahead token and does not have one, it calls the **yylex** subroutine to obtain the next token.
2. Using the current state, and the look-ahead token if needed, the parser decides on its next action and carries it out. As a result, states may be pushed onto or popped off the stack, and the look-ahead token may be processed or left alone.

Shift action

The **shift** action is the most common action that the parser takes. Whenever the parser does a shift, a look-ahead token always exists. For example, consider the following grammar specification rule:

```
IF shift 34
```

If the parser is in the state that contains this rule and the look-ahead token is IF, the parser:

1. Pushes the current state down on the stack.
2. Makes state 34 the current state (puts it at the top of the stack).
3. Clears the look-ahead token.

Reduce action

The **reduce** action keeps the stack from growing too large. The parser uses reduce actions after matching the right side of a rule with the input stream. The parser is then ready to replace the characters in the input stream with the left side of the rule. The parser may have to use the look-ahead token to decide if the pattern is a complete match.

Reduce actions are associated with individual grammar rules. Because grammar rules also have small integer numbers, it is easy to confuse the meanings of the numbers in the two actions, **shift** and **reduce**. For example, the following action refers to grammar rule 18:

```
. reduce 18
```

The following action refers to state 34:

```
IF shift 34
```

For example, to reduce the following rule, the parser pops off the top three states from the stack:

```
A : x y z ;
```

The number of states popped equals the number of symbols on the right side of the rule. These states are the ones put on the stack while recognizing *x*, *y*, and *z*. After popping these states, a state is uncovered, which is the state that the parser was in before beginning to process the rule; that is, the state that needed to recognize rule *A* to satisfy its rule. Using this uncovered state and the symbol on the left side of the rule, the parser performs an action called **goto**, which is similar to a shift of *A*. A new state is obtained, pushed onto the stack, and parsing continues.

The **goto** action is different from an ordinary shift of a token. The look-ahead token is cleared by a shift but is not affected by a **goto** action. When the three states are popped, the uncovered state contains an entry such as the following:

```
A goto 20
```

This entry causes state 20 to be pushed onto the stack and become the current state.

The reduce action is also important in the treatment of user-supplied actions and values. When a rule is reduced, the parser executes the code that you included in the rule before adjusting the stack. Another stack running in parallel with the stack holding the states holds the values returned from the lexical analyzer and the actions. When a shift takes place, the *yyval* external variable is copied onto the stack holding the values. After executing the code that you provide, the parser performs the reduction. When the parser performs the **goto** action, it copies the *yyval* external variable onto the value stack. The **yacc** keywords that begin with \$ refer to the value stack.

Using ambiguous rules in the yacc program

A set of grammar rules is *ambiguous* if any possible input string can be structured in two or more different ways.

For example, the following grammar rule states a rule that forms an arithmetic expression by putting two other expressions together with a minus sign between them.

```
expr : expr '-' expr
```

Unfortunately, this grammar rule does not specify how to structure all complex inputs. For example, if the input is:

```
expr - expr - expr
```

a program could structure this input as either left associative:

```
( expr - expr ) - expr
```

or as right associative:

```
expr - ( expr - expr )
```

and produce different results.

Parser conflicts

When the parser tries to handle an ambiguous rule, confusion can occur over which of its four actions to perform when processing the input. The following major types of conflict develop:

Conflict

shift/reduce conflict

Description

A rule can be evaluated correctly using either a **shift** action or a **reduce** action, but the result is different.

reduce/reduce conflict

A rule can be evaluated correctly using one of two different reduce actions, producing two different actions.

A **shift/shift** conflict is not possible. The **shift/reduce** and **reduce/reduce** conflicts result from a rule that is not completely stated. For example, using the ambiguous rule stated in the previous section, if the parser receives the input:

```
expr - expr - expr
```

after reading the first three parts, the parser has:

```
expr - expr
```

which matches the right side of the preceding grammar rule. The parser can reduce the input by applying this rule. After applying the rule, the input becomes:

```
expr
```

which is the left side of the rule. The parser then reads the final part of the input:

```
- expr
```

and reduces it. This produces a left-associative interpretation.

However, the parser can also look ahead in the input stream. If, when the parser receives the first three parts:

```
expr - expr
```

it reads the input stream until it has the next two parts, it then has the following input:

```
expr - expr - expr
```

Applying the rule to the rightmost three parts reduces them to `expr`. The parser then has the expression:

```
expr - expr
```

Reducing the expression once more produces a right-associative interpretation.

Therefore, at the point where the parser has read only the first three parts, it can take either of two valid actions: a **shift** or a **reduce**. If the parser has no rule to decide between the two actions, a **shift/reduce** conflict results.

A similar situation occurs if the parser can choose between two valid reduce actions, which is called a *reduce/reduce conflict*.

How the parser responds to conflicts

When **shift/reduce** or **reduce/reduce** conflicts occur, the **yacc** command produces a parser by selecting one of the valid steps wherever it has a choice. If you do not provide a rule that makes the choice, the **yacc** program uses the following rules:

- In a **shift/reduce** conflict, choose the shift.
- In a **reduce/reduce** conflict, reduce by the grammar rule that can be applied at the earliest point in the input stream.

Using actions within rules can cause conflicts if the action must be performed before the parser is sure which rule is being recognized. In such cases, the preceding rules result in an incorrect parser. For this reason, the **yacc** program reports the number of **shift/reduce** and **reduce/reduce** conflicts resolved by using the preceding rules.

Turning on debug mode for a parser generated by the yacc command

You can access the debugging code either by invoking the **yacc** command with the **-t** option or compiling the **y.tab.c** file with **-DYYDEBUG**.

For normal operation, the *yydebug* external integer variable is set to 0. However, if you set it to a nonzero value, the parser generates a description of the input tokens it receives and actions it takes for each token while parsing an input stream.

Set this variable in one of the following ways:

- Put the following C language statement in the declarations section of the **yacc** grammar file:

```
int yydebug = 1;
```

- Use the **dbx** program to execute the final parser, and set the variable ON or OFF using **dbx** commands.

Example program for the lex and yacc programs

This section contains example programs for the **lex** and **yacc** commands.

Together, these example programs create a simple, desk-calculator program that performs addition, subtraction, multiplication, and division operations. This calculator program also allows you to assign values to variables (each designated by a single, lowercase letter) and then use the variables in calculations. The files that contain the example **lex** and **yacc** programs are as follows:

File	Content
calc.lex	Specifies the lex command specification file that defines the lexical analysis rules.
calc.yacc	Specifies the yacc command grammar file that defines the parsing rules, and calls the yylex subroutine created by the lex command to provide input.

The following descriptions assume that the **calc.lex** and **calc.yacc** example programs are located in your current directory.

Compiling the example program

To create the desk calculator example program, do the following:

1. Process the **yacc** grammar file using the **-d** optional flag (which informs the **yacc** command to create a file that defines the tokens used in addition to the C language source code):

```
yacc -d calc.yacc
```

2. Use the **ls** command to verify that the following files were created:

y.tab.c

The C language source file that the **yacc** command created for the parser

y.tab.h

A header file containing define statements for the tokens used by the parser

3. Process the **lex** specification file:

```
lex calc.lex
```

4. Use the **ls** command to verify that the following file was created:

lex.yy.c

The C language source file that the **lex** command created for the lexical analyzer

5. Compile and link the two C language source files:

```
cc y.tab.c lex.yy.c
```

6. Use the **ls** command to verify that the following files were created:

y.tab.o

The object file for the **y.tab.c** source file

lex.yy.o

The object file for the **lex.yy.c** source file

a.out

The executable program file

To run the program directly from the **a.out** file, type:

```
$ a.out
```

OR

To move the program to a file with a more descriptive name, as in the following example, and run it, type:

```
$ mv a.out calculate
$ calculate
```

In either case, after you start the program, the cursor moves to the line below the \$ (command prompt). Then, enter numbers and operators as you would on a calculator. When you press the Enter key, the program displays the result of the operation. After you assign a value to a variable, as follows, the cursor moves to the next line.

```
m=4 <enter>
-
```

When you use the variable in subsequent calculations, it will have the assigned value:

```
m+5 <enter>
9
-
```

Parser source code

The following example shows the contents of the **calc.yacc** file. This file has entries in all three sections of a **yacc** grammar file: declarations, rules, and programs.

```
%{
#include<stdio.h>

int regs[26];
int base;

%}

%start list

%union { int a; }

%token DIGIT LETTER

%left '|'
%left '&'
%left '+' '-'
%left '*' '/' '%'
%left UMINUS /*supplies precedence for unary minus */

%%
/* beginning of rules section */

list:
    /*empty */
    |
    list stat '\n'
    |
    list error '\n'
    {
        yyerrok;
    }
    ;
stat:
    expr
    {

```

```

        printf("%d\n", $1);
    }
    |
    LETTER '=' expr
    {
        regs[$1.a] = $3.a;
    }
    ;

expr:  '(' expr ')'
    {
        $$ = $2;
    }
    |
    expr '*' expr
    {
        $$ .a = $1.a * $3.a;
    }
    |
    expr '/' expr
    {
        $$ .a = $1.a / $3.a;
    }
    |
    expr '%' expr
    {
        $$ .a = $1.a % $3.a;
    }
    |
    expr '+' expr
    {
        $$ .a = $1.a + $3.a;
    }
    |
    expr '-' expr
    {
        $$ .a = $1.a - $3.a;
    }
    |
    expr '&' expr
    {
        $$ .a = $1.a & $3.a;
    }
    |
    expr '|' expr
    {
        $$ .a = $1.a | $3.a;
    }
    |
    '-' expr %prec UMINUS
    {
        $$ .a = -$2.a;
    }
    |
    LETTER
    {
        $$ .a = regs[$1.a];
    }
    |
    number
    ;

number: DIGIT
    {
        $$ = $1;
        base = ($1.a==0) ? 8 : 10;
    }
    |
    number DIGIT
    {
        $$ .a = base * $1.a + $2.a;
    }
    ;

%%
main()
{
    return(yyparse());
}

```

```

yyerror(s)
char *s;
{
    fprintf(stderr, "%s\n",s);
}

yywrap()
{
    return(1);
}

```

The file contains the following sections:

- **Declarations section.** This section contains entries that:
 - Include standard I/O header file
 - Define global variables
 - Define the `list` rule as the place to start processing
 - Define the tokens used by the parser
 - Define the operators and their precedence
- **Rules section.** The rules section defines the rules that parse the input stream.
 - **%start** - Specifies that the whole input should match **stat**.
 - **%union** - By default, the values returned by actions and the lexical analyzer are integers. **yacc** can also support values of other types, including structures. In addition, **yacc** keeps track of the types, and inserts appropriate union member names so that the resulting parser will be strictly type checked. The **yacc** value stack is declared to be a union of the various types of values desired. The user declares the union, and associates union member names to each token and nonterminal symbol having a value. When the value is referenced through a `$$` or `$n` construction, **yacc** will automatically insert the appropriate union name, so that no unwanted conversions will take place.
 - **%type** - Makes use of the members of the **%union** declaration and gives an individual type for the values associated with each part of the grammar.
 - **%toksn** - Lists the tokens which come from lex tool with their type.
- **Programs section.** The programs section contains the following subroutines. Because these subroutines are included in this file, you do not need to use the **yacc** library when processing this file.

Subroutine	Description
main	The required main program that calls the yyparse subroutine to start the program.
yyerror(s)	This error-handling subroutine only prints a syntax error message.
yywrap	The wrap-up subroutine that returns a value of 1 when the end of input occurs.

Lexical analyzer source code

This file contains include statements for standard input and output, as well as for the **y.tab.h** file. If you use the **-d** flag with the **yacc** command, the **yacc** program generates that file from the **yacc** grammar file information. The **y.tab.h** file contains definitions for the tokens that the parser program uses. In addition, the **calc.lex** file contains the rules to generate these tokens from the input stream. The following are the contents of the **calc.lex** file.

```

%{
#include <stdio.h>
#include "y.tab.h"
int c;
%}
%%
" "      ;
[a-z]    {
          c = yytext[0];
          yylval.a = c - 'a';

```



```

        return(LETTER);
    }
[0-9]   {
        c = yytext[0];
        yy1val.a = c - '0';
        return(DIGIT);
    }
[^a-z0-9\b] {
        c = yytext[0];
        return(c);
    }
%%

```

make command

This topic provides information about simplifying the recompiling and relinking processes using the **make** command.

It allows you to record, once only, specific relationships among files. You can then use the **make** command to automatically perform all updating tasks.

In any project, you normally link programs from object files and libraries. Then, after modifying a source file, you recompile some of the sources and relink the program as often as required. The **make** command assists in maintaining a set of programs, usually pertaining to a particular software project, by building up-to-date versions of programs. The **make** command is most useful for medium-sized programming projects. It does not solve the problems of maintaining more than one source version and of describing large programs (see **sccs** command).

Using the **make** command to maintain programs, you can:

- Combine instructions for creating a large program in a single file.
- Define macros to use within the **make** command description file.
- Use shell commands to define the method of file creation, or use the **make** command to create many of the basic types of files.
- Create libraries.

The **make** command requires a description file, file names, specified rules to inform the **make** command how to build many standard types of files, and time stamps of all system files.

Related concepts

Tools and utilities

This section provides an overview of the tools and utilities that you can use to develop C compiled language programs.

Creating a description file

The **make** command uses information from a description file, which you create, to build a file containing the completed program, which is then called a *target* file.

The description file informs the **make** command how to build the target file, which files are involved, and what their relationships are to the other files in the procedure. The description file contains the following information:

- Target file name
- Parent file names that make up the target file
- Commands that create the target file from the parent files
- Definitions of macros in the description file
- User-specified rules for building target files

By checking the dates of the parent files, the **make** command determines which files to create to get an up-to-date copy of the target file. If any parent file was changed more recently than the target file, the **make** command creates the files affected by the change, including the target file.

If you name the description file either **makefile** or **Makefile** and are working in the directory containing that description file, type the following to update the first target file and its parent files:

```
make
```

Updating occurs regardless of the number of files changed since the last time the **make** command created the target file. In most cases, the description file is easy to write and does not change often.

To keep many different description files in the same directory, name them differently. Then, enter:

```
make -f Desc-File
```

where *Desc-File* is the name of the description file.

Format of an entry in the make description file

The general form of an entry is:

```
target1 [target2..]:[:] [parent1..][; command]...  
[(tab) commands]
```

Items inside brackets are optional. Targets and parents are file names (strings of letters, numbers, periods, and slashes). The **make** command recognizes wildcard characters such as * (asterisk) and ? (question mark). Each line in the description file that contains a target file name is called a *dependency line*. Lines that contain commands must begin with a tab character.

Note: The **make** command uses the \$ (dollar sign) to designate a macro. Do not use that character in file names of target or parent files, or in commands in the description file unless you are using a predefined **make** command macro.

Begin comments in the description file with a # (pound sign). The **make** command ignores the # and all characters that follow it. The **make** command also ignores blank lines.

Except for comment lines, you can enter lines longer than the line width of the input device. To continue a line on the next line, put a \ (backslash) at the end of the line to be continued.

Using commands in a make description file

A command is any string of characters except a # (pound sign) or a new-line character. A command can use a # only if it is in quotation marks.

Commands can appear either after a semicolon on a dependency line or on lines beginning with a tab that immediately follows a dependency line.

When defining the command sequence for a particular target, specify one command sequence for each target in the description file, or else separate the command sequences for special sets of dependencies. Do not do both.

To use one command sequence for every use of the target file, use a single : (colon) following the target name on the dependency line. For example:

```
test:      dependency list1...  
          command list...  
          .  
          .  
test:      dependency list2...
```

defines a target name, *test*, with a set of parent files and a set of commands to create the file. The target name, *test*, can appear in other places in the description file with another dependency list. However, that name cannot have another command list in the description file. When one of the files that *test* depends on changes, the **make** command runs the commands in that one command list to create the target file named *test*.

To specify more than one set of commands to create a particular target file, enter more than one dependency definition. Each dependency line must have the target name, followed by :: (two colons), a dependency list, and a command list that the **make** command uses if any of the files in the dependency list changes. For example:

```
test::    dependency list1...
          command list1...
test::    dependency list2...
          command list2...
```

defines two separate processes to create the target file, `test`. If any of the files in dependency `list1` changes, the **make** command runs `command list1`. If any of the files in dependency `list2` changes, the **make** command runs `command list2`. To avoid conflicts, a parent file cannot appear in both dependency `list1` and dependency `list2`.

Note: The **make** command passes the commands from each command line to a new shell. Be careful when using commands that have meaning only within a single shell process; for example, `cd` and shell commands. The **make** command discards these results before running the commands on the next line.

To group commands together, use the `\` (backslash) at the end of a command line. The **make** command then continues that command line into the next line in the description file. The shell sends both of these lines to a single new shell.

Calling the make command from a description file

To nest calls to the **make** command within a **make** command description file, include the **\$(MAKE)** macro in one of the command lines in the file.

If the `-n` flag is set when the **\$(MAKE)** macro is found, the new copy of the **make** command does not execute any of its commands, except another **\$(MAKE)** macro. To use this characteristic to test a set of description files that describe a program, enter:

```
make -n
```

The **make** command does not perform any of the program operations. However, it does write all of the steps needed to build the program, including output from lower-level calls to the **make** command.

Preventing the make command from stopping on errors

The **make** command normally stops if any program returns a nonzero error code. Some programs return status that has no meaning.

To prevent the **make** command from stopping on errors, do any of the following:

- Use the `-i` flag with the **make** command on the command line.
- Put the fake target name **.IGNORE** on a dependency line by itself in the description file. Because **.IGNORE** is not a real target file, it is called a fake target. If **.IGNORE** has prerequisites, the **make** command ignores errors associated with them.
- Put a `-` (minus sign) in the first character position of each line in the description file where the **make** command should not stop on errors.

Example of a description file

For example, a program named **prog** is made by compiling and linking three C language files `x.c`, `y.c`, and `z.c`. The `x.c` and `y.c` files share some declarations in a file named `defs`. The `z.c` file does not share those declarations. The following is an example of a description file, which creates the **prog** program:

```
# Make prog from 3 object files
prog: x.o y.o z.o
# Use the cc program to make prog
cc x.o y.o z.o -o prog
# Make x.o from 2 other files
x.o: x.c defs
```

```

# Use the cc program to make x.o
cc -c x.c
# Make y.o from 2 other files
y.o: y.c defs
# Use the cc program to make y.o
cc -c y.c
# Make z.o from z.c
z.o: z.c
# Use the cc program to make z.o
cc -c z.c

```

If this file is called `makefile`, enter the `make` command to update the **prog** program after making changes to any of the source files: `x.c`, `y.c`, `z.c`, or `defs`.

Simplifying the description file

To make this file simpler, use the internal rules of the `make` program.

Based on file-system naming conventions, the `make` command recognizes three `.c` files corresponding to the needed `.o` files. This command can also generate an object from a source file, by issuing a `cc -c` command.

Based on these internal rules, the description file becomes:

```

# Make prog from 3 object files
prog: x.o y.o z.o
# Use the cc program to make prog
cc x.o y.o z.o -o prog
# Use the file defs and the .c file
# when making x.o and y.o
x.o y.o: defs

```

Internal rules for the make program

The internal rules for the `make` command are located in a file that looks like a description file.

When the `-r` flag is specified, the `make` command does not use the internal rules file. You must supply the rules to create the files in your description file. The internal-rules file contains a list of file-name suffixes (such as `.o` or `.a`) that the `make` command understands, and rules that inform the `make` command how to create a file with one suffix from a file with another suffix. If you do not change the list, the `make` command understands the following suffixes:

Suffi Description

x

- .a** Archive library
- .C** C++ source file
- .C\~** Source Code Control System (SCCS) file containing C++ source file
- .c** C source file
- .c~** SCCS file containing C source file
- .f** FORTRAN source file
- .f~** SCCS file containing FORTRAN source file
- .h** C language header file.
- .h~** SCCS file containing C language header file
- .l** **lex** source grammar
- .l~** SCCS file containing **lex** source grammar.
- .o** Object file.
- .s** Assembler source file.

Suffix Description

x

- .s~** SCCS file containing assembler source file.
- .sh** Shell-command source file.
- .sh~** SCCS file containing shell-command source file.
- .y** **yacc-c** source grammar.
- .y~** SCCS file containing **yacc-c** source grammar.

The list of suffixes is similar to a dependency list in a description file and follows the fake target name of **.SUFFIXES**. Because the **make** command looks at the suffixes list in left-to-right order, the order of the entries is important.

The **make** command uses the first entry in the list that satisfies the following requirements:

- The entry matches input and output suffix requirements for the current target and dependency files.
- The entry has a rule assigned to it.

The **make** command creates the name of the rule from the two suffixes of the files that the rule defines. For example, the name of the rule to transform a **.c** file to an **.o** file is **.c.o**.

To add more suffixes to the list, add an entry for the fake target name of **.SUFFIXES** in the description file. For a **.SUFFIXES** line without any suffixes following the target name in the description file, the **make** command erases the current list. To change the order of the names in the list, erase the current list and then assign a new set of values to **.SUFFIXES**.

Example of default rules file

The following example shows a portion of the default rules file:

```
# Define suffixes that make knows.
.SUFFIXES: .o .C .C~ .c .c~ .f .f~ .y .y~ .l .l~ .s .s~ .sh .sh~ .h .h~ .a
#Begin macro definitions for
#internal macros
YACC=yacc
YFLAGS=
ASFLAGS=
LEX=lex
LFLAGS=
CC=cc
CCC=x1C
AS=as
CFLAGS=
CCFLAGS=
# End macro definitions for
# internal macros
# Create a .o file from a .c
# file with the cc program.
c.o:
    $(CC) $(CFLAGS) -c $<
```

```
# Create a .o file from
# a .s file with the assembler.
s.o:
    $(AS)$(ASFLAGS) -o $@ $<
```

```
.y.o:
# Use yacc to create an intermediate file
    $(YACC) $(YFLAGS) $<
# Use cc compiler
    $(CC) $(CFLAGS) -c y.tab.c
# Erase the intermediate file
    rm y.tab.c
# Move to target file
    mv y.tab.o $@.
.y.c:
# Use yacc to create an intermediate file
    $(YACC) $(YFLAGS) $<
```

```
# Move to target file
mv y.tab.c $@
```

Single-suffix rules

The **make** command has a set of single-suffix rules to build source files directly into a target file name that does not have a suffix (for example, command files).

The **make** command also has rules to change the following source files with a suffix to object files without a suffix:

Command	Description
.C:	From a C++ language source file
.C\~:	From an SCCS C++ language source file
.c:	From a C language source file
.c~:	From an SCCS C language source file
.sh:	From a shell file
.sh~:	From an SCCS shell file

For example, if all of the needed source files are in the current directory, to maintain the cat program, enter:

```
make cat
```

Using the make command with archive libraries

The **make** command has a set of single-suffix rules to build source files directly into a target file name that does not have a suffix (for example, command files).

The internal rules for changing source files to library files are:

Rule	Explanation
.C.a	C++ source to archive
.C\~.a	SCCS C++ source to archive
.c.a	C source to archive
.c~.a	SCCS C source to archive
.s~.a	SCCS assembler source to archive
.f.a	Fortran source to archive
.f~.a	SCCS Fortran source to archive

Defining default conditions in a description file

When the **make** command creates a target file but cannot find commands in the description file or internal rules to create a file, it looks at the description file for default conditions.

To define the commands that the **make** command performs in this case, use the **.DEFAULT** target name in the description file as follows:

```
.DEFAULT:
    command
    command
    .
    .
    .
```

Because **.DEFAULT** is not a real target file, it is called a *fake target*. Use the **.DEFAULT** fake target name for an error-recovery routine or for a general procedure to create all files in the command that are not defined by an internal rule of the **make** program.

Including other files in a description file

Include files other than the current description file by using the word **include** as the first word on any line in the description file.

Follow the word with a blank or a tab, and then the file name for the **make** command to include in the operation.

Note: Only one file is supported for each **include** statement.

For example:

```
include /home/tom/temp
```

```
include /home/tom/sample
```

directs the **make** command to read the `temp` and `sample` files and the current description file to build the target file.

When you are using include files, do not use more than 16 levels of nesting.

Defining and using macros in a description file

A *macro* is a name (or label) to use in place of several other names. It is a way of writing the longer string of characters by using just one shorter name.

To define a macro:

1. Start a new line with the name of the macro.
2. Follow the name with an = (equal sign).
3. To the right of the = (equal sign), enter the string of characters that the macro name represents.

The macro definition can contain blanks before and after the = (equal sign) without affecting the result. The macro definition cannot contain a : (colon) or a tab before the = (equal sign).

The following are examples of macro definitions:

```
# Macro "-2" has a value of "xyz"
2 = xyz

# Macro "abc" has a value of "-ll -ly"
abc = -ll -ly

# Macro "LIBES" has a null value
LIBES =
```

A macro that is named, but not defined, has the same value as the null string.

Using macros in a description file

After defining a macro in a description file, use the macro in description file commands by putting a \$ (dollar sign) before the name of the macro.

If the macro name is longer than one character, put () (parentheses) or { } (braces) around the macro name. The following are examples of using macros:

```
$(CFLAGS)
$2
$(xy)
$Z
$(Z)
```

The last two examples in the previous list have the same effect.

The following fragment shows how to define and use some macros:

```
# OBJECTS is the 3 files x.o, y.o and
# z.o (previously compiled)
OBJECTS = x.o y.o z.o
# LIBES is the standard library
LIBES = -lc
# prog depends on x.o y.o and z.o
prog: $(OBJECTS)
# Link and load the 3 files with
# the standard library to make prog
cc $(OBJECTS) $(LIBES) -o prog
```

The **make** command that uses this description file links and loads the three object files (x.o, y.o, and z.o) with the **libc.a** library.

A macro definition entered on the command line overrides any duplicate macro definitions in the description file. Therefore, the following command loads the files with the **lex** (-ll) library:

```
make "LIBES= -ll"
```

Note: When macros contain blanks and you enter them, on the command line, put " " (double quotation marks) around the macro. Without the double quotation marks, the shell interprets the blanks as parameter separators and not as part of the macro.

The **make** command handles up to 10 levels of nested macro expansion. Based on the definitions in the following example, the expression $\$(\$(macro2))$ would evaluate to value1:

```
macro1=value1
```

```
macro2=macro1
```

The evaluation of a macro occurs each time the macro is referenced. It is not evaluated when it is defined. If a macro is defined but never used, it will never be evaluated. This is especially important if the macro is assigned values that will be interpreted by the shell, particularly if the value might change. A variable declaration such as:

```
OBJS = 'ls *.o'
```

could change in value if referenced at different times during the process of building or removing object files. It does not hold the value of the **ls** command at the time the **OBJS** macro is defined.

Internal macros

The **make** command has built-in macro definitions for use in the description file.

These macros help specify variables in the description file. The **make** command replaces the macros with one of the following values:

Macro Value

\$@ Name of the current target file

\$\$@ Label name on the dependency line

\$? Names of the files that have changed more recently than the target

\$< Parent file name of the out-of-date file that caused a target file to be created

\$* Name of the current parent file without the suffix

\$% Name of an archive library member

Target file name

If the **\$@** macro is in the command sequence in the description file, the **make** command replaces the symbol with the full name of the current target file before passing the command to the shell to be run. The **make** command replaces the symbol only when it runs commands from the description file to create the target file.

Label name

If the **\$\$@** macro is on the dependency line in a description file, the **make** command replaces this symbol with the label name that is on the left side of the colon in the dependency line. For example, if the following is included on a dependency line:

```
cat:    $$@.c
```

The **make** command translates it to:

```
cat:    cat.c
```

when the **make** command evaluates the expression. Use this macro to build a group of files, each of which has only one source file. For example, to maintain a directory of system commands, use a description file similar to:

```
# Define macro CMDS as a series
# of command names
CMDS = cat dd echo date cc cmp comm ar ld chown
# Each command depends on a .c file
$(CMDS):    $$@.c
# Create the new command set by compiling the out of
# date files ($?) to the target file name ($@)
$(CC) -o $? -o $@
```

The **make** command changes the **\$\$(@F)** macro to the file part of **\$@** when it runs. For example, use this symbol when maintaining the **usr/include** directory while using a description file in another directory. That description file is similar to the following:

```
# Define directory name macro INCDIR
INCDIR = /usr/include
# Define a group of files in the directory
# with the macro name INCLUDES
INCLUDES = \
    $(INCDIR)/stdio.h \
    $(INCDIR)/pwd.h \
    $(INCDIR)/dir.h \
    $(INCDIR)/a.out.h \
# Each file in the list depends on a file
# of the same name in the current directory
$(INCLUDES):    $$(@F)
# Copy the younger files from the current
# directory to /usr/include
cp $? $@
# Set the target files to read only status
chmod 0444 $@
```

The preceding description file creates a file in the **/usr/include** directory when the corresponding file in the current directory has been changed.

Changed Files

If the **\$\$?** macro is in the command sequence in the description file, the **make** command replaces the symbol with a list of parent files that have been changed since the target file was last changed. The **make** command replaces the symbol only when it runs commands from the description file to create the target file.

First out-of-date file

If the **\$\$<** macro is in the command sequence in the description file, the **make** command replaces the symbol with the name of the file that started the file creation. The Out-of-Date file name is the name of the parent file whose timestamp did not match the timestamp of the target file, and therefore caused the **make** command to create the target file again.

In addition, use a letter (**D** or **F**) after the < (less-than sign) to get either the directory name (**D**) or the file name (**F**) of the first out-of-date file. For example, if the first out-of-date file is:

```
/home/linda/sample.c
```

then the **make** command gives the following values:

```
$(<D) = /home/linda
$(<F) = sample.c
$<    = /home/linda/sample.c
```

The **make** command replaces this symbol only when the program runs commands from its internal rules or from the **.DEFAULT** list.

Current file-name prefix

If the **\$*** macro is in the command sequence in the description file, the **make** command replaces the symbol with the file-name part (without the suffix) of the parent file that the **make** command is currently using to generate the target file. For example, if the **make** command is using the file:

```
test.c
```

then the **\$*** macro represents the file name `test`.

In addition, use a letter (**D** or **F**) after the * (asterisk) to get either the directory name (**D**) or the file name (**F**) of the current file.

For example, the **make** command uses many files (specified either in the description file or in the internal rules) to create a target file. Only one of those files (the current file) is used at any moment. If that current file is:

```
/home/tom/sample.c
```

then the **make** command gives the following values for the macros:

```
$(*D) = /home/tom
$(*F) = sample
$*     = /home/tom/sample
```

The **make** command replaces this symbol only when running commands from its internal rules (or from the **.DEFAULT** list), but not when running commands from a description file.

Archive library member

If the **\$%** macro is in a description file, and the target file is an archive library member, the **make** command replaces the macro symbol with the name of the library member. For example, if the target file is:

```
lib(file.o)
```

then the **make** command replaces the **\$%** macro with the member name `file.o`.

Changing macro definitions in a command

When macros in the shell commands are defined in the description file, you can change the values that the **make** command assigns to the macro.

To change the assignment of the macro, put a : (colon) after the macro name, followed by a replacement string. The form is as follows:

```
$(macro:string1=string2)
```

where `string1` is either a suffix or a word to be replaced in the macro definition and `string2` is the replacement suffix or word.

When the **make** command reads the macro and begins to assign the values to the macro based on the macro definition, the command replaces each `string1` in the macro definition with a value of `string2`. For example, if the description file contains the macro definition:

```
FILES=test.o sample.o form.o defs
```

You can replace the `form.o` file with a new file, `input.o`, by using the macro in the description-file commands, as follows:

```
cc -o $(FILES:form.o=input.o)
```

You can all replace all the suffixes `.o` in the macro with `.c` as follows:

```
cc -c $(FILES:.o=.c)
```

The macro values can also be changed using pattern matching replacements of the following form:

```
$(macro: op%os= np%ns)
```

Where `op` is the existing (old) prefix, `os` is the existing (old) suffix, `np` is the new prefix, and `ns` is the new suffix.

`op`, `os`, `np` and `ns` can be a string of zero or more characters. The pattern matched by the percent sign (%) on the left-hand side of the equal-sign, which is a string of zero or more characters, is used with `np` and `ns` to replace the macro value. The percent sign (%) operator can appear any number of times on the right-hand side of the equal sign (=).

For example:

```
F00=abc def  
BAR=$(F00:%=dir1/%.o dir1/%_cltn.o)
```

sets the value of `BAR` to `dir1/abc.o dir1/abc_cltn.o dir1/def.o dir1/def_cltn.o`

Changing the value of a macro in this manner is useful when maintaining archive libraries. For more information, see the **ar** command.

Creating a target file with the make command

To create a file containing the completed program called a *target* file, the **make** program:

1. Finds the name of the target file in the description file or in the **make** command
2. Ensures that the files on which the target file depends exist and are up-to-date
3. Determines if the target file is up-to-date with the files on which it depends.

If the target file or one of the parent files is out-of-date, the **make** command creates the target file using one of the following:

- Commands from the description file
- Internal rules to create the file (if they apply)
- Default rules from the description file

If all files in the procedure are up-to-date when running the **make** program, the **make** command displays a message to indicate that the file is up-to-date, and then stops. If some files have changed, the **make** command builds only those files that are out-of-date. The command does not rebuild files that are already current.

When the **make** command runs commands to create a target file, it replaces macros with their values, writes each command line, and then passes the command to a new copy of the shell.

Using the make command with source code control system files

The source code control system (SCCS) commands and file system is primarily used to control access to a file, track who altered the file, why it was altered, and what was altered.

An SCCS file is any text file controlled with SCCS commands. Using non-SCCS commands to edit SCCS files can damage the SCCS files.

All SCCS files use the prefix **s.** to indicate that the files are SCCS files and not regular text files. The **make** command does not recognize references to prefixes of file names. Therefore, do not refer to SCCS files directly within the **make** command description file. The **make** command uses a different suffix, the ~ (tilde), to represent SCCS files. Therefore, **.c~.o** refers to the rule that transforms an SCCS C language source file into an object file. The internal rule is:

```
.c~.o:
    $(GET) $(GFLAGS) -p $< >$*.c
    $(CC) $(CFLAGS) -c $*.c
    -rm -f $*.c
```

The ~ (tilde) added to any suffix changes the file search into an SCCS file-name search, with the actual suffix named by the . (period) and all characters up to (but not including) the ~ (tilde). The **GFLAGS** macro passes flags to the SCCS to determine which SCCS file version to use.

The **make** command recognizes the following SCCS suffixes:

Suffi	Description
x	
.C\~	C++ source
.c~	c source
.y~	yacc source grammar
.s~	Assembler source
.sh~	Shell
.h~	Header
.f~	FORTRAN
.l~	lex source

The **make** command has internal rules for changing the following SCCS files:

```
.C\~.a:
.C\~.c:
.C\~.o:
.c~:
.c~.a:
.c~.c:
.c~.o:
.f~:
.f~.a:
.f~.o:
.f~.f:
.h~.h:
.l~.o:
.s~.a:
```

```
.sh~:
.s~.o:
.y~.c:
.y~.o:
```

Understanding makefile considerations in Source Code Control Systems (SCCS)

If you specify a description file, or a file named **makefile** or **Makefile** is in the current directory, the **make** command does not look for a description file within SCCS.

If a description file is not in the current directory and you enter the **make** command, the **make** command looks for an SCCS file named either **s.makefile** or **s.Makefile**. If either of these files is present, the **make** command uses a **get** command to direct SCCS to build the description file from that source file. When the SCCS generates the description file, the **make** command uses the file as a normal description file. When the **make** command finishes executing, it removes the created description file from the current directory.

Suppressing Default SCCS Retrieval

You can suppress the default SCCS retrieval. The commands for retrieving source files from SCCS can be specified in the rules for *SCCS_GET* special target in the description file.

This suppresses the default source retrieval from SCCS.

For example:

```
SCCS_GET:
    get -p $< > $*.c
```

Using the make command with non-source code control system (SCCS) files

Start the **make** command from the directory that contains the description file for the file to create.

The variable name *desc-file* represents the name of that description file. Then, enter the command:

```
make -f desc-file
```

If the name of the description file is *makefile* or *Makefile*, you do not have to use the **-f** flag. Enter macro definitions, flags, description file names, and target file names along with the **make** command on the command line as follows:

```
make [flags] [macro definitions] [targets]
```

The **make** command then examines the command-line entries to determine what to do. First, it looks at all macro definitions on the command line (entries that are enclosed in quotes and contain equal signs) and assigns values to them. If the **make** command finds a definition for a macro on the command line different from the definition for that macro in the description file, it uses the command-line definition for the macro.

Next, the **make** command looks at the flags.

The **make** command expects the remaining command-line entries to be the names of target files to be created. Any shell commands enclosed in back quotes that generate target names are performed by the **make** command. Then the **make** command creates the target files in left-to-right order. Without a target file name, the **make** command creates the first target file named in the description file that does not begin with a period. With more than one description file specified, the **make** command searches the first description file for the name of the target file.

Understanding how the make command uses environment variables

Each time the **make** command runs, it reads the current environment variables and adds them to its defined macros.

Using the **MAKEFLAGS** macro or the **MFLAGS** macro, the user can specify flags to be passed to the **make** command. If both are set, the **MAKEFLAGS** macro overrides the **MFLAGS** macro. The flags specified using these variables are passed to the **make** command along with any command-line options. In the case of recursive calls to the **make** command, using the **\$(MAKE)** macro in the description file, the **make** command passes all flags with each invocation.

When the **make** command runs, it assigns macro definitions in the following order:

1. Reads the **MAKEFLAGS** environment variable.
If the **MAKEFLAGS** environment variable is not present or null, the **make** command checks for a non-null value in the **MFLAGS** environment variable. If one of these variables has a value, the **make** command assumes that each letter in the value is an input flag. The **make** command uses these flags (except for the **-f**, **-p**, and **-d** flags, which cannot be set from the **MAKEFLAGS** or **MFLAGS** environment variable) to determine its operating conditions.
2. Reads and sets the input flags from the command line. The command line adds to the previous settings from the **MAKEFLAGS** or **MFLAGS** environment variable.
3. Reads macro definitions from the command line. The **make** command ignores any further assignments to these names.
4. Reads the internal macro definitions.
5. Reads the environment. The **make** command treats the environment variables as macro definitions and passes them to other shell programs.

Using the make command in parallel run mode

Normally, the **make** command runs commands sequentially for only one target at a time, waiting for the command to finish before running the next.

However, the **make** command can also run in a parallel run mode, where it can run many concurrent jobs to build independent targets.

The **-j** flag tells the **make** command simultaneously to run independent targets.

If you follow the **-j** option with an integer, then the integer specifies the maximum number of concurrent jobs that can be executed to build the targets.

If **-j** flag is not followed by an integer, then there is no limit on the number of jobs that will be invoked to build the targets.

If the **make** command encounters an error while building a target in parallel execution mode, and none of the methods used to ignore errors are used, then further execution of commands to build that target will stop, and the **make** command will wait for child jobs already running to exit before completing.

When you are running more than one job in parallel mode, the output that the jobs generate prints to the screen as it is produced. This can lead to the confusion of messages from different jobs on the screen unless you suppress the messages using redirection, or instructing the **make** command to run silently.

m4 macro processor overview

This topic provides information about the **m4** macro processor, which is a front-end processor for any programming language being used in the operating system environment.

At the beginning of a program, you can define a symbolic name or symbolic constant as a particular string of characters. You can then use the **m4** macro processor to replace unquoted occurrences of the symbolic name with the corresponding string. Besides replacing one string of text with another, the **m4** macro processor provides the following features:

- Arithmetic capabilities
- File manipulation
- Conditional macro expansion
- String and substrings functions

The **m4** macro processor processes strings of letters and digits called *tokens*. The **m4** macro processor reads each alphanumeric token and determines if it is the name of a macro. The program then replaces the name of the macro with its defining text, and pushes the resulting string back onto the input to be rescanned. You can call macros with arguments, in which case the arguments are collected and substituted into the right places in the defining text before the defining text is rescanned.

The **m4** macro processor provides built-in macros such as **define**. You can also create new macros. Built-in and user-defined macros work the same way.

Using the m4 macro processor

To use the **m4** macro processor, enter the following command:

```
m4 [file]
```

The **m4** macro processor processes each argument in order. If there are no arguments or if an argument is - (dash), **m4** macro processor reads standard input as its input file. The **m4** macro processor writes its results to standard output. Therefore, to redirect the output to a file for later use, use a command such as the following:

```
m4 [file] >outputfile
```

Creating a user-defined macro

Macro	Description
define (<i>MacroName</i> , <i>Replacement</i>)	Defines new macro <i>MacroName</i> with a value of <i>Replacement</i> .

For example, if the following statement is in a program:

```
define(name, stuff)
```

The **m4** macro processor defines the string name as *stuff*. When the string name occurs in a program file, the **m4** macro processor replaces it with the string *stuff*. The string name must be ASCII alphanumeric and must begin with a letter or underscore. The string *stuff* is any text, but if the text contains parentheses, the number of open, or left, parentheses must equal the number of closed, or right, parentheses. Use the / (slash) character to spread the text for *stuff* over multiple lines.

The open (left) parenthesis must immediately follow the word **define**. For example:

```
define(N, 100)
if (i > N)
```

defines *N* to be 100 and uses the symbolic constant *N* in a later **if** statement.

Macro calls in a program have the following form:

```
name(arg1,arg2, . . . argn)
```

A macro name is recognized only if it is surrounded by non-alphanumeric characters. In the following example, the variable *NNN* is not related to the defined macro *N*.

```
define(N, 100)
if (NNN > 100)
```

You can define macros in terms of other names. For example:

```
define(N, 100)
define(M, N)
```

defines both M and N to be 100. If you later change the definition of N and assign it a new value, M retains the value of 100, not N.

The **m4** macro processor expands macro names into their defining text as soon as possible. The string N is replaced by 100. Then the string M is also replaced by 100. The overall result is the same as using the following input initially.

```
define(M, 100)
```

The order of the definitions can be interchanged as follows:

```
define(M, N)
define(N, 100)
```

Now M is defined to be the string N, so when the value of M is requested later, the result is the value of N at that time (because the M is replaced by N, which is replaced by 100).

Using the quote characters

To delay the expansion of the arguments of **define**, enclose them in quote characters. If you do not change them, quote characters are ` and ' (left and right single quotes). Any text surrounded by quote characters is not expanded immediately, but quote characters are removed. The value of a quoted string is the string with the quote characters removed. If the input is:

```
define(N, 100)
define(M, `N`)
```

The quote characters around the N are removed as the argument is being collected. The result of using quote characters is to define M as the string N, not 100. The general rule is that the **m4** macro processor always strips off one level of quote characters whenever it evaluates something. This is true even outside of macros. To make the word `define` appear in the output, enter the word in quote characters, as follows:

```
`define' = 1;
```

Another example of using quote characters is redefining N. To redefine N, delay the evaluation by putting N in quote characters. For example:

```
define(N, 100)
define(`N', 200)
```

To prevent problems from occurring, quote the first argument of a macro. For example, the following fragment does not redefine N:

```
define(N, 100)
define(N, 200)
```

The N in the second definition is replaced by 100. The result is the same as the following statement:

```
define(100, 200)
```

The **m4** macro processor ignores this statement because it can define only names, not numbers.

Changing the quote characters

Quote characters are normally ` and ' (left or right single quotes). If those characters are not convenient, change the quote characters with the following built-in macro:

Macro**changequote** (*l*, *r*)**Description**Changes the left and right quote characters to the characters represented by the *l* and *r* variables.

To restore the original quote characters, use **changequote** without arguments as follows:

```
changequote
```

Arguments

The simplest form of macro processing is replacing one string by another (fixed) string. However, macros can also have arguments, so that you can use the macro in different places with different results. To indicate where an argument is to be used within the replacement text for a macro (the second argument of its definition), use the symbol **\$n** to indicate the *n*th argument. When the macro is used, the **m4** macro processor replaces the symbol with the value of the indicated argument. For example, the following symbol:

```
$2
```

Refers to the second argument of a macro. Therefore, if you define a macro called bump as:

```
define(bump, $1 = $1 + 1)
```

The **m4** macro processor generates code to increment the first argument by 1. The bump(*x*) statement is equivalent to $x = x + 1$.

A macro can have as many arguments as needed. However, you can access only nine arguments using the **\$n** symbol (\$1 through \$9). To access arguments past the ninth argument, use the **shift** macro.

Macro**shift** (*ParameterList*)**Description**Returns all but the first element of *ParameterList* to perform a destructive left shift of the list.

This macro drops the first argument and reassigns the remaining arguments to the **\$n** symbols (second argument to \$1, third argument to \$2. . . tenth argument to \$9). Using the **shift** macro more than once allows access to all arguments used with the macro.

The **\$0** macro returns the name of the macro. Arguments that are not supplied are replaced by null strings, so that you can define a macro that concatenates its arguments as follows:

```
define(cat, $1$2$3$4$5$6$7$8$9)
```

Thus:

```
cat(x, y, z)
```

is the same as:

```
xyz
```

Arguments \$4 through \$9 in this example are null because corresponding arguments were not provided.

The **m4** macro processor discards leading unquoted blanks, tabs, or new-line characters in arguments, but keeps all other white space. Thus:

```
define(a, b c)
```

defines a to be b c.

Arguments are separated by commas. Use parentheses to enclose arguments containing commas, so that the comma does not end the argument. For example:

```
define(a, (b,c))
```

has only two arguments. The first argument is *a*, and the second is *(b, c)*. To use a comma or single parenthesis, enclose it in quote characters.

Using a predefined m4 macro

The **m4** macro processor provides a set of predefined macros. This section explains many of the macros and their uses.

Removing a macro definition

Macro	Description
undefine (<i>'MacroName'</i>)	Removes the definition of a user-defined or built-in macro (<i>'MacroName'</i>)

For example:

```
undefine('N')
```

removes the definition of *N*. After you remove a built-in macro with the **undefine** macro, as follows:

```
undefine('define')
```

You cannot use the definition of the built-in macro again.

Single quotes are required in this case to prevent substitution.

Checking for a defined macro

Macro	Description
ifdef (<i>'MacroName'</i> , <i>Argument1</i> , <i>Argument2</i>)	If macro <i>MacroName</i> is defined and is not defined to zero, returns the value of <i>Argument1</i> . Otherwise, it returns <i>Argument2</i> .

The **ifdef** macro permits three arguments. If the first argument is defined, the value of **ifdef** is the second argument. If the first argument is not defined, the value of **ifdef** is the third argument. If there is no third argument, the value of **ifdef** is null.

Using integer arithmetic

The **m4** macro processor provides the following built-in functions for doing arithmetic on integers only:

Macro	Description
incr (<i>Number</i>)	Returns the value of <i>Number</i> + 1.
decr (<i>Number</i>)	Returns the value of <i>Number</i> - 1.
eval	Evaluates an arithmetic expression.

Thus, to define a variable as one more than the *Number* value, use the following:

```
define(Number, 100)  
define(Number1, 'incr(Number)')
```

This defines *Number1* as one more than the current value of *Number*.

The **eval** function can evaluate expressions containing the following operators (listed in descending order of precedence) :

- unary + and -
- ** or ^ (exponentiation)

* / % (modulus)
+ -
== != < <= > >=
!(not)
& or && (logical AND)
| or || (logical OR)

Use parentheses to group operations where needed. All operands of an expression must be numeric. The numeric value of a true relation (for example, $1 > 0$) is 1, and false is 0. The precision of the **eval** function is 32 bits.

For example, define M to be $2==N+1$ using the **eval** function as follows:

```
define(N, 3)
define(M, `eval(2==N+1)')
```

Unless the text is very simple, use quote characters around the text that defines a macro

Manipulating files

To merge a new file in the input, use the built-in **include** function.

Macro	Description
include (<i>File</i>)	Returns the contents of the file <i>File</i> .

For example:

```
include(FileName)
```

inserts the contents of *FileName* in place of the **include** command.

A fatal error occurs if the file named in the **include** macro cannot be accessed. To avoid a fatal error, use the alternate form, the **sinclude** macro (silent include).

Macro	Description
sinclude (<i>File</i>)	Returns the contents of the file <i>File</i> , but does not report an error if it cannot access <i>File</i> .

The **sinclude** (silent include) macro does not write a message, but continues if the file named cannot be accessed.

Redirecting output

The output of the **m4** macro processor can be redirected again to temporary files during processing, and the collected material can be output upon command. The **m4** macro processor maintains nine possible temporary files, numbered 1 through 9. If you use the built-in **divert** macro.

Macro	Description
divert (<i>Number</i>)	Changes output stream to the temporary file <i>Number</i> .

The **m4** macro processor writes all output from the program after the **divert** function at the end of temporary file, *Number*. To return the output to the display screen, use either the **divert** or **divert(0)** function, which resumes the normal output process.

The **m4** macro processor writes all redirected output to the temporary files in numerical order at the end of processing. The **m4** macro processor discards the output if you redirect the output to a temporary file other than 0 through 9.

To bring back the data from all temporary files in numerical order, use the built-in **undivert** macro.

Macro**undivert** (*Number1, Number2...*)**Description**

Appends the contents of the indicated temporary files to the current temporary file.

To bring back selected temporary files in a specified order, use the built-in **undivert** macro with arguments. When using the **undivert** macro, the **m4** macro processor discards the temporary files that are recovered and does not search the recovered data for macros.

The value of the **undivert** macro is not the diverted text.

You can use the **divnum** macro to determine which temporary file is currently in use.

Macro**divnum****Description**

Returns the value of the currently active temporary file.

If you do not change the output file with the **divert** macro, the **m4** macro processor puts all output in a temporary file named 0.

Using system programs in a program

You can run any program in the operating system from a program by using the built-in **syscmd** macro. For example, the following statement runs the **date** program:

```
syscmd(date)
```

Using unique file Nnames

Use the built-in **maketemp** macro to make a unique file name from a program.

Macro**maketemp** (*String...nnnnn...String*)**Description**

Creates a unique file name by replacing the characters *nnnnn* in the argument string with the current process ID.

For example, for the statement:

```
maketemp(myfilennnnn)
```

the **m4** macro processor returns a string that is `myfile` concatenated with the process ID. Use this string to name a temporary file.

Using conditional expressions

Conditional expression evaluation allows process time determination of macro expressions.

Expression**ifelse** (*String1, String2, Argument1, Argument2*)**Description**

If *String1* matches *String2*, returns the value of *Argument1*. Otherwise it returns *Argument2*.

The built-in **ifelse** macro performs conditional testing. In the simplest form:

```
ifelse(a, b, c, d)
```

compares the two strings a and b.

If `a` and `b` are identical, the built-in **ifelse** macro returns the string `c`. If they are not identical, it returns string `d`. For example, you can define a macro called `compare` to compare two strings and return `yes` if they are the same, or `no` if they are different, as follows:

```
define(compare, `ifelse($1, $2, yes, no)`)
```

The quote characters prevent the evaluation of the **ifelse** macro from occurring too early. If the fourth argument is missing, it is treated as empty.

The **ifelse** macro can have any number of arguments, and therefore, provides a limited form of multiple-path decision capability. For example:

```
ifelse(a, b, c, d, e, f, g)
```

This statement is logically the same as the following fragment:

```
if(a == b) x = c;
else if(d == e) x = f;
else x = g;
return(x);
```

If the final argument is omitted, the result is null, so:

```
ifelse(a, b, c)
```

is `c` if `a` matches `b`, and null otherwise.

Manipulating strings

The macros in this section allow you to convert input strings into output strings.

Macro	Description
len	Returns the byte length of the string that makes up its argument

Thus:

```
len(abcdef)
```

is 6, and:

```
len((a,b))
```

is 5.

Macro	Description
dlen	Returns the length of the displayable characters in a string

Characters made up from 2-byte codes are displayed as one character. Thus, if the string contains any 2-byte, international character-support characters, the results of **dlen** will differ from the results of **len**.

Macro	Description
substr (<i>String, Position, Length</i>)	Returns a substring of <i>String</i> that begins at character number <i>Position</i> and is <i>Length</i> characters long.

Using input, **substr** (`s, i, n`) returns the substring of `s` that starts at the `i`th position (origin zero) and is `n` characters long. If `n` is omitted, the rest of the string is returned. For example, the function:

```
substr(`now is the time`,1)
```

returns the following string:

```
ow is the time
```

Macro

index (*String1*, *String2*)

Description

Returns the character position in *String1* where *String2* starts (starting with character number 0), or -1 if *String1* does not contain *String2*.

As with the built-in **substr** macro, the origin for strings is 0.

Macro

translit (*String*, *Set1*, *Set2*)

Description

Searches *String* for characters that are in *Set1*. If it finds any, changes (transliterates) those characters to corresponding characters in *Set2*.

It has the general form:

```
translit(s, f, t)
```

which modifies *s* by replacing any character found in *f* by the corresponding character of *t*. For example, the function:

```
translit('little', aeiou, 12345)
```

replaces the vowels by the corresponding digits and returns the following:

```
l3ttl2
```

If *t* is shorter than *f*, characters that do not have an entry in *t* are deleted. If *t* is not present at all, characters from *f* are deleted from *s*. So:

```
translit('little', aeiou)
```

deletes vowels from string *little* and returns the following:

```
lttl
```

Macro

dnl

Description

Deletes all characters that follow it, up to and including the new-line character.

Use this macro to get rid of empty lines. For example, the function:

```
define(N, 100)  
define(M, 200)  
define(L, 300)
```

results in a new-line at the end of each line that is not part of the definition. These new-line characters are passed to the output. To get rid of the new lines, add the built-in **dnl** macro to each of the lines.

```
define(N, 100) dnl  
define(M, 200) dnl  
define(L, 300) dnl
```

Debugging M4 macros

The macros in this section allow you to report errors and processing information.

Macro	Description
errprint (<i>String</i>)	Writes its argument (<i>String</i>) to the standard error file

For example:

```
errprint ('error')
```

Macro	Description
dumpdef (<i>MacroName'...</i>)	Dumps the current names and definitions of items named as arguments (<i>MacroName'...</i>)

If you do not supply arguments, the **dumpdef** macro prints all current names and definitions. Remember to quote the names.

Additional m4 macros

A list of additional **m4** macros, with a brief explanation of each, follows:

Macro	Description
changecom (<i>l, r</i>)	Changes the left and right comment characters to the characters represented by the <i>l</i> and <i>r</i> variables.
defn (<i>MacroName</i>)	Returns the quoted definition of <i>MacroName</i>
en (<i>String</i>)	Returns the number of characters in <i>String</i> .
m4exit (<i>Code</i>)	Exits m4 macro processor with a return code of <i>Code</i> .
m4wrap (<i>MacroName</i>)	Runs macro <i>MacroName</i> at the end of m4 macro processor .
popdef (<i>MacroName</i>)	Replaces the current definition of <i>MacroName</i> with the previous definition saved with the pushdef macro.
pushdef (<i>MacroName, Replacement</i>)	Saves the current definition of <i>MacroName</i> and then defines <i>MacroName</i> to be <i>Replacement</i> .
sysval	Gets the return code from the last use of the syscmd macro.
traceoff (<i>MacroList</i>)	Turns off trace for any macro in <i>MacroList</i> . If <i>MacroList</i> is null, turns off all tracing.
traceon (<i>MacroName</i>)	Turns on trace for macro <i>MacroName</i> . If <i>MacroName</i> is null, turns trace on for all macros.

Object data manager

Object Data Manager (ODM) is a data manager intended for storing system information. Information is stored and maintained as objects with associated characteristics.

You can also use ODM to manage data for application programs.

System data managed by ODM includes:

- Device configuration information
- Display information for SMIT (menus, selectors, and dialogs)
- Vital product data for installation and update procedures
- Communications configuration information
- System resource information

You can create, add, lock, store, change, get, show, delete, and drop objects and object classes with ODM. ODM commands provide a command-line interface to these functions. ODM subroutines access these functions from within an application program.

Some object classes are provided with the system. These object classes are discussed in the documentation for the specific system products that provide them.

This topic contains the following subtopics:

- [ODM Object Classes and Objects Storage](#)
- [ODM Descriptor](#)
- [ODM Object Searches](#)
- [ODM Commands and Subroutines](#)
- [ODM Example Code and Output](#)

ODM object classes and object storage

The basic components of ODM are object classes and objects. To manage object classes and objects, you use the ODM commands and subroutines. Specifically, you use the create and add features of these interfaces to build object classes and objects for storage and management of your own data.

Term	Description
object class	A group of objects with the same definition. An object class comprises one or more descriptors. It is similar to a table. Each object class you create with an odmcreate command or odm_create_class subroutine is stored in a file as a C language definition of an array of structures.
object	<p>A member of an object class. It is an entity that requires storage and management of data. It is similar to a logical record in a database.</p> <p>Each object you add to the object class with an odmadd command or an odm_add_obj subroutine is stored as a C language structure in the same file. You determine the directory in which to store this file when you create the object class.</p>

An object class is conceptually similar to an array of structures, with each object being a structure that is an element of the array. Values are associated with the descriptors of an object when the object is added to an object class. The descriptors of an object and their associated values can be located and changed with ODM facilities.

The following is an example of manipulating object classes and objects.

1. To create an object class called `Fictional_Characters`, enter:

```
class Fictional_Characters {
    char    Story_Star[20];
    char    Birthday[20];
    short   Age;
```



```
};
    char    Friend[20];
```

In this example, the `Fictional_Characters` object class contains four descriptors: `Story_Star`, `Birthday`, and `Friend`, which have a descriptor type of character and a 20-character maximum length; and `Age`, with a descriptor type of short. To create the object class files required by ODM, you process this file with the **odmcreate** command or the **odm_create_class** subroutine.

2. Once you create an object class, you can add objects to the class using the **odmadd** command or the **odm_add_obj** subroutine. For example, enter the following code with the **odmadd** command to add the objects Cinderella and Snow White to the `Fictional_Characters` object class, along with values for the descriptors they inherit:

```
Fictional_Characters:
    Story_Star = "Cinderella"
    Birthday   = "Once upon a time"
    Age        = 19
    Friend     = "mice"

Fictional_Characters:
    Story_Star = "Snow White"
    Birthday   = "Once upon a time"
    Age        = 18
    Friend     = "Fairy Godmother"
```

The `Fictional_Characters` table shows a conceptual picture of the `Fictional_Characters` object class with the two added objects Cinderella and Snow White.

Story Star (char)	Birthday (char)	Age (short)	Friend (char)
Cinderella	Once upon a time	19	Mice
Snow White	Once upon a time	18	Fairy Godmother

```
Retrieved data for 'Story_Star = "Cinderella"'
Cinderella:
    Birthday = Once upon a time
    Age      = 19
    Friend   = Mice
```

3. After the `Fictional_Characters` object class is created and the objects Cinderella and Snow White are added, the retrieved data for `'Story_Star = "Cinderella"'` is:

```
Cinderella:
    Birthday = Once upon a time
    Age      = 19
    Friend   = mice
```

Using ODM commands

When using the **odmcreate** or **odmdrop** command to create or drop an object class, specify the directory from which the class definition file will be accessed using one of the following:

1. Store the file in the default directory indicated by **\$ODMDIR**, which is the `/etc/objrepos` directory.
2. Use the **export** command to set the **ODMDIR** environment variable to specify a directory for storage.
3. Use the **unset** command to unset the **ODMDIR** environment variable and the **cd** command to change the current directory to the one in which you want the object classes stored. Then, run the ODM commands in that directory. The file defining the object classes will be stored in the current directory.

When using the **odmdelete**, **odmadd**, **odmchange**, **odmshow**, or **odmget** command to work with classes and objects, specify the directory containing the object classes using one of the following:

1. Work with object classes in the default directory indicated by **\$ODMDIR**, which is the `/etc/objrepos` directory.
2. Use the **export** command to set the **ODMDIR** environment variable to specify the directory.

3. From the command line, use the **export** command to set the **ODMPATH** environment variable to a string containing a colon-separated list of directories to be searched for object classes. For example:

```
$ export ODMPATH = /usr/lib/objrepos:/tmp/myrepos
```

The directories in the **\$ODMPATH** are searched only if the directory indicated by **\$ODMDIR** does not contain the object class.

Creating an object class

Attention: Making changes to files that define system object classes and objects can result in system problems. Consult your system administrator before using the **/usr/lib/objrepos** directory as a storage directory for object classes and objects.

1. Create the definition for one or more object classes in an ASCII file. “ODM example code and output” on page 553 shows an ASCII file containing several object class definitions.
2. Specify the directory in which the generated object must be stored.

“ODM Object Class and Object Storage” discusses the criteria used at object-class creation time for determining the directory in which to store generated object classes and objects. Most system object classes and objects are stored in the **/usr/lib/objrepos** directory.

Generate an empty object class by running the **odmcreate** command with the ASCII file of object class definitions specified as the *ClassDescriptionFile* input file.

Adding objects to an object class

Attention: Making changes to files that define system object classes and objects can result in system problems. Consult your system administrator before using the **/usr/lib/objrepos** directory as a storage directory for object classes and objects.

1. Create the object class to which the objects will be added. See [Creating an Object Class](#) for instructions on creating an object class.
2. Create the definitions for one or more objects. “ODM example code and output” on page 553 shows an ASCII file containing several object definitions.
3. Specify the directory which contains the object class in which the generated objects will be stored.

“ODM Object Class and Object Storage” discusses the criteria used at object class creation time for determining the directory in which to store generated object classes and objects. Most system object classes and objects are stored in the **/usr/lib/objrepos** directory.

4. Add objects to an object class by running the **odmadd** command with the ASCII file of object definitions specified as the *InputFile* input file.

Locking object classes

ODM does not implicitly lock object classes or objects. The coordination of locking and unlocking is the responsibility of the applications accessing the object classes. However, ODM provides the **odm_lock** and **odm_unlock** subroutines to control locking and unlocking object classes by application programs.

Subroutine

odm_lock

Description

Processes a string that is a path name and can resolve in an object class file or a directory of object classes. It returns a lock identifier and sets a flag to indicate that the specified object class or classes defined by the path name are in use.

When the **odm_lock** subroutine sets the lock flag, it does not disable use of the object class by other processes. If usage collision is a potential problem, an application program should explicitly wait until it is granted a lock on a class before using the class.

Another application cannot acquire a lock on the same path name while a lock is in effect. However, a lock on a directory name does not prevent another application from acquiring a lock on a subdirectory or the files within that directory.

To unlock a locked object class, use an **odm_unlock** subroutine called with the lock identifier returned by the **odm_lock** subroutine.

Storing object classes and objects

Each object class you create with an **odmcreate** command or **odm_create_class** subroutine is stored in a file as a C language definition of an array of structures. Each object you add to the object class with an **odmadd** command or an **odm_add_obj** subroutine is stored as a C language structure in the same file.

You determine the directory in which to store this file when you create the object class.

Storage methods vary according to whether commands or subroutines are used to create object classes and objects.

Attention: Making changes to files that define system object classes and objects can result in system problems. Consult your system administrator before using the **/usr/lib/objrepos** directory as a storage directory for object classes and objects.

Using the **odm_create_class** or **odm_add_obj** subroutines

The **odm_create_class** or **odm_add_obj** subroutine is used to create object classes and objects:

- If there is a specific requirement for your application to store object classes other than specified by the **ODMDIR** environment variable, use the **odm_set_path** subroutine to reset the path. It is strongly recommended that you use this subroutine to set explicitly the storage path whenever creating object classes or objects from an application.

OR

- Before running your application, use the **set** command from the command line to set the **ODMDIR** environment variable to specify a directory for storage.

OR

- Store the file in the object repository used to store many of the system object classes, the **/usr/lib/objrepos** directory.

ODM Descriptors

An Object Data Manager (ODM) descriptor is conceptually similar to a variable with a name and type. When an object class is created, its descriptors are defined like variable names with associated ODM descriptor types. When an object is added to an object class, it gets a copy of all of the descriptors of the object class. Values are also associated with object descriptors already stated.

ODM supports several descriptor types:

Descriptor	Definition
terminal descriptor	Defines a character or numerical data type.
link descriptor	Defines a relationship between object classes
method descriptor	Defines an operation or method for an object

Use the descriptors of an object and their associated values to define criteria for retrieving individual objects from an object class. Format the selection criteria you pass to ODM as defined in [ODM Object Searches](#). Do not use the **binary** terminal descriptor in search criteria because of its arbitrary length.

ODM terminal descriptors

Terminal descriptors define the most primitive data types used by ODM. A terminal descriptor is basically a variable defined with an ODM terminal descriptor type. The terminal descriptor types provided by ODM are:

Descriptor	Definition
short	Specifies a signed 2-byte number.
long	Specifies a signed 4-byte number.
ulong	Specifies an unsigned 4-byte number.
binary	Specifies a fixed-length bit string. The binary terminal descriptor type is defined by the user at ODM creation time. The binary terminal descriptor type cannot be used in selection criteria.
char	Specifies a fixed-length, null-terminated string.
vchar	Specifies variable-length, null-terminated string. The vchar terminal descriptor type can be used in selection criteria.
long64/ODM_LONG_LONG/int64	Specifies a signed 8-byte number.
ulong64/ODM_ULONG_LONG/uint64	Specifies an unsigned 8-byte number.

ODM link descriptor

The ODM *link descriptor* establishes a relationship between an object in an object class and an object in another object class. A link descriptor is a variable defined with the ODM link descriptor type.

For example, the following code can be processed by the ODM create facilities to generate the **Friend_Table** and **Fictional_Characters** object classes:

```
class Friend_Table {
    char    Friend_of[20];
    char    Friend[20];
};

class Fictional_Characters {
    char    Story_Star[20];
    char    Birthday[20];
    short   Age;
    link    Friend_Table Friend_Table Friend_of Friends_of;
};
```

The **Fictional_Characters** object class uses a link descriptor to make the **Friends_of** descriptors link to the **Friend_Table** object class. To resolve the link, the **Friends_of** descriptor retrieves objects in the **Friend_Table** object class with matching data in its **Friend_of** descriptors. The link descriptor in the **Fictional_Characters** object class defines the class being linked to (**Friend_Table**), the descriptor being linked to (**Friend_of**), and the name of the link descriptor (**Friends_of**) in the **Fictional_Characters** object class.

The following code could be used to add objects to the **Fictional_Characters** and **Friend_Table** object classes:

```
Fictional_Characters:
    Story_Star      = "Cinderella"
    Birthday       = "Once upon a time"
    Age            = 19
    Friends_of     = "Cinderella"
```

```
Fictional_Characters:
    Story_Star      = "Snow White"
    Birthday       = "Once upon a time"
    Age            = 18
    Friends_of     = "Snow White"
```

```
Friend_Table:
    Friend_of      = "Cinderella"
    Friend         = "Fairy Godmother"
```

```
Friend_Table:
  Friend_of      = "Cinderella"
  Friend         = "mice"
```

```
Friend_Table:
  Friend_of      = "Snow White"
  Friend         = "Sneezy"
```

```
Friend_Table:
  Friend_of      = "Snow White"
  Friend         = "Sleepy"
```

```
Friend_Table:
  Friend_of      = "Cinderella"
  Friend         = "Prince"
```

```
Friend_Table:
  Friend_of      = "Snow White"
  Friend         = "Happy"
```

The following tables show a conceptual picture of the **Fictional_Characters** and **Friend_Table** object classes, the objects added to the classes, and the link relationship between them.

Story_Star (char)	Birthday (char)	Age (short)	Friends_of (link)
Cinderella	Once upon a time	19	Cinderella
Snow White	Once upon a time	18	Snow White

```
Retrieved data for 'Story_Star = "Cinderella"
Cinderella:
  Birthday = Once upon a time
  Age      = 19
  Friends_of = Cinderella
  Friend_of = Cinderella
```

There is a direct link between the "**Friends_of**" and "**Friend_of**" columns of the two tables. The following table provides a conceptual picture of a link relationship between two object classes.

Friend_of (char)	Friend (char)
Cinderella	Fairy Godmother
Cinderella	mice
Snow White	Sneezy
Snow White	Sleepy
Cinderella	Prince
Snow White	Happy

After the **Fictional_Characters** and **Friend_Table** object classes are created and the objects are added, the retrieved data for `Story_Star = 'Cinderella'` would be:

```
Cinderella:
  Birthday      = Once upon a time
  Age           = 19
  Friends_of    = Cinderella
  Friend_of     = Cinderella
```

To see the expanded relationship between the linked object classes, use the **odmget** command on the **Friend_Table** object class. The retrieved data for the `Friend_of = 'Cinderella'` object class would be:

```

Friend_Table:
  Friend_Of = "Cinderella"
  Friend = "Fairy Godmother"

Friend_Table:
  Friend_of = "Cinderella"
  Friend= "mice"

Friend_Table:
  Friend_of = "Cinderella"
  Friend = "Prince"

```

ODM method descriptor

The ODM *method descriptor* gives the definition of an object class with objects that can have associated methods or operations. A method descriptor is a variable defined with the ODM method descriptor type.

The operation or method descriptor value for an object is a character string that can be any command, program, or shell script run by method invocation. A different method or operation can be defined for each object in an object class. The operations themselves are not part of ODM; they are defined and coded by the application programmer.

The method for an object is called by a call to the [odm_run_method](#) subroutine. The invocation of a method is a synchronous event, causing ODM operation to pause until the operation is completed.

For example, the following code can be input to the ODM create facilities to generate the **Supporting_Cast_Ratings** object class:

```

class Supporting_Cast_Ratings {
  char   Others[20];
  short  Dexterity;
  short  Speed;
  short  Strength;
  method Do_This;
};

```

In the example, the Do_This descriptor is a method descriptor defined for the **Supporting_Cast_Ratings** object class. The value of the method descriptor can be a string specifying a command, program, or shell script for future invocation by an [odm_run_method](#) subroutine.

The following code is an example of how to add objects to the **Supporting_Cast_Ratings** object class:

```

Supporting_Cast_Ratings:
  Others      = "Sleepy"
  Dexterity   = 1
  Speed       = 1
  Strength    = 3
  Do_This     = "echo Sleepy has speed of 1"

```

```

Supporting_Cast_Ratings:
  Others      = "Fairy Godmother"
  Dexterity   = 10
  Speed       = 10
  Strength    = 10
  Do_This     = "odmget -q "Others='Fairy Godmother'" Supporting_Cast_Ratings"

```

The following table shows a conceptual picture of the Supporting_Cast_Ratings object class with the Do_This method descriptor and operations associated with individual objects in the class.

Others (char)	Dexterity (short)	Speed (short)	Stength (short)	Do_This (method)
Sleepy	1	1	3	echo Sleepy has speed of 1
Fairy Godmother	10	10	10	odmget -q "Others='Fairy Godmother'" Supporting_Cast_Ratings"

```

odm_run_method run of Sleepy's method displays
(using echo):
"Sleepy has speed of 1"

```

After the `Supporting_Cast_Ratings` object class is created and the objects are added, an invocation (by the `odm_run_method` subroutine) of the method defined for `Sleepy` would cause the `echo` command to display:

```
Sleepy has speed of 1
```

ODM object searches

Many ODM routines require that one or more objects in a specified object class be selected for processing. You can include search criteria in the form of qualifiers when you select objects with certain routines.

qualifier

A null-terminated string parameter on ODM subroutine calls that gives the qualification criteria for the objects to retrieve

The descriptor names and qualification criteria specified by this parameter determine which objects in the object class are selected for later processing. Each qualifier contains one or more *predicates* connected with logical operators. Each predicate consists of a descriptor name, a comparison operator, and a constant.

A qualifier with three predicates joined by two logical operators follows:

```
SUPPNO=30 AND (PARTNO>0 AND PARTNO<101)
```

In this example, the entire string is considered the qualifier. The three predicates are `SUPPNO=30`, `PARTNO>0`, and `PARTNO<101`, and the `AND` logical operator is used to join the predicates. In the first predicate, `SUPPNO` is the name of the descriptor in an object, the `=` (equal sign) is a comparison operator, and `30` is the constant against which the value of the descriptor is compared.

Each predicate specifies a test applied to a descriptor that is defined for each object in the object class. The test is a comparison between the value of the descriptor of an object and the specified constant. The first predicate in the example shows an `=` (equal to) comparison between the value of a descriptor (`SUPPNO`) and a constant (`30`).

The part of the qualifier within parentheses

```
PARTNO>0 AND PARTNO<101
```

contains two predicates joined by the `AND` logical operator. The `PARTNO` descriptor is tested for a value greater than `0` in the first predicate, then tested for a value less than `101` in the second predicate. Then the two predicates are logically concatenated to determine a value for that part of the qualifier. For example, if `PARTNO` is the descriptor name for a part number in a company inventory, then this part of the qualifier defines a selection for all products with part numbers greater than `0` and less than `101`.

In another example, the qualifier:

```
Iname='Smith' AND Company.Dept='099' AND Salary<2500
```

can be used to select everyone (in ODM, every object) with a last name of Smith who is in Department `099` and has a salary less than `$2500`. Note that the `Dept` descriptor name is qualified with its `Company` object class to create a unique descriptor.

Descriptor names in ODM predicates

In ODM, a descriptor name is not necessarily unique. You can use a descriptor name in more than one object class. When this is the case, you specify the object class name along with the descriptor name in a predicate to create a unique reference to the descriptor.

Comparison operators in ODM predicates

The following are valid comparison operators:

Operator	Definition
<code>=</code>	Equal to

Operator	Definition
!=	Not equal to
>	Greater than
>=	Greater than or equal to
<	Less than
<=	Less than or equal to
LIKE	Similar to; finds patterns in character string data

Comparisons can be made only between compatible data types.

LIKE comparison operator

The LIKE operator enables searching for a pattern within a char descriptor type. For example, the predicate:

```
NAME LIKE 'ANNE'
```

defines a search for the value ANNE in the NAME descriptor in each object in the specified object class. In this case, the example is equivalent to:

```
NAME = 'ANNE'
```

You can also use the LIKE operator with the following pattern-matching characters and conventions:

- Use the ? (question mark) to represent any single character. The predicate example:

```
NAME LIKE '?A?'
```

defines a search for any three-character string that has A as a second character in the value of the NAME descriptor of an object. The descriptor values PAM, DAN, and PAT all satisfy this search criterion.

- Use the * (asterisk) to represent any string of zero or more characters. The predicate example:

```
NAME LIKE '*ANNE*'
```

defines a search for any string that contains the value ANNE in the NAME descriptor of an object. The descriptor values LIZANNE, ANNETTE, and ANNE all satisfy this search criterion.

- Use [] (brackets) to match any of the characters enclosed within the brackets. The predicate example:

```
NAME LIKE '[ST]*'
```

defines a search for any descriptor value that begins with S or T in the NAME descriptor of an object.

Use a - (minus sign) to specify a range of characters. The predicate example:

```
NAME LIKE '[AD-GST]*'
```

defines a search for any descriptor value that begins with any of the characters A, D, E, F, G, S, or T.

- Use [!] (brackets enclosing an exclamation mark) to match any single character except one of those enclosed within the brackets. The predicate example:

```
NAME LIKE '[!ST]*'
```

defines a search for any descriptor value except those that begin with S or T in the NAME descriptor of an object.

You can use the pattern-matching characters and conventions in any combination in the string.

Constants in ODM predicates

The specified constant can be either a numeric constant or a character string constant:

1. Numeric constants in ODM predicates consist of an optional sign followed by a number (with or without a decimal point), optionally followed by an exponent marked by the letter *E* or *e*. If used, the letter *E* or *e* must be followed by an exponent that can be signed.

Some valid numeric constants are:

```
2          2.545   0.5   -2e5   2.11E0
+4.555e-10 4E0    -10   999   +42
```

The E0 exponent can be used to specify no exponent.

2. Character string constants must be enclosed in single quotation marks:

```
'smith'   '91'
```

All character string constants are considered to have a variable length. To represent a single quotation mark inside a string constant, use two single quotation marks. For example:

```
'DON' 'T GO'
```

is interpreted as:

```
DON'T GO
```

AND logical operator for predicates

The AND logical operator can be used with predicates. Use AND or and for the AND logical operator.

The AND logical operator connects two or more predicates. The qualifier example:

```
predicate1 AND predicate2 AND predicate3
```

specifies predicate1 logically concatenated with predicate2 followed by the result, which is logically concatenated with predicate3.

ODM commands and subroutines

You can create, add, change, retrieve, display, delete, and remove objects and object classes with ODM. You enter ODM commands on the command line.

You can put ODM subroutines in a C language program to handle objects and object classes. An ODM subroutine returns a value of -1 if the subroutine is unsuccessful. The specific error diagnostic is returned as the **odmerrno** external variable (defined in the **odmi.h** include file). ODM error-diagnostic constants are also included in the **odmi.h** include file.

Note: If you are writing a C language program using the subroutines, use option:
-binitfini: __odm_initfini_init: __odm_initfini_fini.

Commands

ODM commands are:

Command	Description
<u>odmadd</u>	Adds objects to an object class. The odmadd command takes an ASCII stanza file as input and populates object classes with objects found in the stanza file.
<u>odmchange</u>	Changes specific objects in a specified object class.
<u>odmcreate</u>	Creates empty object classes. The odmcreate command takes an ASCII file describing object classes as input and produces C language .h and .c files to be used by the application accessing objects in those object classes.
<u>odmdelete</u>	Removes objects from an object class.

Command	Description
<u>odmdrop</u>	Removes an entire object class.
<u>odmshow</u>	Displays the description of an object class. The odmshow command takes an object class name as input and puts the object class information into odmcreate command format.
<u>odmget</u>	Retrieves objects from object classes and puts the object information into odmadd command format.

Subroutines

ODM subroutines are:

Subroutine	Description
<u>odm_add_obj</u>	Adds a new object to the object class.
<u>odm_change_obj</u>	Changes the contents of an object.
<u>odm_close_class</u>	Closes an object class.
<u>odm_create_class</u>	Creates an empty object class.
<u>odm_err_msg</u>	Retrieves a message string.
<u>odm_free_list</u>	Frees memory allocated for the odm_get_list subroutine.
<u>odm_get_by_id</u>	Retrieves an object by specifying its ID.
<u>odm_get_first</u>	Retrieves the first object that matches the specified criteria in an object class.
<u>odm_get_list</u>	Retrieves a list of objects that match the specified criteria in an object class.
<u>odm_get_next</u>	Retrieves the next object that matches the specified criteria in an object class.
<u>odm_get_obj</u>	Retrieves an object that matches the specified criteria from an object class.
<u>odm_initialize</u>	Initializes an ODM session.
<u>odm_lock</u>	Locks an object class or group of classes.
<u>odm_mount_class</u>	Retrieves the class symbol structure for the specified object class.
<u>odm_open_class</u>	Opens an object class.
<u>odm_rm_by_id</u>	Removes an object by specifying its ID.
<u>odm_rm_obj</u>	Removes all objects that match the specified criteria from the object class.
<u>odm_run_method</u>	Invokes a method for the specified object.
<u>odm_rm_class</u>	Removes an object class.
<u>odm_set_path</u>	Sets the default path for locating object classes.
<u>odm_unlock</u>	Unlocks an object class or group of classes.
<u>odm_terminate</u>	Ends an ODM session.

ODM example code and output

The following Fictional_Characters, Friend_Table, and Enemy_Table Object Classes and Relationships tables list the object classes and objects created by the example code in this section.

Story_Star (char)	Birthday (char)	Age (short)	Friends_of (link)	Enemies_of (link)	Do_This (method)
Cinderella	Once upon a time	19	Cinderella	Cinderella	echo Cleans House
Snow White	Once upon a time	18	Snow White	Snow White	echo Cleans House

Friend_of (char)	Friend (char)
Cinderella	Fairy Godmother
Cinderella	mice
Snow White	Sneezy
Snow White	Sleepy
Cinderella	Prince
Snow White	Happy

Enemy_of (char)	Enemy (char)
Cinderella	midnight
Cinderella	Mean Stepmother
Snow White	Mean Stepmother

ODM example input code for creating object classes

The following example code in the **MyObjects.cre** file creates three object classes when used as an input file with the **odmcreate** command:

```
*      MyObjects.cre
*      An input file for ODM create utilities.
*      Creates three object classes:
*          Friend_Table
*          Enemy_Table
*          Fictional_Characters
```

```
class  Friend_Table {
char   Friend_of[20];
char   Friend[20];
};
```

```
class  Enemy_Table {
char   Enemy_of[20];
char   Enemy[20];
};
```

```
class  Fictional_Characters {
char   Story_Star[20];
char   Birthday[20];
```

```

short   Age;
link    Friend_Table Friend_Table Friend_of Friends_of;
link    Enemy_Table Enemy_Table Enemy_of Enemies_of;
method  Do_This;
};

```

```
* End of MyObjects.cre input file for ODM create utilities. *
```

The `Fictional_Characters` object class contains six descriptors:

- `Story_Star` and `Birthday`, each with a descriptor type of `char` and a 20-character maximum length.
 - `Age` with a descriptor type of `short`.
 - `Friends_of` and `Enemies_of` are both from the link class and link to the two previously defined object classes.
- Note:** The object class `link` is repeated twice.
- `Do_This` with a descriptor type of `method`.

To generate the object class files required by ODM, the file containing this code must be processed with the `odmcreate` command

ODM example output for object class definitions

Processing the code in the `MyObjects.cre` file with the `odmcreate` command generates the following structures in a `.h` file:

```

* MyObjects.h
* The file output from ODM processing of the MyObjects.cre input
* file. Defines structures for the three object classes:
*   Friend_Table
*   Enemy_Table
*   Fictional_Characters
#include <odmi.h>

```

```

struct Friend_Table {
    long   _id;          * unique object id within object class *
    long   _reserved;   * reserved field *
    long   _scratch;    * extra field for application use *
    char   Friend_of[20];
    char   Friend[20];
};

```

```

#define Friend_Table_Descs 2
extern struct Class Friend_Table_CLASS[];
#define get_Friend_Table_list(a,b,c,d,e) (struct Friend_Table * )odm_get_list (a,b,c,d,e)

```

```

struct Enemy_Table {
    long   _id;
    long   _reserved;
    long   _scratch;
    char   Enemy_of[20];
    char   Enemy[20];
};
#define Enemy_Table_Descs 2
extern struct Class Enemy_Table_CLASS[];
#define get_Enemy_Table_list(a,b,c,d,e) (struct Enemy_Table * )odm_get_list (a,b,c,d,e)

```

```

struct Fictional_Characters {
    long   _id;
    long   _reserved;
    long   _scratch;
    char   Story_Star[20];
    char   Birthday[20];
    short  Age;
    struct Friend_Table *Friends_of;      * link *
    struct listinfo *Friends_of_info;    * link *
    char   Friends_of_Lvalue[20];        * link *
    struct Enemy_Table *Enemies_of;      * link *
    struct listinfo *Enemies_of_info;    * link *
    char   Enemies_of_Lvalue[20];        * link *
};

```

```

};          char    Do_This[256];          * method *

#define Fictional_Characters_Descs 6

extern struct Class Fictional_Characters_CLASS[];
#define get_Fictional_Characters_list(a,b,c,d,e) (struct Fictional_Characters * )odm_get_list
(a,b,c,d,e)

* End of MyObjects.h structure definition file output from ODM * processing.

```

ODM example code for adding objects to object classes

The following code can be processed by the **odmadd** command to populate the object classes created by the processing of the **MyObjects.cre** input file:

```

* MyObjects.add
* An input file for ODM add utilities.
* Populates three created object classes:
*     Friend_Table
*     Enemy_Table
*     Fictional_Characters

```

```

Fictional_Characters:
Story_Star = "Cinderella" #a comment for the MyObjects.add file.
Birthday   = "Once upon a time"
Age        = 19
Friends_of = "Cinderella"
Enemies_of = "Cinderella"
Do_This    = "echo Cleans house"

```

```

Fictional_Characters:
Story_Star = "Snow White"
Birthday   = "Once upon a time"
Age        = 18
Friends_of = "Snow White"
Enemies_of = "Snow White"
Do_This    = "echo Cleans house"

```

```

Friend_Table:
Friend_of    = "Cinderella"
Friend      = "Fairy Godmother"

```

```

Friend_Table:
Friend_of    = "Cinderella"
Friend      = "mice"

```

```

Friend_Table:
Friend_of    = "Snow White"
Friend      = "Sneezy"

```

```

Friend_Table:
Friend_of    = "Snow White"
Friend      = "Sleepy"

```

```

Friend_Table:
Friend_of    = "Cinderella"
Friend      = "Prince"

```

```

Friend_Table:
Friend_of    = "Snow White"
Friend      = "Happy"

```

```

Enemy_Table:
Enemy_of     = "Cinderella"
Enemy       = "midnight"

```

```
Enemy_Table:
Enemy_of      =      "Cinderella"
Enemy        =      "Mean Stepmother"
```

```
Enemy_Table:
Enemy_of      =      "Snow White"
Enemy        =      "Mean Stepmother"
```

```
* End of MyObjects.add input file for ODM add utilities. *
```

Note: The * (asterisk) or the # (pound sign) comment in the previous example code will not go into the object file. If the beginning of the line causes the line to be a comment, the command does not go in the object file. The comment will be included in the file and treated as a string if it is included inside the " " (double quotation marks).

Simultaneous multithreading

Simultaneous multithreading is the ability of a single physical processor to simultaneously dispatch instructions from more than one hardware thread context. Because there are two hardware threads per physical processor, additional instructions can run at the same time.

Simultaneous multithreading allows you to take advantage of the superscalar nature of the processor by scheduling two applications at the same time on the same processor. No single application can fully saturate the processor.

Benefitting from Simultaneous Multithreading

It is primarily beneficial in commercial environments where the speed of an individual transaction is not as important as the total number of transactions that are performed. Simultaneous multithreading is expected to increase the throughput of workloads with large or frequently changing working sets, such as database servers and Web servers.

Workloads that see the greatest simultaneous multithreading benefit are those that have a high Cycles Per Instruction (CPI) count. These workloads tend to use processor and memory resources poorly. Large CPIs are usually caused by high cache-miss rates from a large working set. Large commercial workloads are somewhat dependent upon whether the two hardware threads share instructions or data, or the hardware threads are completely distinct. Large commercial workloads typically have this characteristic. Workloads that share instructions or data, including those that run extensively in the operating system or within a single application, might see increased benefits from simultaneous multithreading.

Workloads that do not benefit much from simultaneous multithreading are those in which the majority of individual software threads use a large amount of any resource in the processor or memory. For example, workloads that are floating-point intensive are likely to gain little from simultaneous multithreading and are the ones most likely to lose performance. These workloads heavily use either the floating-point units or the memory bandwidth. Workloads with low CPI and low cache miss rates might see a some small benefit.

Measurements taken on a dedicated partition with commercial workloads indicated a 25%-40% increase in throughput. Simultaneous multithreading should help shared processor partition processing. The extra threads give the partition a boost after simultaneous multithreading is dispatched because the partition recovers its working set more quickly. Subsequently, the threads perform like they would in a dedicated partition. Although it might be somewhat counterintuitive, simultaneous multithreading performs best when the performance of the cache is at its worst.

Setting the mode with the `smtctl` command

AIX allows you to control the mode of the partition for simultaneous multithreading with the `smtctl` command. With this command, you can turn simultaneous multithreading on or off system-wide, either immediately or the next time the system boots. The simultaneous multithreading mode persists across system boots. By default, AIX enables simultaneous multithreading.

The syntax for the `smtctl` command is as follows:

```
smtctl [ -m { off | on } [ { -boot | -now } ] ]
```

Hardware Management Console Configuration for Simultaneous Multithreading

When you configure shared processor partitions at the Hardware Management Console (HMC), you specify the minimum, desired, and maximum number of virtual processors. For dedicated partitions, you specify the same type of parameters, but the processor terminology is different. For dedicated partitions, the processors are always called processors.

Both partitioning models require that you specify a range of processors that control the boot and runtime assignment of processors to the partition. If possible, the desired processor setting is granted when the system starts. If this is not possible, the POWER Hypervisor chooses a different value based on the set of available resources that is greater than, or equal to, the minimum value.

The number of processors specified at the HMC impacts the number of logical processors that AIX allocates. If the partition is capable of simultaneous multithreading, AIX allocates twice as many logical processors as the maximum processor value because there are two hardware threads per processor and AIX configures each hardware thread as a separate logical processor. This allows AIX to enable or disable simultaneous multithreading without rebooting the partition.

Dynamic Logical Partitioning for Simultaneous Multithreading

While a partition is running, you can change the number of processors that are assigned to a partition through Dynamic Logical Partitioning (DLPAR) procedures at the HMC. You can add or remove processors within the constraints of the processor range defined for the partition. When a processor is added to a partition that is enabled for simultaneous multithreading, AIX starts both hardware threads and two logical processors are brought online. When a processor is removed from a partition that is enabled for simultaneous multithreading, AIX stops both hardware threads and two logical processors are taken offline.

Two DLPAR events are generated when simultaneous multithreading is enabled. One event is generated for each of the logical processors that is added or removed. The API for DLPAR scripts is based on logical processors, so the number of DLPAR events parallels the addition and removal of logical processors. If simultaneous multithreading is not enabled in the partition, there is only one DLPAR event. AIX automatically translates the DLPAR request that is sent from the HMC into the appropriate number of DLPAR events presented to DLPAR-aware applications.

Micro-Partitioning[®] and Simultaneous Multithreading

The POWER Hypervisor[™] saves and restores all necessary processor states when preempting or dispatching virtual processors. For processors that are enabled for simultaneous multithreading, this means two active thread contexts. Each hardware thread is supported as a separate logical processor by AIX. For this reason, a dedicated partition that is created with one physical processor is configured by AIX as a logical 2-way processor. Because this is independent of the partition type, a shared partition with two virtual processors is configured by AIX as a logical 4-way processor and a shared partition with four virtual processors is configured by AIX as a logical 8-way processor. Paired threads are always scheduled together at the same time in the same partition.

Shared processor capacity is always delivered in terms of whole physical processors. Without simultaneous multithreading, AIX configures a 4-way virtual processor partition with 200 units of processor entitlement as a 4-way logical processor partition where each logical processor has the power of 50% of a physical processor. With simultaneous multithreading, the 4-way logical processor partition becomes an 8-way logical processor partition, where each logical processor has the approximate power of 25% of a physical processor. However, with simultaneous multithreading, latency concerns normally associated with a virtual processor's fractional capacity do not apply linearly to the threads. Because both threads are dispatched together, they are active for the duration of a 50% dispatch window and they share the underlying physical processor to achieve the logical 25%. This means that each of the logical processors are able to field interrupts for twice as long as their individual capacities allow.

Hardware thread priorities

The processor allows priorities to be assigned to hardware threads. The difference in priority between sibling threads determines the ratio of physical processor decode slots allotted to each thread. More slots provide better thread performance. Normally, AIX maintains sibling threads at the same priority, but raises or lowers thread priorities in key places to optimize performance. For example, AIX lowers thread priorities when the thread is doing nonproductive work spinning in the idle loop or on a kernel lock. Thread priorities are raised when a thread is holding a critical kernel lock. These priority adjustments do not persist in user mode. AIX does not consider a software thread's dispatching priority when it is choosing the hardware thread priority.

Work is distributed across all primary threads before work is dispatched to secondary threads. The performance of a thread is best when its paired thread is idle. Thread affinity is also considered in idle stealing and in periodic run queue load balancing.

Dynamic logical partitioning

Partitioning your system is similar to partitioning a hard disk drive. When you partition a hard disk drive, you divide a single physical hard disk drive so that the operating system recognizes it as a number of separate logical hard disk drives.

On each of these divisions, called partitions, you can install an operating system and use each partition as you would a separate physical system.

A *logical partition* (LPAR) is the division of a computer's processors, memory, and hardware resources into multiple environments so that each environment can be operated independently with its own operating system and applications. The number of logical partitions that can be created depends on the system. Typically, partitions are used for different purposes, such as database operation, client/server operations, Web server operations, test environments, and production environments. Each partition can communicate with the other partitions as if each partition were a separate machine.

Dynamic logical partitioning (DLPAR) provides the ability to logically attach and detach a managed system's resources to and from a logical partition's operating system without rebooting. Some of the features of DLPAR include:

- The *Capacity Upgrade on Demand (CUoD)* feature of the IBM System p, which allows customer to activate preinstalled but inactive processors as resource requirements change.
- The Dynamic Processor Deallocation feature of the IBM Power servers, and on some SMP models. Dynamic Processor Deallocation enables a processor to be taken offline dynamically when an internal threshold of recoverable errors is exceeded. DLPAR enhances the Dynamic Processor Deallocation feature by allowing an inactive processor, if one exists, to be substituted for the processor that is suspected of being defective. This online switch does not affect applications or kernel extensions.
- DLPAR enables cross-partition workload management, which is particularly important for server consolidation in that it can be used to manage system resources across partitions.

DLPAR requests are built from simple add and remove requests that are directed to logical partitions. The user can execute these commands as move requests at the Hardware Management Console (HMC), which manages all DLPAR operations. DLPAR operations are enabled by System p firmware and AIX.

Related information

[cpupstat](#)

[drmgr](#)

[dr_reconfig](#)

[reconfig](#)

DLPAR-safe and aware programs

A DLPAR-safe program is one that does not fail as a result of DLPAR operations.

Its performance might suffer when resources are removed and it might not scale with the addition of new resources, but the program still works as expected. In fact, a DLPAR-safe program can prevent a DLPAR operation from succeeding because it has a dependency that the operating system is obligated to honor.

A *DLPAR-aware* program is one that has DLPAR code that is designed to adjust its use of system resources as the actual capacity of the system varies over time. This can be accomplished in the following ways:

- By regularly polling the system in an attempt to discover changes in the system topology
- By registering application specific code that is notified when a change is occurring to the system topography

DLPAR-aware programs must be designed, at minimum, to avoid introducing conditions that might cause DLPAR operations to fail. At maximum, DLPAR-aware programs are concerned with performance and scalability. This is a much more complicated task because buffers might need to be drained and resized to maintain expected levels of performance when memory is added or removed. In addition, the number of threads must be dynamically adjusted to account for changes in the number of online processors. These thread-based adjustments are not necessarily limited to processor-based decisions. For example, the best way to reduce memory consumption in Java programs might be to reduce the number of threads, because this reduces the number of active objects that need to be processed by the Java Virtual Machine's garbage collector.

Most applications are DLPAR-safe by default.

Making programs DLPAR-safe

The following types of errors, which represent binary compatibility exposures, can be introduced by DLPAR:

Note: These errors are a result of processor addition.

- If a program has code that is optimized for uniprocessor systems and the number of processors in the partition is increased from one to two, programs that make runtime checks might take an unexpected path if a processor is added during one of these checks. Potential problems can also occur in programs that implement their own locking primitives, but do so using uniprocessor serialization techniques; that is, the `sync` and `isync` instructions are not included. The use of these instructions is also required for self-modifying and generated code, and are thus necessary on DLPAR-enabled systems. Be sure to look for uniprocessor-based logic. Programs that make uniprocessor assertions must include logic that identifies the number of online processors.

Programs can determine the number of online processors by:

- Loading the `_system_configuration.ncpus` field
 - `var.v_ncpus`
 - Using the **sysconf** system call with the **_SC_NPROCESSORS_ONLN** flag.
- Programs that index data by processor number typically use the **mycpu** system call to determine the identity of the currently executing processor, in order to index into their data structures. The problem potentially occurs when a new processor is added because the path to the data might not be properly initialized or allocated. Programs that preallocate processor-based lists using the number of online CPUs are broken, because this value changes with DLPAR.

Avoid this problem by preallocating processor-based data using the maximum possible number of processors that can be brought online at the same time. The operating system can be said to be configured to support a maximum of N processors, not that there are N processors active at any given time. The maximum number of processors is constant, while the number of online processors is incremented and decremented as processors are brought online and taken offline. When a partition is created, the minimum, desired, and maximum numbers of processors are specified. The maximum value is reflected in the following variables:

- `_system_configuration.max_ncpus`
- `_system_configuration.original_ncpus`
- `var.v_ncpus_cfg`
- **sysconf** (**_SC_NPROCESSORS_CONF**)

The `_system_configuration.original_ncpus` and `var.v_ncpus_cfg` variables are preexisting variables. On DLPAR-enabled systems they represent a potential maximum value. On systems not enabled for DLPAR,

the value is dictated by the number of processors that are configured at boot time. Both represent the conceptual maximum value that can be supported, even though a processor might have been taken offline by Dynamic Processor Deallocation. The use of these preexisting fields is recommended for applications that are built on AIX 4.3, because this facilitates the use of the same binary on AIX 4.3 and later. If the application requires runtime initialization of its processor-based data, it can register a DLPAR handler that is called before a processor is added.

Making programs DLPAR-aware

A *DLPAR-aware* program is one that is designed to recognize and dynamically adapt to changes in the system configuration. This code need not subscribe to the DLPAR model of awareness, but can be structured more generally in the form of a system resource monitor that regularly polls the system to discover changes in the system configuration. This approach can be used to achieve some limited performance-related goals, but because it is not tightly integrated with DLPAR, it cannot be effectively used to manage large-scale changes to the system configuration. For example, the polling model might not be suitable for a system that supports processor hot plug, because the hot-pluggable unit might be composed of several processor and memory boards. Nor can it be used to manage application-specific dependencies, such as processor bindings, that need to be resolved before the DLPAR processor remove event is started.

The following types of applications can exploit DLPAR technology:

- Applications that are designed to scale with the system configuration, including those:
 - That detect the number of online processors or the size of physical memory when the application starts
 - That are externally directed to scale based on an assumed configuration of processors and memory, which usually translates into the use of a maximum number of threads, maximum buffer sizes, or a maximum amount of pinned memory
- Applications that are aware of the number of online processors and the total quantity of system memory, including the following types of applications:
 - Performance monitors
 - Debugging tools
 - System crash tools
 - Workload managers
 - License managers

Note: Not all license managers are candidates for DLPAR, especially user-based license managers.

- Applications that pin their application data, text, or stack using the **plock** system call
- Applications that use System V Shared Memory Segments with the **PinvOption (SHM_PIN)**
- Applications that bind to processors using the **bindprocessor** system call

Dynamic logical partitioning of large memory pages is not supported. The amount of memory that is preallocated to the large page pool can have a material impact on the DLPAR capabilities of the partition regarding memory. A memory region that contains a large page cannot be removed. Therefore, application developers might want to provide an option to not use large pages.

Making programs DLPAR-aware using DLPAR APIs

Application interfaces are provided to make programs DLPAR-aware. The **SIGRECONFIG** signal is sent to the applications at the various phases of dynamic logical partitioning. The DLPAR subsystem defines check, pre and post phases for a typical operation. Applications can watch for this signal and use the DLPAR-supported system calls to learn more about the operation in progress and to take any necessary actions.

Note: When using signals, the application might inadvertently block the signal, or the load on the system might prevent the thread from running in a timely fashion. In the case of signals, the system will wait a short period of time, which is a function of the user-specified timeout, and proceed to the next phase.

It is not appropriate to wait indefinitely because a non-privileged rogue thread could prevent all DLPAR operations from occurring.

The issue of timely signal delivery can be managed by the application by controlling the signal mask and scheduling priority. The DLPAR-aware code can be directly incorporated into the algorithm. Also, the signal handler can be cascaded across multiple shared libraries so that notification can be incorporated in a more modular way.

To integrate the DLPAR event using APIs, do the following:

1. Catch the **SIGRECONFIG** signal by using the **sigaction** system call. The default action is to ignore the signal.
2. Control the signal mask in at least one of the threads so that the signal can be delivered in real time.
3. Ensure that the scheduling priority for the thread that is to receive the signal is sufficient so that it will run quickly after the signal has been sent.
4. Run the **dr_reconfig** system call to obtain the type of resource, type of action, phase of the event, as well as other information that is relevant to the current request.

Note: The **dr_reconfig** system call is used inside the signal handler to determine the nature of the DLPAR request.

Managing an application's DLPAR dependencies

A DLPAR remove request can fail for a variety of reasons. The most common of these is that the resource is busy, or that there are not enough system resources currently available to complete the request.

In these cases, the resource is left in a normal state as if the DLPAR event never happened.

The primary cause of *processor removal* failure is processor bindings. The operating system cannot ignore processor bindings and carry on DLPAR operations or applications might not continue to operate properly. To ensure that this does not occur, release the binding, establish a new one, or terminate the application. The specific process or threads that are impacted is a function of the type of binding that is used.

The primary cause of *memory removal* failure is that there is not enough pinned memory available in the system to complete the request. This is a system-level issue and is not necessarily the result of a specific application. If a page in the memory region to be removed has a pinned page, its contents must be migrated to another pinned page, while atomically maintaining its virtual to physical mappings. The failure occurs when there is not enough pinnable memory in the system to accommodate the migration of the pinned data in the region that is being removed. To ensure that this does not occur, lower the level of pinned memory in the system. This can be accomplished by destroying pinned shared memory segments, terminating programs that implement the **plock** system call, or removing the **plock** on the program.

The primary cause of *PCI slot removal* failure is that the adapters in the slot are busy. Note that device dependencies are not tracked. For example, the device dependency might extend from a slot to one of the following: an adapter, a device, a volume group, a logical volume, a file system, or a file. In this case, resolve the dependencies manually by stopping the relevant applications, unmounting file systems, and varying off volume groups.

Related concepts

[Processor bindings](#)

Applications can bind to a processor by using the **bindprocessor** system call. This system call assumes a processor-numbering scheme starting with zero (0) and ending with $N-1$, where N is the number of online CPUs.

Processor bindings

Applications can bind to a processor by using the **bindprocessor** system call. This system call assumes a processor-numbering scheme starting with zero (0) and ending with $N-1$, where N is the number of online CPUs.

N is determined programmatically by reading the **_system_configuration.ncpus** system variable. As processors are added and removed, this variable is incremented and decremented using dynamic logical partitioning.

Note that the numbering scheme does not include holes. Processors are always added to the N th position and removed from the $N-1$ position. The numbering scheme used by the **bindprocessor** system call cannot be used to bind to a specific logical processor, because any processor can be removed and this is not reflected in the numbering scheme, because the $N-1$ CPU is always deallocated. For this reason, the identifiers used by the **bindprocessor** system call are called *bind CPU IDs*.

Changes to the **_system_configuration.ncpus** system variables have the following implications:

- Applications must be prepared to receive an error from the **bindprocessor** system call if the last processor is removed after the applications have read the variable. This error condition was first introduced by the Dynamic Processor Deallocation (runtime deallocation of defective processors).
- Applications that are designed to scale with the number of processors must reread the **_system_configuration.ncpus** system variable when the number of processors changes.

Applications can also bind to a set of processors using a feature of Workload Manager (WLM) called *Software Partitioning*. It assumes a numbering scheme that is based on logical CPU IDs, which also start with zero (0) and end with $N-1$. However, N in this case is the maximum number of processors that can be supported architecturally by the partition. The numbering scheme reflects both online and offline processors.

Therefore, it is important to note the type of binding that is used so that the correct remedy can be applied when removing a processor. The **bindprocessor** command can be used to determine the number of online processors. The **ps** command can be used to identify the processes and threads that are bound to the last online processor. After the targets have been identified, the **bindprocessor** command can be used again to define new attachments.

WLM-related dependencies can be resolved by identifying the particular software partitions that are causing problems. To resolve these dependencies, do the following:

Note: The system schedules bound jobs around offline or pending offline processors, so no change is required if the particular software partition has another online CPU.

1. Use the **lsrset** command to view the set of software partitions that are used by WLM.
2. Identify these software partitions by using the **lsclass** command.
3. Identify the set of classes that use these software partitions by using the **chclass** command.
4. Reclassify the system using the **wlmctrl** command.

At this point, the new class definitions take effect and the system automatically migrates bound jobs away from the logical processor that is being removed.

Related concepts

[Managing an application's DLPAR dependencies](#)

A DLPAR remove request can fail for a variety of reasons. The most common of these is that the resource is busy, or that there are not enough system resources currently available to complete the request.

Integrating the DLPAR operation into the application

The DLPAR operation can be integrated into the application in the following ways:

- Using a script-based approach, where the user installs a set of DLPAR scripts into a directory. These scripts are invoked while the DLPAR operation is being run. The scripts are designed to externally reconfigure the application.
- Using the **SIGRECONFIG** signal, which is used to catch the signal of every process that has registered. The signal method assumes that the application has been coded to catch the signal, and that the signal handler will reconfigure the application. The signal handler invokes an interface to determine the nature of the DLPAR operation.

Both of these methods follow the same high level structure. Either method can be used to provide support for DLPAR, although only the script-based mechanism can be used to manage DLPAR dependencies related to Workload Manager software partitions (processor sets). No APIs are associated with Workload Manager, so the use of a signal handler is not a suitable vehicle for dealing with the Workload Manager-imposed scheduling constraints. The applications themselves are not Workload Manager-aware. In this case, the system administrator might want to provide a script that invokes Workload Manager commands to manage DLPAR interactions with Workload Manager.

The decision of which method to use should be based on how the system or resource-specific logic was introduced into the application. If the application was externally directed to use a specific number of threads or to size its buffers, use the script-based approach. If the application is directly aware of the system configuration and uses this information accordingly, use the signal-based approach.

The DLPAR operation itself is divided into the following phases:

- **check phase**

The *check phase* is invoked first and it enables applications to fail the current DLPAR request before any state in the system is changed. For example, the check phase could be used by a CPU-based license manager to fail the integration of a new processor if that CPU addition makes the number of processors in the system exceed the number of licensed processors. It could also be used to preserve the DLPAR safeness of a program that is not DLPAR-safe. In the latter case, consideration must be given to services provided by the application, because it might be better to stop the program, complete the request, and then restart the program.

- **pre phase** and **post phase**

The *pre phase* and *post phase* are provided to stop the program, complete the request, and then restart the program.

The system attempts to ensure that all of the check code across the different mediums is executed in its entirety at the system level before the DLPAR event advances to the next phase.

Actions taken by DLPAR scripts

Application scripts are started for both add and remove operations.

When removing resources, scripts are provided to resolve conditions imposed by the application that prevent the resource from being removed. The presence of particular processor bindings and the lack of pinnable memory might cause a remove request to fail. A set of commands is provided to identify these situations, so that scripts can be written to resolve them.

To identify and resolve the DLPAR dependencies, the following commands can be used:

- The **ps** command displays bindprocessor attachments and **block** system call status at the process level.
- The **bindprocessor** command displays online processors and makes new attachments.
- The **kill** command sends signals to processes.
- The **ipcs** command displays pinned shared-memory segments at the process level.

- The **lsrset** command displays processor sets.
- The **lsclass** command displays Workload Manager classes, which might include processor sets.
- The **chclass** command is used to change class definitions.

Scripts can also be used for scalability and general performance issues. When resources are removed, you can reduce the number of threads that are used or the size of application buffers. When the resources are added, you can increase these parameters. You can provide commands that can be used to dynamically make these adjustments, which can be triggered by these scripts. Install the scripts to start these commands within the context of the DLPAR operation.

High-level structure of DLPAR scripts

This section provides an overview of the scripts, which can be Perl scripts, shell scripts, or commands. Application scripts are required to provide the following commands:

- **scriptinfo**

Identifies the version, date, and vendor of the script. It is called when the script is installed.

- **register**

Identifies the resources managed by the script. If the script returns the resource name *cpu*, *mem*, *capacity*, or *var_weight*, the script is automatically started when DLPAR attempts to reconfigure processors, memory, entitled capacity, or variable weight. The **register** command is called when the script is installed with the DLPAR subsystem.

- **usage** *resource_name*

Returns information describing how the resource is being used by the application. The description should be relevant so that the user can determine whether to install or uninstall the script. It should identify the software capabilities of the application that are impacted. The **usage** command is called for each resource that was identified by the **register** command.

- **checkrelease** *resource_name*

Indicates whether the DLPAR subsystem should continue with the removal of the named resource. A script might indicate that the resource should not be removed if the application is not DLPAR-aware and the application is considered critical to the operation of the system.

- **prerelease** *resource_name*

Reconfigures, suspends, or terminates the application so that its hold on the named resource is released.

- **postrelease** *resource_name*

Resumes or restarts the application.

- **undoprerelease** *resource_name*

Started if an error is encountered and the resource is not released.

- **checkacquire** *resource_name*

Indicates whether the DLPAR subsystem should proceed with the resource addition. It might be used by a license manager to prevent the addition of a new resource, for example *cpu*, until the resource is licensed.

- **preacquire** *resource_name*

Used to prepare for a resource addition.

- **undopreacquire** *resource_name*

Starts if an error is encountered in the preacquire phase or when the event is acted on.

- **postacquire** *resource_name*

Resumes or starts the application.

- **preacvent** *resource_name*

Used to prepare a DLPAR update.

- **postaccevent** *resource_name*

Resumes or starts the application.

- **undopreaccevent** *resource_name*

Started if an error is encountered in the **preaccevent** phase or when the event is acted upon.

- **pretopolgyupdate** *resource_name*

Used to prepare for a system topology update.

- **postopolgyupdate** *resource_name*

Resumes or starts the application.

Installing application scripts using the **drmgr** command

The **drmgr** command maintains an internal database of installed-script information. This information is collected when the system is booted and is refreshed when new scripts are installed or uninstalled. The information is derived from the **scriptinfo**, **register**, and **usage** commands. The installation of scripts is supported through in the **drmgr** command, which copies the named script to the script repository where it can be later accessed. The default location for this repository is **/usr/lib/dr/scripts/all**. Within workload partitions, the default script repository location is **/var/dr/scripts**. You can specify an alternate location for this repository. To determine the machine upon which a script is used, specify the target host name when installing the script.

To specify the location of the base repository, use the following command:

```
drmgr -R base_directory_path
```

To install a script, use the following command:

```
drmgr -i script_name [-f] [-w mins] [-D hostname]
```

The following flags are defined:

- The **-i** flag is used to name the script.
- The **-f** flag must be used to replace a registered script.
- The **-w** flag is used to specify the number of minutes that the script is expected to execute. This is provided as an override option to the value specified by the vendor.
- The **-D** flag is used to register a script to be used on a particular host.

To uninstall a script, use the following command:

```
drmgr -u script_name [-D hostname]
```

The following flags are defined:

- The **-u** flag is used to indicate which script should be uninstalled.
- The **-D** flag is used to uninstall a script that has been registered for a specific directory.

To display information about scripts that have already been installed, use the following command:

```
drmgr -l
```

Naming conventions for scripts

It is suggested that the script names be built from the vendor name and the subsystem that is being controlled. System administrators should name their scripts with the *sysadmin* prefix. For example, a system administrator who wanted to provide a script to control Workload Manager assignments might name the script *sysadmin_wlm*.

Script execution environment and input parameters

Scripts are started with the following execution environment:

- Process UID is set to the UID of the script.
- Process GID is set to the GID of the script.
- **PATH** environment variable is set to the **/usr/bin:/etc:/usr/sbin** directory.
- **LANG** environment variable might or might not be set.
- Current working directory is set to **/tmp**.
- Command arguments and environment variables are used to describe the DLPAR event.

Scripts receive input parameters through command arguments and environment variables, and provide output by writing *name=value* pairs to standard output, where *name=value* pairs are delimited by new lines. The *name* is defined to be the name of the return data item that is expected, and *value* is the value associated with the data item. Text strings must be enclosed by parentheses; for example, `DR_ERROR="(text)"`. All environment variables and *name=value* pairs must begin with `DR_`, which is reserved for communicating with application scripts.

The scripts use **DR_ERROR** *name=value* environment variable pair to provide error descriptions.

You can examine the command arguments to the script to determine the phase of the DLPAR operation, the type of action, and the type of resource that is the subject of the pending DLPAR request. For example, if the script command arguments are `checkrelease mem`, then the phase is `check`, the action is `remove`, and the type of resource is `memory`. The specific resource that is involved can be identified by examining environment variables.

The following environment variables are set for memory add and remove:

Note: In the following description, one frame is equal to 4 KB.

- **DR_FREE_FRAMES=0xFFFFFFFF**

The number of free frames currently in the system, in hexadecimal format.

- **DR_MEM_SIZE_COMPLETED=*n***

The number of megabytes that were successfully added or removed, in decimal format.

- **DR_MEM_SIZE_REQUEST=*n***

The size of the memory request in megabytes, in decimal format.

- **DR_PINNABLE_FRAMES=0xFFFFFFFF**

The total number of pinnable frames currently in the system, in hexadecimal format. This parameter provides valuable information when removing memory in that it can be used to determine when the system is approaching the limit of pinnable memory, which is the primary cause of failure for memory remove requests.

- **DR_TOTAL_FRAMES=0xFFFFFFFF**

The total number of frames currently in the system, in hexadecimal format.

The following environment variables are set for processor add and remove:

- **DR_BCPUID=*N***

The bind CPU ID of the processor that is being added or removed in decimal format. A **bindprocessor** attachment to this processor does not necessarily mean that the attachment has to be undone. This is only true if it is the *N*th processor in the system, because the *N*th processor position is the one that is always removed in a Central Processing Unit (CPU) remove operation. Bind IDs are consecutive in nature, ranging from 0 to *N* and are intended to identify only online processors. Use the **bindprocessor** command to determine the number of online CPUs.

- **DR_LCPUID=*N***

The logical Central Processor Unit (CPU) ID of the processor that is being added or removed in decimal format.

The following environment variables are set for Micro-Partitioning.

DR_CPU_CAPACITY=N

The partition's percentage of shared physical processors.

DR_VAR_WEIGHT=N

The partition's relative priority for determining how to allocate shared pool idle cycles.

DR_CPU_CAPACITY_DELTA=N

The difference between the current value of the partition's percentage of shared physical processors and the value to which it will be changed when this operation is complete.

DR_VAR_WEIGHT_DELTA=N

The difference between the current value of the partition's variable weight and the value to which it will be changed when this operation is complete.

The operator can display the information about the current DLPAR request using the detail level at the HMC to observe events as they occur. This parameter is specified to the script using the **DR_DETAIL_LEVEL=N** environment variable, where *N* can range from 0 to 5. The default value is zero (0) and signifies no information. A value of one (1) is reserved for the operating system and is used to present the high-level flow. The remaining levels (2-5) can be used by the scripts to provide information with the assumption that larger numbers provide greater detail.

Scripts provide detailed data by writing the following *name=value* pairs to standard output:

name=value pair	Description
DR_LOG_ERR = <i>message</i>	Logs the message with the syslog level of the LOG_ERR environment variable.
DR_LOG_WARNING = <i>message</i>	Logs the message with the syslog level of the LOG_WARNING environment variable.
DR_LOG_INFO = <i>message</i>	Logs the message with the syslog level of the LOG_INFO environment variable.
DR_LOG_EMERG = <i>message</i>	Logs the message with the syslog level of the LOG_EMERG environment variable.
DR_LOG_DEBUG = <i>message</i>	Logs the message with the syslog level of the LOG_DEBUG environment variable.

In addition, the operator can also set up a log of information that is preserved by using the **syslog** facility, in which case, the above information is routed to that facility as well. You must configure the **syslog** facility in this case.

DLPAR script commands

This section describes the script commands for DLPAR:

scriptinfo

Provides information about the installed scripts, such as their creation date and resources.

register

started to collect the list of resources that are managed by the script. The **drmgr** command uses these lists to start scripts based on the type of resource that is being reconfigured.

usage

Provides human-readable strings describing the service provided by the named resource. The context of the message should help the user decide the implications on the application and the services that it provides when named resource is reconfigured. This command is started when the script is installed, and the information provided by this command is maintained in an internal database that is used by the **drmgr** command. Display the information using the *-l* list option of the **drmgr** command.

checkrelease

When removing resources, the **drmgr** command assesses the impacts of the removal of the resource. This includes execution of DLPAR scripts that implement the **checkrelease** command. Each DLPAR script in turn will be able to evaluate the peculiarities of its application and indicate to the **drmgr**

command that is using the script's return code whether the resource removal will affect the associated application. If it finds that the removal of the resource can be done safely, an exit status of success is returned. If the application is in a state that the resource is critical to its execution and cannot be reconfigured without interrupting the execution of the application, then the script indicates that the resource should not be removed by returning an error. When the *FORCE* option is specified by the user, which applies to the entire DLPAR operation including its phases, the **drmgr** command skips the **checkrelease** command and begins with the **prerelease** commands.

prerelease

Before a resource is released, the DLPAR scripts are directed to assist in the release of the named resource by reducing or eliminating the use of the resource from the application. However, if the script detects that the resource cannot be released from the application, it should indicate that the resource will not be removed from the application by returning an error. This does not prevent the system from attempting to remove the resource in either the forced or non-forced mode of execution, and the script will be called in the post phase, regardless of the actions or inactions that were taken by the **prerelease** command. The actions taken by the operating system are safe. If a resource cannot be cleanly removed, the operation will fail.

The DLPAR script is expected to internally record the actions that were taken by the **prerelease** command, so that they can be restored in the post phase should an error occur. This can also be managed in post phase if rediscovery is implemented. The application might need to take severe measures if the *force* option is specified.

postrelease

After a resource is successfully released, the **postrelease** command for each installed DLPAR script is started. Each DLPAR script performs any post processing that is required during this step. Applications that were halted should be restarted.

The calling program will ignore any errors reported by the **postrelease** commands, and the operation will be considered a success, although an indication of any errors that may have occurred will also be reported to the user. The **DR_ERROR** environment variable message is provided for this purpose, so the message should identify the application that was not properly reconfigured.

undoprerelease

After a **prerelease** command is issued by the **drmgr** command to the DLPAR script, if the **drmgr** command fails to remove or release the resource, it will try to revert to the old state. As part of this process, the **drmgr** command will issue the **undoprerelease** command to the DLPAR script. The **undoprerelease** command will only be started if the script was previously called to release the resource in the current DLPAR request. In this case, the script should undo any actions that were taken by the **prerelease** command of the script. To this end, the script might need to document its actions, or otherwise provide the capability of rediscovering the state of the system and reconfiguring the application, so that in effect, the DLPAR event never occurred.

checkacquire

This command is the first DLPAR script-based command that is called in the acquire-new-resource sequence. It is called for each installed script that previously indicated that it supported the particular type of resource that is being added. One of the primary purposes of the checkacquire phase is to enable processor-based license managers, which might want to fail the addition of a processor. The **checkacquire** command is always started, regardless of the value of the **FORCE** environment variable, and the calling program honors the return code of the script. The user cannot force the addition of a new processor if a script or DLPAR-aware program fails the DLPAR operation in the check phase.

In short, the **FORCE** environment variable does not really apply to the **checkacquire** command, although it does apply to the other phases. In the preacquire phase, it dictates how far the script should go when reconfiguring the application. The *force* option can be used by the scripts to control the policy by which applications are stopped and restarted similar to when a resource is released, which is mostly a DLPAR-safe issue.

preacquire

Assuming that no errors were reported in the checkacquire phase, the system advances to the preacquire phase, where the same set of scripts are invoked to prepare for the acquisition of a new resource, which is supported through the **preacquire** command. Each of these scripts are called,

before the system actually attempts to integrate the resource, unless an error was reported and the **FORCE** environment variable was not specified by the user. If the **FORCE** environment variable was specified, the system proceeds to the integrate stage, regardless of the script's stated return code. No errors are detected when the **FORCE** environment variable is specified, because all errors are avoidable by unconfiguring the application, which is an accepted practice when the **FORCE** environment variable is specified. If an error is encountered and the **FORCE** environment variable is not specified, the system will proceed to the undopreacquire phase, but only the previously executed scripts in the current phase are rerun. During this latter phase, the scripts are directed to perform recovery actions.

undopreacquire

The undopreacquire phase is provided so that the scripts can perform recovery operations. If a script is called in the undopreacquire phase, it can assume that it successfully completed the **preacquire** command.

postacquire

The **postacquire** command is executed after the resource has been successfully integrated by the system. Each DLPAR script that was previously called in the check and pre phases is called again. This command is used to incorporate the new resource into the application. For example, the application might want to create new threads, expands its buffers, or the application might need to be restarted if it was previously halted.

checkmigrate

This command is the first DLPAR script-based command that is called in the migration sequence. It is called for each installed script that previously indicated that it supported the particular type of resource that is being added. The **checkmigrate** command is always started, regardless of the value of the **FORCE** environment variable, and the calling program honors the return code of the script. The user cannot force the partition migration if a script or DLPAR-aware program fails the DLPAR operation in the check phase.

premigrate

Assuming that no errors were reported in the **checkmigrate** phase, the system advances to the **premigrate** phase, where the same set of scripts are started to prepare for partition takes place. Each of these scripts are called, before the system actually attempts to migrate the partition. Regardless of the script's stated return code, the system proceeds to the migration stage. If an error is encountered, the system will proceed to the **undopremigrate** phase, but only the previously executed scripts in the current phase are rerun. During this latter phase, the scripts are directed to perform recovery actions.

undopremigrate

The **undopremigrate** phase is provided so that the scripts can perform recovery operations. If a script is called in the **undopremigrate** phase, it can assume that it successfully completed the **premigrate** command.

postmigrate

The **postmigrate** command is executed after the partition has been successfully migrated. Each DLPAR script that was previously called in the check and pre phases is called again.

pretopologyupdate

The **pretopologyupdate** command is executed before an action takes place that will affect the topology of the partition, such as the addition or removal of processors or memory. This command is meant to inform the scripts a topology action has taken place and can not fail. The system proceeds to the integrate stage, regardless of the script's stated return code.

posttopologyupdate

The **posttopologyupdate** command is executed after the partition has been successfully completed the topology action. Each DLPAR script that was previously called in the pre phase is called again.

checkhibernate

This command is the first DLPAR script-based command that is called in the hibernation sequence. It is called for each installed script that previously indicated that it supported the particular type of resource that is being added. The **checkhibernate** command is always started, regardless of the value of the **FORCE** environment variable, and the calling program honors the return code of the script. The user cannot force the partition hibernation if a script or DLPAR-aware program fails the DLPAR operation in the check phase.

prehibernate

Assuming that no errors were reported in the **checkhibernate** phase, the system advances to the **prehibernate** phase, where the same set of scripts are started to prepare for partition takes place. Each of these scripts are called, before the system actually attempts to hibernate the partition. Regardless of the script's stated return code, the system proceeds to the hibernation stage. If an error is encountered, the system will proceed to the **undohibernate** phase, but only the previously executed scripts in the current phase are rerun. During this latter phase, the scripts are directed to perform recovery actions.

undohibernate

The **undohibernate** phase is provided so that the scripts can perform recovery operations. If a script is called in the **checkhibernate** phase, it can assume that it successfully completed the **checkhibernate** command.

posthibernate

The **posthibernate** command is executed after the partition has been successfully hibernated. Each DLPAR script that was previously called in the check and pre phases is called again.

preaccevent

This command is the first DLPAR script-based command that is called in the encryption accelerator DLPAR sequence. It is called for each installed script that previously indicated that it supported the particular type of resource that is being added or released. It is unknown at the time of this event if the following action will be an add or release of the encryption accelerator. That action will be provided during one of the following post phases.

postaccevent

The **postaccevent** command is executed after the resource has been successfully processed by the system. Each DLPAR script that was previously called in the pre phase is called again. This command is used to incorporate the new resource state into the application.

undoaccevent

The **undoaccevent** phase is provided so that the scripts can perform recovery operations. If a script is called in the **undoaccevent** phase, it successfully completed the **preaccevent** command.

Making kernel extensions DLPAR-aware

Like applications, most kernel extensions are DLPAR-safe by default.

However, some are sensitive to the system configuration and might need to be registered with the DLPAR subsystem. Some kernel extensions partition their data along processor lines, create threads based on the number of online processors, or provide large pinned memory buffer pools. These kernel extensions must be notified when the system topology changes. The mechanism and the actions that need to be taken parallel those of DLPAR-aware applications.

Registering reconfiguration handlers

The following kernel services are provided to register and unregister reconfiguration handlers:

```
#include sys/dr.h

int reconfig_register(int (*handler)(void *, void *, int, dr_info_t *),
                    int actions, void * h_arg, ulong *h_token, char *name);

void reconfig_unregister(ulong h_token);

int (*handler)(void *event, void *h_arg, unsigned long long req, void *resource_info);
void reconfig_unregister(ulong h_token);

int reconfig_register_ext (int (*handler)(void *, void *, unsigned long long, dr_info_t *),
                          unsigned long long actions, void * h_arg, ulong *h_token, char *name);
int (*handler)(void *event, void *h_arg, unsigned long long req, void *resource_info);

kernno_t reconfig_register_list(int (*handler)(void *, void *, dr_kevent_t, void *),
                               dr_kevent_t event_list[], size_t list_size, void *h_arg, ulong *h_token, char *name);
```

```
int (*handler)(void *event, void *h_arg, dr_kevent_t event_in_prog, void *resource_info);
```

Note: You are encouraged to use the kernel service **reconfig_register_list**. This service supports more events to be notified to the kernel extensions. The previous kernel services (**reconfig_register** and **reconfig_register_ext**) are limited to 32 and 64 events, respectively, making the kernel extensions using this service not portable to future systems supporting more than 32 and 64 events.

The parameters for the **reconfig_register**, **reconfig_register_ext**, and **reconfig_register_list** subroutines are as follows:

- The *handler* parameter is the kernel extension function to be invoked.
- The *actions* parameter allows the kernel extension to specify which events require notification. For a list of the events, see the **reconfig_register**, **reconfig_register_ext**, and **reconfig_unregister** kernel services.
- The *h_arg* parameter is specified by the kernel extension, remembered by the kernel along with the function descriptor for the handler, and then passed to the handler when it is invoked. It is not used directly by the kernel, but is intended to support kernel extensions that manage multiple adapter instances. In practice, this parameter points to an adapter control block.
- The *h_token* parameter is an output parameter and is intended to be used when the handler is unregistered.
- The *name* parameter is provided for information purposes and can be included within an error log entry if the driver returns an error. It is provided by the kernel extension and should be limited to 15 ASCII characters.
- The *event_list* parameter is an array of `dr_kevent_t` values for which the kernel extension should be notified of when they occur. For a list of defined events, see the **reconfig_register_list** kernel service.
- The *list_size* parameter is the size of the memory consumed by the *event_list* parameter.

The **reconfig_register** and **reconfig_register_ext** functions return 0 for success, and the appropriate `errno` value otherwise.

The **reconfig_unregister** function is called to remove a previously installed handler.

The **reconfig_register**, **reconfig_register_ext**, and **reconfig_unregister** functions can only be called in the process environment.

If a kernel extension registers for the pre-phase, it is advisable that it register for the check phase to avoid partial unconfiguration of the system when removing resources.

Reconfiguration Handlers

The interface to the reconfiguration handler used with the **reconfig_register_list** kernel service is as follows:

```
int (*handler)(void *event, void *h_arg, dr_kevent_t event_in_prog, void *resource_info);
```

The parameters to the reconfiguration handler are as follows:

- The *event* parameter is passed to the handler and is intended to be used only when calling the **reconfig_handler_complete** subroutine.
- The *h_arg* parameter is specified at registration time by the handler.
- The *event_in_prog* parameter indicates the DLPAR operation performed by the handler. For a list of the events, see the **reconfig_register_list** kernel service.
- The *resource_info* parameter identifies the resource-specific information for the current DLPAR request. If the request is processor-based, then the *resource_info* data is provided through a *dri_cpu* structure. If the request is memory-based, a *dri_mem* structure is used. On a Micro-Partitioning partition, if the request is processor-capacity based, the *resource_info* data is provided through a *dri_cpu_capacity*

structure. For more information, and for the format of the `dri_cpu_capacity` structure, refer to [reconfig Kernel Service](#).

```
struct dri_cpu {
    cpu_t      lcpu;          /* Logical CPU Id of target CPU */
    cpu_t      bcpu;          /* Bind Id of target CPU */
};

struct dri_mem {
    size64_t   req_memsz_change; /* user requested mem size */
    size64_t   sys_memsz;        /* system mem size at start */
    size64_t   act_memsz_change; /* mem added/removed so far */
    rpn64_t    sys_free_frames; /* Number of free frames */
    rpn64_t    sys_pinnable_frames; /* Number of pinnable frames */
    rpn64_t    sys_total_frames; /* Total number of frames */
    unsigned long long lmb_addr; /* start addr of logical memory block */
    size64_t   lmb_size;        /* Size of logical memory block being added */
};
```

If the current DLPAR request is migration of a partition, the handler provides `resource_info` data to the `resource_info` data kernel extensions, but the kernel extensions do not need to access the contents of the `resource_info` data because this data is not used by the kernel extensions.

Reconfiguration handlers are invoked in the process environment.

Kernel extensions can assume the following:

- Only a single type of resource is being configured or removed at a time.
- Multiple processors will not be specified at the same time. However, kernel extensions should be coded to support the addition or removal of multiple logical memory blocks. You can initiate a request to add or remove gigabytes of memory.

The check phase provides the ability for DLPAR-aware applications and kernel extensions to react to the user's request before it has been applied. Therefore, the check-phase kernel extension handler is invoked once, even though the request might devolve to multiple logical memory blocks. Unlike the check phase, the pre-phase, post phase, and post-error phase are applied at the logical memory block level. This is different for application notification, where the pre-phase, post phase, or alternatively the post-error phase are invoked once for each user request, regardless of the number of underlying logical memory blocks. Another difference is that the post-error phase for kernel extensions is used when a specific logical memory block operation fails, whereas the post-error phase for applications is used when the operation, which in this case is the entire user request, fails.

In general, during the check phase, the kernel extension examines its state to determine whether it can comply with the impending DLPAR request. If this operation cannot be managed, or if it would adversely effect the proper execution of the extension, then the handler returns **DR_FAIL**. Otherwise the handler returns **DR_SUCCESS**.

During the pre-remove phase, kernel extensions attempt to remove any dependencies that they might have on the designated resource. An example is a driver that maintains per-processor buffer pools. The driver might mark the associated buffer pool as pending delete, so that new requests are not allocated from it. In time, the pool is drained and might be freed. Other items that must be considered in the pre-remove phase are timers and bound threads, which need to be stopped and terminated, respectively. Alternatively, bound threads can be unbound.

During the post-remove phase, kernel extensions attempt to free any resources through garbage collection, assuming that the resource was actually removed. If it was not, timers and threads must be reestablished. The **DR_resource_POST_ERROR** request is used to signify that an error occurred.

During the pre-add phase, kernel extensions must pre-initialize any data paths that are dependent on the new resource, so that when the new resource is configured, it is ready to be used. The system does not guarantee that the resource will not be used prior to the handler being called again in the post phase.

During the post-add phase, kernel extensions can assume that the resource has been properly added and can be used. This phase is a convenient place to start bound threads, schedule timers, and increase the size of buffers.

Kernel extensions can also be notified of memory removals or additions on a per-operation basis, much like applications, by registering for one or more of the **_OP_** notification types. This enables a kernel

extension to modify its resource usage in response to a memory DR operation only once per operation, rather than once per logical memory block (LMB).

DR_MEM_REMOVE_OP_PRE notification is sent before a memory remove. Reconfiguration handlers can start adjusting their resources in anticipation of the memory remove at this time.

DR_MEM_REMOVE_OP_POST and **DR_MEM_ADD_OP_POST** notifications are sent after a memory remove or add operation, respectively, whether the operation failed or not. If the operation failed, **act_memsz_change** is 0.

If possible, within a few seconds, the reconfiguration handlers return **DR_SUCCESS** to indicate successful reconfiguration, or **DR_FAIL** to indicate failure. If more time is required, the handler returns **DR_WAIT**.

Extended DR handlers

If a kernel extension expects that the operation is likely to take a long time, that is, several seconds, the handler returns **DR_WAIT** to the caller, but proceed with the request asynchronously. In the following case, the handler indicates that it has completed the request by invoking the **reconfig_handler_complete** routine.

```
void reconfig_handler_complete(void *event, int rc);
```

The *event* parameter is the same parameter that was passed to the handler when it was invoked by the kernel. The *rc* parameter must be set to either **DR_SUCCESS** or **DR_FAIL** to indicate the completion status of the handler.

The **reconfig_handler_complete** kernel service can be invoked in the process or interrupt environments.

Using the **xmemdma** kernel service

On systems that are capable of DLPAR, such as the dynamic removal of memory, calls to the **xmemdma** kernel service without the **XMEM_DR_SAFE** flag result in the specified memory being flagged as not removable. This is done to guarantee the integrity of the system, because the system has no knowledge of how the caller intends to use the real memory address that was returned. Dynamic memory removal operations are still possible for other memory, but not for the memory that the **xmemdma** call specifies.

If the caller is using the real memory address only for informational purposes, such as for trace buffers or debug information, then the caller can set the **XMEM_DR_SAFE** flag. This is an indication to the system that the real memory address can be exposed to the caller without any risk of data corruption. When this flag is present, the system will still permit the specified memory to be dynamically removed.

If the caller is using the real memory address to perform actual data access, either by turning off data translation and performing CPU load or store access to the real memory, or by programming direct memory access (DMA) controllers to target the real memory, the **XMEM_DR_SAFE** flag must not be set. If the flag is set, the system's data integrity could be compromised when the memory is dynamically removed. For information on converting a kernel extension that uses real memory addresses in this way to be DLPAR-aware, contact your IBM Service Representative.

For more information, see the **xmemdma** kernel service.

Controlling memory DLPAR notification for applications

Dynamic addition or removal of memory from an LPAR running multiple DLPAR-aware programs can result in conflict for resources. By default, each program gets notified equally about the resource change. For example, if 1 GB of memory is removed from an LPAR running two DR-aware programs, then, by default, each program is notified that 1 GB of memory has been removed. Because the two programs are generally unaware of each other, both of them will scale down their memory use by 1 GB, leading to inefficiency. A similar efficiency problem can also occur when new memory is added.

To overcome this problem, AIX allows application scripts to be installed with a percentage factor that indicates the percentage of the actual memory resource change. The system then notifies the application in the event of a memory DLPAR. While installing the application scripts using the **drmgr** command, you can specify this percentage factor using the **DR_MEM_PERCENT name=value** pair. The application script will need to output this name=value pair when it is invoked by the **drmgr** command with the **scriptinfo** subcommand. The value must be an integer between 1 and 100. Any value outside of this range is

ignored, and the default value, which is 100, is used. Additionally, you can also set this **name=value** pair as an environment variable at the time of installation. During installation, the value from the environment variable, if set, will override the value provided by the application script.

Similarly, in applications using the **SIGRECONFIG** signal handler and **dr_reconfig()** system call, you can control the memory DLPAR notification by setting the **DR_MEM_PERCENT name=value** pair as an environment variable before the application begins running. This value, however, cannot be changed without restarting the application.

sed program information

The **sed** program is a text editor that has similar functions to those of **ed**, the line editor.

Unlike **ed**, however, the **sed** program performs its editing without interacting with the person requesting the editing.

Manipulating strings with sed

The **sed** program performs its editing without interacting with the person requesting the editing.

This method of operation allows **sed** to do the following:

- Edit very large files
- Perform complex editing operations many times without requiring extensive retyping and cursor positioning (as interactive editors do)
- Perform global changes in one pass through the input.

The editor keeps only a few lines of the file being edited in memory at one time, and does not use temporary files. Therefore, the file to be edited can be any size as long as there is room for both the input file and the output file in the file system.

Related concepts

[Tools and utilities](#)

This section provides an overview of the tools and utilities that you can use to develop C compiled language programs.

Starting the editor

The following example description file could maintain the **make** program.

Each command in the command file must be on a separate line. Once the command file is created, enter the following command on the command line:

```
sed -fCommandFile >Output <Input
```

In this command the parameters mean the following:

Parameter	Definition
<i>CommandFile</i>	The name of the file containing editing commands.
<i>Output</i>	The name of the file to contain the edited output.
<i>Input</i>	The name of the file, or files, to be edited.

The **sed** program then makes the changes and writes the changed information to the output file. The contents of the input file are not changed.

How sed works

The **sed** program is a stream editor that receives its input from standard input, changes that input as directed by commands in a command file, and writes the resulting stream to standard output.

If you do not provide a command file and do not use any flags with the **sed** command, the **sed** program copies standard input to standard output without change. Input to the program comes from two sources:

Program	Description
Input stream	A stream of ASCII characters either from one or more files or entered directly from the keyboard. This stream is the data to be edited.
Commands	A set of addresses and associated commands to be performed, in the following general form:

```
[Line1 [,Line2] ] command [argument]
```

The parameters *Line1* and *Line2* are called addresses. Addresses can be either patterns to match in the input stream, or line numbers in the input stream.

You can also enter editing commands along with the **sed** command by using the **-e** flag.

When **sed** edits, it reads the input stream one line at a time into an area in memory called the pattern space. When a line of data is in the pattern space, **sed** reads the command file and tries to match the addresses in the command file with characters in the pattern space. If it finds an address that matches something in the pattern space, **sed** then performs the command associated with that address on the part of the pattern space that matched the address. The result of that command changes the contents of the pattern space, and thus becomes the input for all following commands.

When **sed** has tried to match all addresses in the command file with the contents of the pattern space, it writes the final contents of the pattern space to standard output. Then it reads a new input line from standard input and starts the process over at the start of the command file.

Some editing commands change the way the process operates.

Flags used with the **sed** command can also change the operation of the command.

Using regular expressions

A regular expression is a string that contains literal characters, pattern-matching characters and/or operators that define a set of one or more possible strings.

The stream editor uses a set of pattern-matching characters that is different from the shell pattern-matching characters, but the same as the line editor, **ed**.

Using the sed command summary

All **sed** commands are single letters plus some parameters, such as line numbers or text strings.

The commands summarized below make changes to the lines in the pattern space.

The following symbols are used in the syntax diagrams:

Symbol	Meaning
[]	Square brackets enclose optional parts of the commands
<i>italics</i>	Parameters in italics represent general names for a name that you enter. For example, <i>FileName</i> represents a parameter that you replace with the name of an actual file.
<i>Line1</i>	This symbol is a line number or regular expression to match that defines the starting point for applying the editing command.
<i>Line2</i>	This symbol is a line number or regular expression to match that defines the ending point to stop applying the editing command.

Line manipulation

This section describes the line manipulation.

Function	Syntax/Description
append lines	<p>[<i>Line1</i>]a\\n<i>Text</i></p> <p>Writes the lines contained in <i>Text</i> to the output stream after <i>Line1</i>. The a command must appear at the end of a line.</p>
change lines	<p>[<i>Line1</i> [,<i>Line2</i>]]c\\n<i>Text</i></p> <p>Deletes the lines specified by <i>Line1</i> and <i>Line2</i> as the <i>delete lines</i> command does. Then it writes <i>Text</i> to the output stream in place of the deleted lines.</p>
delete lines	<p>[<i>Line1</i> [,<i>Line2</i>]]d</p> <p>Removes lines from the input stream and does not copy them to the output stream. The lines not copied begin at line number <i>Line1</i>. The next line copied to the output stream is line number <i>Line2</i> + 1. If you specify only one line number, then only that line is not copied. If you do not specify a line number, the next line is not copied. You cannot perform any other functions on lines that are not copied to the output.</p>
insert lines	<p>[<i>Line1</i>] i \\n<i>Text</i></p> <p>Writes the lines contained in <i>Text</i> to the output stream before <i>Line1</i>. The i command must appear at the end of a line.</p>
next line	<p>[<i>Line1</i> [,<i>Line2</i>]]n</p> <p>Reads the next line, or group of lines from <i>Line1</i> to <i>Line2</i> into the pattern space. The current contents of the pattern space are written to the output if it has not been deleted.</p>

Substitution

This section describes the substitution.

Function	Syntax/Description
substitution for pattern	<p>[<i>Line1</i> [,<i>Line2</i>]] s/<i>Pattern/String/Flags</i></p> <p>Searches the indicated line(s) for a set of characters that matches the regular expression defined in <i>Pattern</i>. When it finds a match, the command replaces that set of characters with the set of characters specified by <i>String</i>.</p>

Input and output

This section describes the input output.

Function	Syntax/Description
print lines	<p>[<i>Line1</i> [,<i>Line2</i>]] p</p> <p>Writes the indicated lines to STDOUT at the point in the editing process that the p command occurs.</p>

Function	Syntax/Description
write lines	<p><i>[Line1 [,Line2]]w FileName</i></p> <p>Writes the indicated lines to a <i>FileName</i> at the point in the editing process that the w command occurs.</p> <p>If <i>FileName</i> exists, it is overwritten; otherwise, it is created. A maximum of 10 different files can be mentioned as input or output files in the entire editing process. Include exactly one space between w and <i>FileName</i>.</p>
read file	<p><i>[Line1]r FileName</i></p> <p>Reads <i>FileName</i> and appends the contents after the line indicated by <i>Line1</i>.</p> <p>Include exactly one space between r and <i>FileName</i>. If <i>FileName</i> cannot be opened, the command reads it as a null file without giving any indication of an error.</p>

Matching across lines

This section describes the matching across lines.

Function	Syntax/Description
join next line	<p><i>[Line1 [,Line2]]N</i></p> <p>Joins the indicated input lines together, separating them by an embedded new-line character. Pattern matches can extend across the embedded new-lines(s).</p>
delete first line of pattern space	<p><i>[Line1 [,Line2]]D</i></p> <p>Deletes all text in the pattern space up to and including the first new-line character. If only one line is in the pattern space, it reads another line. Starts the list of editing commands again from the beginning.</p>
print first line of pattern space	<p><i>[Line1 [,Line2]]P</i></p> <p>Prints all text in the pattern space up to and including the first new-line character to STDOUT.</p>

Pick up and put down

This section describes the pick up and put down.

Function	Syntax/Description
pick up copy	<p><i>[Line1 [,Line2]]h</i></p> <p>Copies the contents of the pattern space indicated by <i>Line1</i> and <i>Line2</i> if present, to the holding area.</p>
pick up copy, appended	<p><i>[Line1 [,Line2]]H</i></p> <p>Copies the contents of the pattern space indicated by <i>Line1</i> and <i>Line2</i> if present, to the holding area, and appends it to the end of the previous contents of the holding area.</p>

Function	Syntax/Description
put down copy	<i>[Line1 [,Line2]]g</i> Copies the contents of the holding area to the pattern space indicated by <i>Line1</i> and <i>Line2</i> if present. The previous contents of the pattern space are destroyed.
put down copy, appended	<i>[Line1 [,Line2]]G</i> Copies the contents of the holding area to the end of the pattern space indicated by <i>Line1</i> and <i>Line2</i> if present. The previous contents of the pattern space are not changed. A new-line character separates the previous contents from the appended text.
exchange copies	<i>[Line1 [,Line2]]x</i> Exchanges the contents of the holding area with the contents of the pattern space indicated by <i>Line1</i> and <i>Line2</i> if present.

Control

This section describes the pick up and put down.

Function	Syntax/Description
negation	<i>[Line1 [,Line2]]!</i> The ! (exclamation point) applies the command that follows it on the same line to the parts of the input file that are <i>not</i> selected by <i>Line1</i> and <i>Line2</i> .
command groups	<i>[Line1 [,Line2]]{</i> <i>grouped commands</i> <i>}</i> The { (left brace) and the } (right brace) enclose a set of commands to be applied as a set to the input lines selected by <i>Line1</i> and <i>Line2</i> . The first command in the set can be on the same line or on the line following the left brace. The right brace must be on a line by itself. You can nest groups within groups.
labels	<i>:Label</i> Marks a place in the stream of editing command to be used as a destination of each branch. The symbol <i>Label</i> is a string of up to 8 bytes. Each <i>Label</i> in the editing stream must be different from any other <i>Label</i> .

Function	Syntax/Description
branch to label, unconditional	<i>[Line1 [,Line2]]bLabel</i> Branches to the point in the editing stream indicated by <i>Label</i> and continues processing the current input line with the commands following <i>Label</i> . If <i>Label</i> is null, branches to the end of the editing stream, which results in reading a new input line and starting the editing stream over. The string <i>Label</i> must appear as a <i>Label</i> in the editing stream.
test and branch	<i>[Line1 [,Line2]]tLabel</i> If any successful substitutions were made on the current input line, branches to <i>Label</i> . If no substitutions were made, does nothing. Clears the flag that indicates a substitution was made. This flag is cleared at the start of each new input line.
wait	<i>[Line1]q</i> Stops editing in an orderly fashion by writing the current line to the output, writing any appended or read test to the output, and stopping the editor.
find line number	<i>[Line1]=</i> Writes to standard output the line number of the line that matches <i>Line1</i> .

Using text in commands

The **append**, **insert** and **change** lines commands all use a supplied text string to add to the output stream.

This text string conforms to the following rules:

- Can be one or more lines long.
- Each `\n` (new-line character) inside *Text* must have an additional `\` character before it (`\\n`).
- The *Text* string ends with a new-line that does not have an additional `\` character before it (`\n`).
- Once the command inserts the *Text* string, the string:
 - Is always written to the output stream, regardless of what other commands do to the line that caused it to be inserted.
 - Is not scanned for address matches.
 - Is not affected by other editing commands.
 - Does not affect the line number counter.

Using string replacement

The **s** command performs string replacement in the indicated lines in the input file.

If the command finds a set of characters in the input file that satisfies the regular expression *Pattern*, it replaces the set of characters with the set of characters specified in *String*.

The *String* parameter is a literal set of characters (digits, letters and symbols). Two special symbols can be used in *String*:

Symbol	Use
&	This symbol in <i>String</i> is replaced by the set of characters in the input lines that matched <i>Pattern</i> . For example, the command:

```
s/boy/&s/
```

tells **sed** to find a pattern boy in the input line, and copy that pattern to the output with an appended s. Therefore, it changes the input line:

From:

The boy look at the game.

To:

The boys look at the game.

Symbol	Use
\d	d is a single digit. This symbol in <i>String</i> is replaced by the set of characters in the input lines that matches the dth substring in <i>Pattern</i> . Substrings begin with the characters \ (and end with the characters) . For example, the command: <pre>s/\(stu\) \(dy\) \1r\2/</pre> <p>From: The study chair</p> <p>To: The sturdy chair</p>

The letters that appear as flags change the replacement as follows:

Symbol	Use
g	Substitutes <i>String</i> for all instances of <i>Pattern</i> in the indicated line(s). Characters in <i>String</i> are not scanned for a match of <i>Pattern</i> after they are inserted. For example, the command: <pre>s/r/R/g</pre> <p>changes:</p> <p>From: the red round rock</p> <p>To: the Red Round Rock</p>
p	Prints (to STDOUT) the line that contains a successfully matched <i>Pattern</i> .
w <i>FileName</i>	Writes to <i>FileName</i> the line that contains a successfully matched <i>Pattern</i> . if <i>FileName</i> exists, it is overwritten; otherwise, it is created. A maximum of 10 different files can be mentioned as input or output files in the entire editing process. Include exactly one space between w and <i>FileName</i> .

Shared libraries and shared memory

This topic provides information about the operating system facilities provided for sharing libraries and memory allocation.

The operating system provides facilities for the creation and use of dynamically bound shared libraries. Dynamic binding allows external symbols referenced in user code and defined in a shared library to be resolved by the loader at run time.

The shared library code is not present in the executable image on disk. Shared code is loaded into memory once in the shared library segment and shared by all processes that reference it. The advantages of shared libraries are:

- Less disk space is used because the shared library code is not included in the executable programs.
- Less memory is used because the shared library code is only loaded once.
- Load time may be reduced because the shared library code may already be in memory.
- Performance may be improved because fewer page faults will be generated when the shared library code is already in memory. However, there is a performance cost in calls to shared library routines of one to eight instructions.

The symbols defined in the shared library code that are to be made available to referencing modules must be explicitly exported using an exports file, unless the `-bexpall` options is used. The first line of the file optionally contains the path name of the shared library. Subsequent lines contain the symbols to be exported.

Related information

[ar](#)

[as](#)

[dump](#)

[ipcs](#)

[ipcrm](#)

[id](#)

[pagesize](#)

[rtl_enable](#)

[update](#)

[vmstat](#)

[XCOFF](#)

Shared objects and run time linking

By default, programs are linked so that a reference to a symbol imported from a shared object is bound to that definition at load time.

This is true even if the program, or another shared object required by the program, defines the same symbol.

run time linker

A shared object that allows symbols to be rebound for appropriately linked programs

You include the run time linker in a program by linking the program with the **-brtl** option. This option has the following effects:

- A reference to the run time linker is added to your program. When program execution begins, the startup code (`/lib/crt0.o`) will call the run time linker before the main function is called.
- All input files that are shared objects are listed as dependents of your program in your program's loader section. The shared objects are listed in the same order as they were specified on the command line. This causes the system loader to load all these shared objects so that the run time linker can use their definitions. If the **-brtl** option is not used, a shared object that is not referenced by the program is not listed, even if it provides definitions that might be needed by another shared object used by the program.
- A shared object contained in an archive is only listed if the archive specifies automatic loading for the shared object member. You specify automatic loading for an archive member **foo.o** by creating a file with the following lines:

```
# autoload
#! (foo.o)
```

and adding the file as a member to the archive.

- In dynamic mode, input files specified with the **-l** flag may end in **.so**, as well as in **.a**. That is, a reference to **-lfoo** is satisfied by the first **libfoo.so** or **libfoo.a** found in any of the directories being searched. Dynamic mode is in effect by default unless the **-bstatic** option is used.

The run time linker mimics the behavior of the **ld** command when static linking is used, except that only exported symbols can be used to resolve symbols. Even when run time linking is used, the system loader must be able to load and resolve all symbol references in the main program and any module it depends on. Therefore, if a definition is removed from a module, and the main program has a reference to this definition, the program will not execute, even if another definition for the symbol exists in another module.

The run time linker can rebind all references to symbols imported from another module. A reference to a symbol defined in the same module as the reference can only be rebound if the module was built with run time linking enabled for that symbol.

Shared modules shipped with AIX 4.2 or later have run time linking enabled for most exported variables. run time linking for functions is only enabled for functions called through a function pointer. For example, as shipped, calls to the **malloc** subroutine within shared object **shr.o** in **/lib/libc.a** cannot be rebound, even if a definition of **malloc** exists in the main program or another shared module. You can link most shipped shared modules to enable run time linking for functions as well as variables by running the **rtl_enable** command.

Operation of the run time linker

The main program is loaded and resolved by the system loader in the usual way. If the executable program cannot be loaded for any reason, the **exec()** subroutine fails and the run time linker is not invoked at all. If the main program loads successfully, control passes to the run time linker, which rebinds symbols as described below. When the run time linker completes, initialization routines are called, if appropriate, and then the main function is called.

The run time linker processes modules in breadth-first search order, starting with the main executable and continuing with the direct dependents of the main executable, according to the order of dependent modules listed in each module's loader section. This order is also used when searching for the defining instance of a symbol. The "defining instance" of a symbol is usually the first instance of a symbol, but there are two exceptions. If the first instance of a symbol is an unresolved, deferred import, no defining instance exists. If the first instance is a BSS symbol (that is, with type **XTY_CM**, indicating an uninitialized variable), and there is another instance of the symbol that is neither a BSS symbol nor an unresolved, deferred import, the first such instance of the symbol is the defining instance.

The loader section of each module lists imported symbols, which are usually defined in another specified module, and exported symbols, which are usually defined in the module itself. A symbol that is imported and exported is called a "passed-through" import. Such a symbol appears to be defined in one module, although it is actually defined in another module.

Symbols can also be marked as "deferred imports." References to deferred import symbols are never rebound by the run time linker. Resolution of these symbols must be performed by the system loader, either by calling **loadbind()** or by loading a new module explicitly with **load()** or **dlopen()**.

References to imported symbols (other than deferred imports) can always be rebound. The system loader will have already resolved most imports. References to each imported symbol are rebound to the symbol's defining instance. If no defining instance exists, an error message will be printed to standard error. In addition, if the typechecking hash string of an imported symbol does not match the hash string of the defining symbol, an error message is printed.

References to exported symbols are also rebound to their defining instances, as long as the references appear in the relocation table of the loader section. (Passed-through imports are processed along with other imports, as described above.) Depending on how the module was linked, some references to exported symbols are bound at link time and cannot be rebound. Since exported symbols are defined in the exporting module, a defining instance of the symbol will always exist, unless the first instance is a deferred import, so errors are unlikely, but still possible, when rebinding exported symbols. As with imports, errors are printed if the typechecking hash strings do not match when a symbol is rebound.

Whenever a symbol is rebound, a dependency is added from the module using the symbol to the module defining the symbol. This dependency prevents modules from being removed from the address space prematurely. This is important when a module loaded by the **dlopen** subroutine defines a symbol that is still being used when an attempt is made to unload the module with the **dlclose** subroutine.

The loader section symbol table does not contain any information about the alignment or length of symbols. Thus, no errors are detected when symbols are rebound to instances that are too short or improperly aligned. Execution errors may occur in this case.

Once all modules have been processed, the run time linker calls the **exit** subroutine if any run time linking errors occurred, passing an exit code of 144 (0x90). Otherwise, execution continues by calling initialization routines or **main()**.

Creating a shared object with run time linking enabled

To create a shared object enabled for run time linking, you link with the **-G** flag. When this flag is used, the following actions take place:

1. Exported symbols are given the `nosymbolic` attribute, so that all references to the symbols can be rebound by the run time linker.
2. Undefined symbols are permitted (see the **-berok** option). Such symbols are marked as being imported from the symbolic module name `..`. Symbols imported from `..` must be resolved by the run time linker before they can be used because the system loader will not resolve these symbols.
3. The output file is given a module type of SRE, as if the **-bM:SRE** option had been specified.
4. All shared objects listed on the command line are listed as dependents of the output module, in the same manner as described when linking a program with the **-brtl** option.
5. Shared objects in archives are listed if they have the `autoLoad` attribute.

Using the **-berok** option, implied by the **-G** flag, can mask errors that could be detected at link time. If you intend to define all referenced symbols when linking a module, you should use the **-bernotok** option after the **-G** flag. This causes errors to be reported for undefined symbols.

Shared libraries and lazy loading

By default, when a module is loaded, the system loader automatically loads all of the module's dependents at the same time. Loading of dependents occurs because when a module is linked, a list of the module's dependent modules is saved in the loader section of the module.

Module	Description
dump -H	Command that allows viewing of dependent modules list.
-blazy	In AIX 4.2.1 and later, linker option that links a module so that only some of its dependents are loaded when a function in the module is first used.

When you use lazy loading, you can improve the performance of a program if most of a module's dependents are never actually used. On the other hand, every function call to a lazily loaded module has an additional overhead of about 7 instructions, and the first call to a function requires loading the defining module and modifying the function call. Therefore, if a module calls functions in most of its dependents, lazy loading may not be appropriate.

When a function defined in a lazily loaded module is called for the first time, an attempt is made to load the defining module and find the desired function. If the module cannot be found or if the function is not exported by the module, the default behavior is to print an error message to standard error and exit with a return code of 1. An application can supply its own error handler by calling the function **_lazySetErrorHandler** and supplying the address of an error handler. An error handler is called with 3 arguments: the name of the module, the name of the symbol, and an error value indicating the cause of the error. If the error handler returns, its return value should be the address of a substitute function for the desired function. The return value for **_lazySetErrorHandler** is NULL if no error handler exists, and the address of a previous handler if one exists.

Using lazy loading does not usually change the behavior of a program, but there are a few exceptions. First, any program that relies on the order that modules are loaded is going to be affected, because modules can be loaded in a different order, and some modules might not be loaded at all.

Second, a program that compares function pointers might not work correctly when lazy loading is used, because a single function can have multiple addresses. In particular, if module A calls function `f` in module B, and if lazy loading of module B was specified when module A was linked, then the address of `f` computed in module A differs from the address of `f` computed in other modules. Thus, when you use lazy loading, two function pointers might not be equal, even if they point to the same function.

Third, if any modules are loaded with relative path names and if the program changes working directories, the dependent module might not be found when it needs to be loaded. When you use lazy loading, you should use only absolute path names when referring to dependent modules at link time.

The decision to enable lazy loading is made at link time on a module-by-module basis. In a single program, you can mix modules that use lazy loading with modules that do not. When linking a single module, a reference to a variable in a dependent module prevents that module from being loaded lazily. If all references to a module are to function symbols, the dependent module can be loaded lazily.

The lazy loading facility can be used in both threaded and non-threaded applications.

Lazy loading execution tracing

A runtime feature is provided that allows you to view the loading activity as it takes place. This is accomplished using the environment variable **LDLAZYDEBUG**. The value of this variable is a number, in decimal, octal (leading 0), or hexadecimal (leading 0x) that is the sum of one or more of the following values:

Variable	Description
1	Show load or look-up errors. If a required module cannot be found, a message displays and the lazy load error handler is called. If a requested symbol is not available in the loaded referenced module, a message displays before the error handler is called.
2	Write tracing messages to stderr instead of stdout . By default, these messages are written to the standard output file stream. This value selects the standard error stream.
4	Display the name of the module being loaded. When a new module is required to resolve a function call, the name of the module that is found and loaded displays. This only occurs at the first reference to a function within that module; that is, once a module is loaded, it remains available for subsequent references to functions within that module. Additional load operations are not required.
8	Display the name of the called function. The name of the required function, along with the name of the module from which the function is expected, displays. This information displays before the module is loaded.

Named shared library areas

By default, AIX shares libraries among processes using a global set of segments, referred to as the global shared library area.

For 32-bit processes, this area consists of one segment for shared library text (segment 0xD) and one segment for pre-relocated library data (segment 0xF). Sharing text and pre-relocating data improves performance on systems where a large number of processes use common shared libraries.

Because the global shared library area is a single fixed-size resource, attempts to share a set of libraries that exceed the capacity of the area cannot succeed. In this situation, a portion of a process libraries are loaded privately. Loading libraries privately, as opposed to shared, consumes private address space in the process and places greater demands on paging space, leading to a degradation in overall system performance.

To address this limitation of the global shared library area, AIX 5.3 supports *named shared library areas*, which include the following benefits:

- A named shared library area replaces the global shared library area for a group of processes.
- A named shared library area enables a group of processes to have the full shared library capacity available to them at the same location in the effective address space as the global shared library area (segments 0xD and 0xF).
- The named shared library area feature is enabled via the **LDR_CNTRL** environment variable and no changes are required to existing binaries.
- Multiple named shared library areas can be active on the system simultaneously.
- Processes specify a particular named shared library area by a unique name. This name is chosen by the process that causes the area's creation.
- Named shared library areas are available for use only by 32-bit processes.

Because the use of a specific named shared library area is restricted to processes that request it, none of its space will be consumed by processes using the global shared library area or a different named shared library area. This decrease in competition for space in the named shared library area benefits processes using the area. These processes experience lower private address space consumption and a greater ability to share libraries among themselves. Use of a named shared library area among processes that utilize common shared libraries can optimize the processes' address space usage and decrease demands on paging space, resulting in enhanced overall system performance.

Alternate memory model (doubletext32)

In addition to the default shared library area memory model (one segment dedicated to shared library text and one segment dedicated to pre-relocated library data), named shared library areas support an alternate memory model that dedicates both segments to shared library text. This model is useful for process groups that share greater than 256 MB of library text. Note that since this alternate memory model performs no pre-relocation of library data, some performance degradation during module loading (for both exec-time dependencies and dynamically loaded modules) may be experienced. Therefore, the actual performance benefits of increased shared library text capacity should be considered on a case by case basis.

Interface

Access

A process requests the use of a named shared library area by having the **LDR_CNTRL** environment variable with the **NAMEDSHLIB** option in its environment at run time. The syntax of the new option is as follows:

```
NAMEDSHLIB=name[,attribute][,attribute2]...[,attributeN]
```

A valid name string can be any string matching the regular expression, `[A-Za-z0-9_\.]+` (containing only alphanumeric, underbar, and period characters).

A valid name string must be terminated by one of the following characters:

- @ (at sign): The delimiter for multiple **LDR_CNTRL** options
- , (comma): The delimiter for **NAMEDSHLIB** attributes
- \0 (null): The terminator of the **LDR_CNTRL** environment string

If an invalid name string is specified, the entire **NAMEDSHLIB** option is ignored. If an invalid attribute is specified, only that attribute is ignored. Currently, there is only one supported attribute: **doubletext32**.

There are no access restrictions for using named shared library areas. All requests for use of an area are granted.

Creation

There is no explicit interface to create a named shared library area. When a process requests the use of a named shared library area that does not exist, the area is automatically created.

Purging

The system removes unused libraries from a named shared library area using the same mechanisms that apply to the global shared library area:

- Automatic removal of unused libraries occurs when the area becomes full.
- Forced removal of unused libraries can be accomplished using the **slibclean** command.

Destruction

There is no explicit interface to destroy a named shared library area. When the last process using a named shared library area exits (the usecount of the area drops to zero), the area is automatically destroyed.

Attributes

The **NAMEDSHLIB** attributes are examined by the system loader only during named shared library area creation. Therefore, requests to use an existing named shared library area are not strictly required to specify attributes matching those specified at creation (the request will not fail because of an attribute mismatch). However, because the system automatically destroys unused named shared library areas, it is good practice to always specify attributes, even when you are requesting the use of an existing named shared library area.

Examples

1. Run a pair of applications using the named shared library area named *XYZ* with one segment dedicated to shared library text and one segment dedicated to pre-relocated library data by running the following commands:

```
$ export LDR_CNTRL=NAMEDSHLIB=XYZ
$ xyz_app
$ xyz_app2
```

2. Run a pair of applications using the named shared library area named *more_shtext* with both segments dedicated to shared library text by running the following commands:

```
$ export LDR_CNTRL=NAMEDSHLIB=more_shtext,doubletext32
$ mybigapp
$ mybigapp2
```

Creating a shared library

This section describes how to create a shared library.

Prerequisite tasks

1. Create one or more source files that are to be compiled and linked to create a shared library. These files contain the exported symbols that are referenced in other source files.

For the examples in this article, two source files, *share1.c* and *share2.c*, are used. The *share1.c* file contains the following code:

```
/*
*****
* share1.c: shared library source.
*****
*/

#include <stdio.h>

void func1 ()
{
    printf("func1 called\n");
}

void func2 ()
{
    printf("func2 called\n");
}
```

The `share2.c` file contains the following code:

```
/******  
 * share2.c: shared library source.  
******/  
  
void func3 ()  
{  
    printf("func3 called\n");  
}
```

The exported symbols in these files are `func1`, `func2`, and `func3`.

2. Create a main source file that references the exported symbols that will be contained in the shared library.

For the examples in this article the main source file named `main.c` is used. The `main.c` file contains the following code:

```
/******  
 * main.c: contains references to symbols defined  
 * in share1.c and share2.c  
******/  
  
#include <stdio.h>  
  
extern void func1 (),  
          func2 (),  
          func3 ();  
  
main ()  
{  
    func1 ();  
    func2 ();  
    func3 ();  
}
```

3. Create the exports file necessary to explicitly export the symbols in the shared library that are referenced by other object modules.

For the examples in this article, an exports file named `shrsub.exp` is used. The `shrsub.exp` file contains the following code:

```
#!/home/sharelib/shrsub.o  
* Above is full pathname to shared library object file  
func1  
func2  
func3
```

The `#!` line is meaningful only when the file is being used as an import file. In this case, the `#!` line identifies the name of the shared library file to be used at run time.

Procedure

1. Compile and link the two source code files to be shared. (This procedure assumes you are in the `/home/sharedlib` directory.) To compile and link the source files, enter the following commands:

```
cc -c share1.c  
cc -c share2.c  
cc -o shrsub.o share1.o share2.o -bE:shrsub.exp -bM:SRE -bnoentry
```

This creates a shared library name `shrsub.o` in the `/home/sharedlib` directory.

-bM:SRE flag

Marks the resultant object file `shrsub.o` as a re-entrant, shared library

Each process that uses the shared code gets a private copy of the data in its private process area.

flag

Sets the dummy entry point `_nostart` to override the default entry point, `_start`

-bnoentry flag

Tells the linkage editor that the shared library does not have an entry point

A shared library may have an entry point, but the system loader does not make use of an entry point when a shared library is loaded.

2. Use the following command to put the shared library in an archive file:

```
ar qv libsub.a shsub.o
```

This step is optional. Putting the shared library in an archive makes it easier to specify the shared library when linking your program, because you can use the **-l** and **-L** flags with the **ld** command.

3. Compile and link the main source code with the shared library to create the executable file. (This step assumes your current working directory contains the **main.c** file.) Use the following command:

```
cc -o main main.c -lsub -L/home/sharedlib
```

If the shared library is not in an archive, use the command:

```
cc -o main main.c /home/sharedlib/shsub.o -L/home/sharedlib
```

The program **main** is now executable. The **func1**, **func2**, and **func3** symbols have been marked for load-time deferred resolution. At run time, the system loader loads the module in to the shared library (unless the module is already loaded) and dynamically resolves the references.

Note: When creating a shared library from a C++ object, you must use the mangled C++ symbol names in the export file. However, your C++ compiler might provide an option to create the shared library for you. Refer to your compiler documentation for more information.

-L flag

Adds the specified directory (in this case, **/home/sharedlib**) to the library search path, which is saved in the loader section of the program.

At run time the library search path is used to tell the loader where to find shared libraries.

LIBPATH environment variable

A colon-separated list of directory paths that can also be used to specify a different library search path. Its format is identical to that of the **PATH** environment variable.

The directories in the list are searched to resolve references to shared objects. The **/usr/lib** and **/lib** directories contain shared libraries and should normally be included in your library search path.

Program address space overview

The Base Operating System provides a number of services for programming application program memory use.

Tools are available to assist in allocating memory, mapping memory and files, and profiling application memory usage. As background, this section describes the system's memory management architecture and memory management policy.

System memory architecture introduction

The system employs a memory management scheme that uses software to extend the capabilities of the physical hardware. Because the address space does not correspond one-to-one with real memory, the address space (and the way the system makes it correspond to real memory) is called virtual memory.

The subsystems of the kernel and the hardware that cooperate to translate the virtual address to physical addresses make up the memory management subsystem. The actions the kernel takes to ensure that processes share main memory fairly comprise the memory management policy. The following sections describe the characteristics of the memory management subsystem in greater detail.

The physical address space of 64-bit systems

The hardware provides a continuous range of virtual memory addresses, from **0x0000000000000000** to **0xFFFFFFFFFFFFFFFF**, for accessing data. The total addressable space is more than 1000000000000 terabytes. Memory access instructions generate an address of 64

bits: 36 bits to select a segment register and 28 bits to give an offset within the segment. This addressing scheme provides access to more than 64 million segments of up to 256M bytes each. Each segment register contains a 52-bit segment ID that becomes a prefix to the 28-bit offset, which together form the virtual memory address. The resulting 80-bit virtual address refers to a single, large, systemwide virtual memory space.

The process space is a 64-bit address space; that is, programs use 64-bit pointers. However, each process or interrupt handler can address only the systemwide virtual memory space (segment) whose segment IDs are in the segment register.

Segment register addressing

The system kernel loads some segment registers in the conventional way for all processes, implicitly providing the memory addressability needed by most processes. These registers include two kernel segments, and a shared-library segment, and an I/O device segment, that are shared by all processes and whose contents are read-only to non-kernel programs. There is also a segment for the **exec** system call of a process, which is shared on a read-only basis with other processes executing the same program, a private shared-library data segment that contains read-write library data, and a read-write segment that is private to the process. The remaining segment registers may be loaded using memory mapping techniques to provide more memory, or through memory access to files according to access permissions imposed by the kernel.

The system's 32-bit addressing and the access provided through indirection capabilities gives each process an interface that does not depend on the actual size of the systemwide virtual memory space. Some segment registers are shared by all processes, others by a subset of processes, and yet others are accessible to only one process. Sharing is achieved by allowing two or more processes to load the same segment ID.

Paging space

To accommodate the large virtual memory space with a limited real memory space, the system uses real memory as a work space and keeps inactive data and programs that are not mapped on disk. The area of disk that contains this data is called the paging space. A page is a unit of virtual memory that holds 4K bytes of data and can be transferred between real and auxiliary storage. When the system needs data or a program in the page space, it:

1. Finds an area of memory that is not currently active.
2. Ensures that an up-to-date copy of the data or program from that area of memory is in the paging space on disk.
3. Reads the new program or data from the paging space on disk into the newly freed area of memory.

Memory management policy

The real-to-virtual address translation and most other virtual memory facilities are provided to the system transparently by the Virtual Memory Manager (VMM). The VMM implements virtual memory, allowing the creation of segments larger than the physical memory available in the system. It accomplishes this by maintaining a list of free pages of real memory that it uses to retrieve pages that need to be brought into memory.

The VMM occasionally must replenish the pages on the free list by removing some of the current page data from real memory. The process of moving data between memory and disk as the data is needed is called "paging." To accomplish paging, the VMM uses page-stealing algorithms that categorize pages into three classes, each with unique entry and exit criteria:

- working storage pages
- local file pages
- remote file pages

In general, working pages have highest priority, followed by local file pages, and then remote file pages.

In addition, the VMM uses a technique known as the clock algorithm to select pages to be replaced. This technique takes advantage of a referenced bit for each page as an indication of what pages have been recently used (referenced). When a page-stealer routine is called, it cycles through a page frame table,

examining each page's referenced bit. If the page was unreferenced and is stealable (that is, not pinned and meets other page-stealing criteria), it is stolen and placed on the free list. Referenced pages may not be stolen, but their reference bit is reset, effectively "aging" the reference so that the page may be stolen the next time a page-stealing algorithm is issued.

Memory allocation

Version 3 of the operating system uses a delayed paging slot technique for storage allocated to applications. This means that when storage is allocated to an application with a subroutine such as **malloc**, no paging space is assigned to that storage until the storage is referenced.

Related concepts

Understanding memory mapping

The speed at which application instructions are processed on a system is proportionate to the number of access operations required to obtain data outside of program-addressable memory.

Paging space programming requirements

The amount of paging space required by an application depends on the type of activities performed on the system. If paging space runs low, processes may be lost.

System memory allocation using the malloc subsystem

Memory is allocated to applications using the **malloc** subsystem.

Understanding memory mapping

The speed at which application instructions are processed on a system is proportionate to the number of access operations required to obtain data outside of program-addressable memory.

The system provides two methods for reducing the transactional overhead associated with these external read and write operations. You can map file data into the process address space. You can also map processes to anonymous memory regions that may be shared by cooperating processes.

Memory mapped files provide a mechanism for a process to access files by directly incorporating file data into the process address space. The use of mapped files can significantly reduce I/O data movement since the file data does not have to be copied into process data buffers, as is done by the **read** and **write** subroutines. When more than one process maps the same file, its contents are shared among them, providing a low-overhead mechanism by which processes can synchronize and communicate.

Mapped memory regions, also called shared memory areas, can serve as a large pool for exchanging data among processes. The available subroutines do not provide locks or access control among the processes. Therefore, processes using shared memory areas must set up a signal or semaphore control method to prevent access conflicts and to keep one process from changing data that another is using. Shared memory areas can be most beneficial when the amount of data to be exchanged between processes is too large to transfer with messages, or when many processes maintain a common large database.

The system provides two methods for mapping files and anonymous memory regions. The following subroutines, known collectively as the **shmat** services, are typically used to create and use shared memory segments from a program:

Subroutine	Definition
-------------------	-------------------

<u>shmctl</u>	Controls shared memory operations
----------------------	-----------------------------------

<u>shmget</u>	Gets or creates a shared memory segment
----------------------	---

<u>shmat</u>	Attaches a shared memory segment from a process. Does not allow you to map block devices.
---------------------	---

<u>shmdt</u>	Detaches a shared memory segment from a process
---------------------	---

<u>mprotect</u>	Modifies the access protections of a specified address range within a shared memory segment.
------------------------	--

<u>disclaim</u>	Removes a mapping from a specified address range within a shared memory segment
------------------------	---

The **ftok** subroutine provides the key that the **shmget** subroutine uses to create the shared segment

The second set of services, collectively known as the **mmap** services, is typically used for mapping files, although it may be used for creating shared memory segments as well.

All operations valid on memory resulting from `mmap()` of a file are valid on memory resulting from `mmap()` of a block device. A block device is a special file that provides access to a device driver that presents a block interface. A block interface to a device driver requires data access in blocks of a fixed size. The interface is typically used for data storage devices.

The `mmap` services include the following subroutines:

Subroutine	Definition
-------------------	-------------------

<u>madvise</u>	Advises the system of a process' expected paging behavior
-----------------------	---

<u>mincore</u>	Determines residency of memory pages
-----------------------	--------------------------------------

<u>mmap</u>	Maps an object file into virtual memory. Allows you to map block devices one process at a time.
--------------------	---

<u>mprotect</u>	Modifies the access protections of memory mapping
------------------------	---

<u>msync</u>	Synchronizes a mapped file with its underlying storage device
---------------------	---

<u>munmap</u>	Unmaps a mapped memory region
----------------------	-------------------------------

The **msem_init**, **msem_lock**, **msem_unlock**, **msem_remove**, **msleep**, and **mwakeup** subroutines provide access control for the processes mapped using the **mmap** services.

Refer to the following sections to learn more about memory mapping:

Comparing `mmap` with `shmat`

As with the `shmat` services, the portion of the process address space available for mapping files with the `mmap` services is dependent on whether a process is a 32-bit process or a 64-bit process. For 32-bit processes, the portion of address space available for mapping consists of addresses in the range of `0x30000000-0xCFFFFFFF`, for a total of 2.5G bytes of address space. The portion of address space available for mapping files consists of addresses in the ranges of `0x30000000-0xCFFFFFFF` and `0xE0000000-0xEFFFFFFF` for a total of 2.75G bytes of address space. In AIX 5.2 and later, a 32-bit process run with the very large address-space model has the range `0x30000000-0xFFFFFFFF` available for mappings, with a total of up to 3.25GB of address space.

All available ranges within the 32-bit process address space are available for both fixed-location and variable-location mappings. Fixed-location mappings occur when applications specify that a mapping be placed at a fixed location within the address space. Variable-location mappings occur when applications specify that the system should decide the location at which a mapping should be placed.

For 64-bit processes, two sets of address ranges with the process address space are available for **mmap** or **shmat** mappings. The first, consisting of the single range `0x07000000_00000000-0x07FFFFFF_FFFFFFFF`, is available for both fixed-location and variable-location mappings. The second set of address ranges is available for fixed-location mappings only and consists of the ranges `0x30000000-0xCFFFFFFF`, `0xE0000000-0xEFFFFFFF`, and `0x10_00000000-0x06FFFFFF_FFFFFFFF`. The last range of this set, consisting of `0x10_00000000-0x06FFFFFF_FFFFFFFF`, is also made available to system loader to hold program text, data and heap, so only unused portions of the range are available for fixed-location mappings.

Both the **mmap** and **shmat** services provide the capability for multiple processes to map the same region of an object such that they share addressability to that object. However, the **mmap** subroutine extends this capability beyond that provided by the **shmat** subroutine by allowing a relatively unlimited number of such mappings to be established. While this capability increases the number of mappings supported per file object or memory segment, it can prove inefficient for applications in which many processes map the same file data into their address space.

The **mmap** subroutine provides a unique object address for each process that maps to an object. The software accomplishes this by providing each process with a unique virtual address, known as an alias. The **shmat** subroutine allows processes to share the addresses of the mapped objects.

Because only one of the existing aliases for a given page in an object has a real address translation at any given time, only one of the **mmap** mappings can make a reference to that page without incurring a page fault. Any reference to the page by a different mapping (and thus a different alias) results in a page fault that causes the existing real-address translation for the page to be invalidated. As a result, a new translation must be established for it under a different alias. Processes share pages by moving them between these different translations.

For applications in which many processes map the same file data into their address space, this toggling process may have an adverse affect on performance. In these cases, the **shmat** subroutine may provide more efficient file-mapping capabilities.

Note: On systems with PowerPC processors, multiple virtual addresses can exist for the same real address. A real address can be aliased to different effective addresses in different processes without toggling. Because there is no toggling, there is no performance degradation.

Use the **shmat** services under the following circumstances:

- For 32-bit application, eleven or fewer files are mapped simultaneously, and each is smaller than 256MB.
- When mapping files larger than 256MB.
- When mapping shared memory regions which need to be shared among unrelated processes (no parent-child relationship).
- When mapping entire files.

Use **mmap** under the following circumstances:

- Portability of the application is a concern.
- Many files are mapped simultaneously.
- Only a portion of a file needs to be mapped.
- Page-level protection needs to be set on the mapping.
- Private mapping is required.

An "extended **shmat**" capability is available for 32-bit applications with their limited address spaces. If you define the environment variable **EXTSHM=ON**, then processes executing in that environment can create and attach more than eleven shared memory segments. The process can attach these segments into the address space for the size of the segment. Another segment can be attached at the end of the first one in the same 256M byte region. The address at which a process can attach is at page boundaries, which is a multiple of **SHMLBA_EXTSHM** bytes.

Some restrictions exist on the use of the extended **shmat** feature. These shared memory regions cannot be used as I/O buffers where the unpinning of the buffer occurs in an interrupt handler. The restrictions on the use of extended **shmat** I/O buffers is the same as that of **mmap** buffers.

The environment variable provides the option of executing an application with either the additional functionality of attaching more than 11 segments when **EXTSHM=ON**, or the higher-performance access to 11 or fewer segments when the environment variable is not set. Again, the "extended **shmat**" capability only applies to 32-bit processes.

mmap Compatibility Considerations

The **mmap** services are specified by various standards and commonly used as the file-mapping interface of choice in other operating system implementations. However, the system's implementation of the **mmap** subroutine may differ from other implementations. The **mmap** subroutine incorporates the following modifications:

- Mapping into the process private area is not supported.

- Mappings are not implicitly unmapped. An **mmap** operation which specifies **MAP_FIXED** will fail if a mapping already exists within the range specified.
- For private mappings, the copy-on-write semantic makes a copy of a page on the first write reference.
- Mapping of I/O or device memory is not supported.
- Mapping of character devices or use of an **mmap** region as a buffer for a read-write operation to a character device is not supported.
- The **madvise** subroutine is provided for compatibility only. The system takes no action on the advice specified.
- The **mprotect** subroutine allows the specified region to contain unmapped pages. In operation, the unmapped pages are simply skipped over.
- The OSF/AES-specific options for default exact mapping and for the **MAP_INHERIT**, **MAP_HASSEMAPHORE**, and **MAP_UNALIGNED** flags are not supported.

Using the semaphore subroutines

The **msem_init**, **msem_lock**, **msem_unlock**, **msem_remove**, **msleep** and **mwakeup** subroutines conform to the OSF Application Environment specification. They provide an alternative to IPC interfaces such as the **semget** and **semop** subroutines. Benefits of using the semaphores include an efficient serialization method and the reduced overhead of not having to make a system call in cases where there is no contention for the semaphore.

Semaphores should be located in a shared memory region. Semaphores are specified by **msemaphore** structures. All of the values in a **msemaphore** structure should result from a **msem_init** subroutine call. This call may or may not be followed by a sequence of calls to the **msem_lock** subroutine or the **msem_unlock** subroutine. If a **msemaphore** structure values originated in another manner, the results of the semaphore subroutines are undefined.

The address of the **msemaphore** structure is significant. You should be careful not to modify the structure's address. If the structure contains values copied from a **msemaphore** structure at another address, the results of the semaphore subroutines are undefined.

The semaphore subroutines may prove less efficient when the semaphore structures exist in anonymous memory regions created with the **mmap** subroutine, particularly in cases where many processes reference the same semaphores. In these instances, the semaphore structures should be allocated out of shared memory regions created with the **shmget** and **shmat** subroutines.

Mapping files with the shmat subroutine

Mapping can be used to reduce the overhead involved in writing and reading the contents of files. Once the contents of a file are mapped to an area of user memory, the file may be manipulated as if it were data in memory, using pointers to that data instead of input/output calls. The copy of the file on disk also serves as the paging area for that file, saving paging space.

A program can use any regular file as a mapped data file. You can also extend the features of mapped data files to files containing compiled and executable object code. Because mapped files can be accessed more quickly than regular files, the system can load a program more quickly if its executable object file is mapped to a file.

To create a program as a mapped executable file, compile and link the program using the **-K** flag with the **cc** or **ld** command. The **-K** flag tells the linker to create an object file with a page-aligned format. That is, each part of the object file starts on a page boundary (an address that can be divided by 2K bytes with no remainder). This option results in some empty space in the object file but allows the executable file to be mapped into memory. When the system maps an object file into memory, the text and data portions are handled differently.

Copy-on-write mapped files

To prevent changes made to mapped files from appearing immediately in the file on disk, map the file as a copy-on-write file. This option creates a mapped file with changes that are saved in the system paging space, instead of to the copy of the file on disk. You must choose to write those changes to the copy on disk to save the changes. Otherwise, you lose the changes when closing the file.

Because the changes are not immediately reflected in the copy of the file that other users may access, use copy-on-write mapped files only among processes that cooperate with each other.

The system does not detect the end of files mapped with the **shmat** subroutine. Therefore, if a program writes beyond the current end of file in a copy-on-write mapped file by storing into the corresponding memory segment (where the file is mapped), the actual file on disk is extended with blocks of zeros in preparation for the new data. If the program does not use the **fsync** subroutine before closing the file, the data written beyond the previous end of file is not written to disk. The file appears larger, but contains only the added zeros. Therefore, always use the **fsync** subroutine before closing a copy-on-write mapped file to preserve any added or changed data.

Mapping shared memory segments with the shmat subroutine

The system uses shared memory segments similarly to the way it creates and uses files. Defining the terms used for shared memory with respect to the more familiar file-system terms is critical to understanding shared memory. A definition list of shared memory terms follows:

Term	Definition
key	The unique identifier of a particular shared segment. It is associated with the shared segment as long as the shared segment exists. In this respect, it is similar to the <i>file name</i> of a file.
shmid	The identifier assigned to the shared segment for use within a particular process. It is similar in use to a <i>file descriptor</i> for a file.
attach	Specifies that a process must attach a shared segment in order to use it. Attaching a shared segment is similar to opening a file.
detach	Specifies that a process must detach a shared segment once it is finished using it. Detaching a shared segment is similar to closing a file.

Related concepts

[Program address space overview](#)

The Base Operating System provides a number of services for programming application program memory use.

[Creating a shared memory segment with the shmat subroutine](#)

This section describes how to create a shared memory segment with the **shmat** subroutine

[Interprocess communication limits](#)

This topic describes system limits for Interprocess communication (IPC) mechanisms.

Interprocess communication limits

This topic describes system limits for Interprocess communication (IPC) mechanisms.

On some UNIX systems, the system administrator can edit the **/etc/master** file and define limits for IPC mechanisms (semaphores, shared memory segments, and message queues). The problem with this method is that the higher the limits, the more memory the operating system uses, and performance can be adversely affected.

AIX uses a different method. In AIX, upper limits are defined for the IPC mechanisms, which are not configurable. The individual IPC data structures are allocated and deallocated as needed, so memory requirements depend on the current system usage of IPC mechanisms.

This difference in methods sometimes confuses users who are installing or using databases. The limit that is the most confusing is the maximum number of shared memory segments that can be attached simultaneously per process. For 64-bit processes, the maximum number of shared memory segments is 268435456. For 32-bit processes, the maximum number of shared memory segments is 11, unless the extended **shmat** capability is used.

The following tables summarize the semaphore limits on IPC mechanisms.

Semaphores	4.3.0	4.3.1	4.3.2	5.1	5.2	5.3	7.1
Maximum number of semaphore IDs for 32-bit kernel	4096	4096	131072	131072	131072	131072	N/A
Maximum number of semaphore IDs for 64-bit kernel	4096	4096	131072	131072	131072	1048576	1048576
Maximum semaphores per semaphore ID	65535	65535	65535	65535	65535	65535	65535
Maximum operations per semop call	1024	1024	1024	1024	1024	1024	1024
Maximum undo entries per process	1024	1024	1024	1024	1024	1024	1024
Size in bytes of undo structure	8208	8208	8208	8208	8208	8208	8208
Semaphore maximum value	32767	32767	32767	32767	32767	32767	32767
Adjust on exit maximum value	16384	16384	16384	16384	16384	16384	16384

The following tables summarize the message queue limits on IPC mechanisms.

Message queue	4.3.0	4.3.1	4.3.2	5.1	5.2	5.3	7.1
Maximum message size	4 MB	4 MB	4 MB	4 MB	4 MB	4 MB	4 MB
Maximum bytes on queue	4 MB	4 MB	4 MB	4 MB	4 MB	4 MB	4 MB
Maximum number of message queue IDs for 32-bit kernel	4096	4096	131072	131072	131072	131072	131072
Maximum number of message queue IDs for 64-bit kernel	4096	4096	131072	131072	131072	1048576	1048576
Maximum messages per queue ID	524288	524288	524288	524288	524288	524288	524288

The following tables summarize the shared memory limits on IPC mechanisms.

Shared memory	4.3.0	4.3.1	4.3.2	5.1	5.2	5.3	7.1
Maximum segment size (32-bit process)	256 MB	2 GB	2 GB	2 GB	2 GB	2 GB	2 GB

Shared memory	4.3.0	4.3.1	4.3.2	5.1	5.2	5.3	7.1
Maximum segment size (64-bit process) for 32-bit kernel	256 MB	2 GB	2 GB	64 GB	1 TB	1 TB	N/A
Maximum segment size (64-bit process) for 64-bit kernel	256 MB	2 GB	2 GB	64 GB	1 TB	32 TB	32 TB
Minimum segment size	1	1	1	1	1	1	1
Maximum number of shared memory IDs (32-bit kernel)	4096	4096	131072	131072	131072	131072	131072
Maximum number of shared memory IDs (64-bit kernel)	4096	4096	131072	131072	131072	1048576	1048576
Maximum number of segments per process (32-bit process)	11	11	11	11	11	11	11
Maximum number of segments per process (64-bit process)	268435456	268435456	268435456	268435456	268435456	268435456	268435456

Note: For 32-bit processes, the maximum number of segments per process is limited only by the size of the address space when the extended **shmat** capability is used.

IPC limits on AIX4.3

- For semaphores and message queues, the table shows the system limits
- For shared memory, the maximum shared memory segment size is 256GB.

- For shared memory without the extended **shmat** capability:
 - A process can attach a maximum of 11 shared memory segments.
- For shared memory with the extended **shmat** capability:
 - When a shared memory segments is attached, its size is rounded to a multiple of 4096 bytes
 - A process can attach as many shared memory segments as will fit in the available address space. The maximum available address space size is 11 segments, or 11 times 256 MB.
- The extended **shmat** capability is used if the environment variable **EXTSHM** has the value ON when the process starts executing.
- The available address space for attaching shared memory segments is reduced if the large or very large address-space model is used.

IPC limits on AIX 4.3.1

- The maximum size of a shared memory segment increases from 256 MB to 2 GB. When a shared memory segment larger than 256 MB is attached, its size is rounded to a multiple of 256 MB, even if the extended **shmat** capability is being used.

IPC limits on AIX 4.3.2

- The maximum number of message queues, semaphore IDs, and shared memory segments is 131072.
- The maximum number of messages per queue is 524288.

IPC limits on AIX 5.1

- The maximum size of a shared memory segment for 64-bit processes is 64 GB. A 32-bit process cannot attach a shared memory segment larger than 2 GB.

IPC limits on AIX 5.2

- The maximum size of a shared memory segment for 64-bit processes is 1 TB. A 32-bit process cannot attach a shared memory segment larger than 2 GB.
- 32-bit applications can use the **shmat** capability to obtain more than 11 segments when using the very large address space model without having to use extended **shmat**. For more information on the very large address space model.
- Applications can query the IPC limits on the system using the **vmgetinfo** system call.

IPC limits on AIX 5.3

- The maximum size of a shared memory segment for 64-bit processes is 32 TB. A 32-bit process cannot attach a shared memory segment larger than 2 GB.
- 32-bit applications can use the **shmat** capability to obtain more than 11 segments when using the very large address space model without having to use extended **shmat**. For more information on the very large address space model.
- Applications can query the IPC limits on the system using the **vmgetinfo** system call.

IPC limits on AIX 6.1

The 32-bit kernel is no longer supported on AIX 6.1. All other values are the same.

Related concepts

Understanding memory mapping

The speed at which application instructions are processed on a system is proportionate to the number of access operations required to obtain data outside of program-addressable memory.

Creating a mapped data file with the shmat subroutine

This section describes how to create a mapped data file using **shmat** subroutine.

Prerequisite condition

The file to be mapped is a regular file.

Procedure

The creation of a mapped data file is a two-step process. First, you create the mapped file. Then, because the **shmat** subroutine does not provide for it, you must program a method for detecting the end of the mapped file.

1. To create the mapped data file:
 - a. Open (or create) the file and save the file descriptor:

```
if( ( fildes = open( filename , 2 ) ) < 0 )
{
    printf( "cannot open file\n" );
    exit(1);
}
```

- b. Map the file to a segment with the **shmat** subroutine:

```
file_ptr=shmat (fildes, 0, SHM_MAP);
```

The SHM_MAP constant is defined in the **/usr/include/sys/shm.h** file. This constant indicates that the file is a mapped file. Include this file and the other shared memory header files in a program with the following directives:

```
#include <sys/shm.h>
```

2. To detect the end of the mapped file:
 - a. Use the **lseek** subroutine to go to the end of file:

```
eof = file_ptr + lseek(fildes, 0, 2);
```

This example sets the value of eof to an address that is 1 byte beyond the end of file. Use this value as the end-of-file marker in the program.

- b. Use **file_ptr** as a pointer to the start of the data file, and access the data as if it were in memory:

```
while ( file_ptr < eof)
{
    :
    :
    (references to file using file_ptr)
}
```

Note: The **read** and **write** subroutines also work on mapped files and produce the same data as when pointers are used to access the data.

- c. Close the file when the program is finished working with it:

```
close (fildes );
```

Creating a copy-on-write mapped data file with the shmat subroutine

This section describes how to create a copy-on-write mapped data file with the **shmat** subroutine.

Prerequisite Condition

The file to be mapped is a regular file.

Procedure

1. Open (or create) the file and save the file descriptor:

```
if( ( fildes = open( filename , 2 ) ) < 0 )
{
    printf( "cannot open file\n" );
    exit(1);
}
```


2. Map the file to a segment as copy-on-write, with the **shmat** subroutine:

```
file_ptr = shmat( fildes, 0, SHM_COPY );
```

The SHM_COPY constant is defined in the `/usr/include/sys/shm.h` file. This constant indicates that the file is a copy-on-write mapped file. Include this header file and other shared memory header files in a program with the following directives:

```
#include <sys/shm.h>
```

3. Use `file_ptr` as a pointer to the start of the data file, and access the data as if it were in memory.

```
while ( file_ptr < eof)
{
    .
    .
    . (references to file using file_ptr)
}
```

4. Use the **fsync** subroutine to write changes to the copy of the file on disk to save the changes:

```
fsync( fildes );
```

5. Close the file when the program is finished working with it:

```
close( fildes );
```

Creating a shared memory segment with the shmat subroutine

This section describes how to create a shared memory segment with the **shmat** subroutine

Prerequisite tasks or conditions

None.

Procedure

1. Create a key to uniquely identify the shared segment. Use the **ftok** subroutine to create the key. For example, to create the key `mykey` using a project ID of `R` contained in the variable `proj` (type **char**) and a file name of `null_file`, use a statement like:

```
mykey = ftok( null_file, proj );
```

2. Either:

- Create a shared memory segment with the **shmget** subroutine. For example, to create a shared segment that contains 4096 bytes and assign the **shmid** to an integer variable `mem_id`, use a statement like:

```
mem_id = shmget(mykey, 4096, IPC_CREAT | 0666 );
```

- Get a previously created shared segment with the **shmget** subroutine. For example, to get a shared segment that is already associated with the key `mykey` and assign the **shmid** to an integer variable `mem_id`, use a statement like:

```
mem_id = shmget( mykey, 4096, IPC_ACCESS );
```

3. Attach the shared segment to the process with the **shmat** subroutine. For example, to attach a previously created segment, use a statement like:

```
ptr = shmat( mem_id, 0, 0 );
```

In this example, the variable `ptr` is a pointer to a structure that defines the fields in the shared segment. Use this template structure to store and retrieve data in the shared segment. This template should be the same for all processes using the segment.

4. Work with the data in the segment using the template structure.
5. Detach from the segment using the [shmdt](#) subroutine:

```
shmdt( ptr );
```

6. If the shared segment is no longer needed, remove it from the system with the [shmctl](#) subroutine:

```
shmctl( mem_id, IPC_RMID, ptr );
```

Note: You can also use the [ipcs](#) command to get information about a segment, and the [ipcrm](#) command to remove a segment.

Related concepts

Understanding memory mapping

The speed at which application instructions are processed on a system is proportionate to the number of access operations required to obtain data outside of program-addressable memory.

Paging space programming requirements

The amount of paging space required by an application depends on the type of activities performed on the system. If paging space runs low, processes may be lost.

If paging space runs out, the system may panic. When a paging space low condition is detected, additional paging space should be defined.

The system monitors the number of free paging space blocks and detects when a paging space shortage exists. The [vmstat](#) command obtains statistics related to this condition. When the number of free paging space blocks falls below a threshold known as the paging space warning level, the system informs all processes (excepts [kprocs](#)) of the low condition by sending the **SIGDANGER** signal.

Note: If the shortage continues and falls below a second threshold known as the paging space kill level, the system sends the SIGKILL signal to processes that are the major users of paging space and that do not have a signal handler for the SIGDANGER signal. (The default action for the SIGDANGER signal is to ignore the signal). The system continues sending SIGKILL signals until the number of free paging space blocks is above the paging space kill level. If the `low_ps_handling` parameter is set to 2 (under the `vmo` command) and if no process is found to kill (without the SIGDANGER handler), the system sends the SIGKILL signal to the earliest processes that have a signal handler for the SIGDANGER signal.

Processes that dynamically allocate memory can ensure that sufficient paging space exists by monitoring the paging space levels with the [psdanger](#) subroutine or by using special allocation routines. Processes can avoid being ended when the paging space kill level is reached by defining a signal handler for the **SIGDANGER** signal and by using the [disclaim](#) subroutine to release memory and paging space resources allocated in the data and stack areas, and in shared memory segments.

Other subroutines that can assist in dynamically retrieving paging information from the VMM include the following:

Subroutine	Description
------------	-------------

mincore	Determines the residency of memory pages.
madvise	Permits a process to advise the system about its expected paging behavior.
swapqry	Returns paging device status.
swapon	Activates paging or swapping to a designated block device.

Related concepts

[Program address space overview](#)

The Base Operating System provides a number of services for programming application program memory use.

List of memory manipulation services

The memory functions operate on arrays of characters in memory called memory areas.

These subroutines enable you to:

- Locate a character within a memory area
- Copy characters between memory areas
- Compare contents of memory areas
- Set a memory area to a value.

You do not need to specify any special flag to the compiler in order to use the memory functions. However, you must include the header file for these functions in your program. To include the header file, use the following statement:

```
#include <memory.h>
```

The following memory services are provided:

Service	Description
<u>compare_and_swap</u>	Compares and swaps data
<u>fetch_and_add</u>	Updates a single word variable atomically
<u>fetch_and_and</u> or <u>fetch_and_or</u>	Set or clear bits in a single word variable atomically
<u>malloc</u>, <u>free</u>, <u>realloc</u>, <u>calloc</u>, <u>mallopt</u>, <u>mallinfo</u>, or <u>alloca</u>	Allocate memory
<u>memccpy</u>, <u>memchr</u>, <u>memcmp</u>, <u>memcpy</u>, <u>memset</u> or <u>memmove</u>	Perform memory operations.

Service	Description
<u>moncontrol</u>	Starts and stops execution profiling after initialization by the monitor subroutine
<u>monitor</u>	Starts and stops execution profiling using data areas defined in the function parameters
<u>monstartup</u>	Starts and stops execution profiling using default-sized data areas
<u>mprotect</u>	Modifies the access protections of a specified address range within a shared memory segment.
<u>msem_init</u>	Initializes a semaphore in a mapped file or shared memory region
<u>msem_lock</u>	Locks a semaphore
<u>msem_remove</u>	Removes a semaphore
<u>msem_unlock</u>	Unlocks a semaphore
<u>msleep</u>	Puts a process to sleep when a semaphore is busy
<u>mwakeup</u>	Wakes up a process that is waiting on a semaphore
<u>disclaim</u>	Disclaims the content of a memory address range
<u>ftok</u>	Generates a standard interprocess communication key
<u>getpagesize</u>	Gets the system page size
<u>psdanger</u>	Defines the amount of free paging space available

Service	Description
<u>shmat</u>	Attaches a shared memory segment or a mapped file to the current process
<u>shmctl</u>	Controls shared memory operations
<u>shmdt</u>	Detaches a shared memory segment
<u>shmget</u>	Gets a shared memory segment
<u>swapon</u>	Activates paging or swapping to a designated block device
<u>swapqry</u>	Returns device status

List of memory mapping services

The memory mapping subroutines operate on memory regions that have been mapped with the **mmap** subroutine.

These subroutines enable you to:

- Map an object file into virtual memory
- Synchronize a mapped file
- Determine residency of memory pages
- Determine access protections to a mapped memory region
- Unmap mapped memory regions.

You do not need to specify any special flag to the compiler to use the memory functions. However, you must include the header file for some of these subroutines. If the subroutine description specifies a header file, you can include it with the following statement:

```
#include <HeaderFile.h>
```

The following memory mapping services are provided:

Service	Description
<u>madvise</u>	Advises the system of a process' expected paging behavior.
<u>mincore</u>	Determines residency of memory pages.
<u>mmap</u>	Maps an object file onto virtual memory.
<u>mprotect</u>	Modifies access protections of memory mapping.
<u>msync</u>	Synchronizes a mapped file with its underlying storage device.
<u>munmap</u>	Unmaps a mapped memory region.

AIX vector programming

Some PowerPC processors implement a Single Instruction Multiple Data (SIMD)-style vector extension.

Often referred to as AltiVec or VMX, the vector extension to the PowerPC architecture provides an additional instruction set for performing vector and matrix mathematical functions.

The Vector Arithmetic Logic Unit is an SIMD-style arithmetic unit, in which a single instruction performs the same operation on all the data elements of each vector. AIX 5.3 with recommended technology level 5300-30 is the first AIX release to enable vector programming. The IBM PowerPC 970 processor is the first processor supported by AIX that implements the vector extension. These processors are currently found in the JS20 blade servers offered with the BladeCenter.

Vector extension overview

The vector extension consists of an additional set of 32 128-bit registers that can contain a variety of vectors including signed or unsigned 8-bit, 16-bit, or 32-bit integers, or 32-bit IEEE single-precision floats. There is a vector status and control register that contains a sticky status bit indicating saturation, as well as a control bit for enabling Java or non-Java mode for floating-point operations.

The default mode initialized by AIX for every new process is Java-mode enabled, which provides IEEE-compliant floating-point operations. The alternate non-Java mode results in a less precise mode for floating point computations, which might be significantly faster on some implementations and for specific operations. For example, on the PowerPC 970 processor running in Java mode, some vector floating-point instructions will encounter an exception if the input operands or result are denormal, resulting in costly emulation by the operating system. For this reason, you are encouraged to consider explicitly enabling the non-Java mode if the rounding is acceptable, or to carefully attempt to avoid vector computations on denormal values.

The vector extension also includes more than 160 instructions providing load and store access between vector registers and memory, in register manipulation, floating point arithmetic, integer arithmetic and logical operations, and vector comparison operations. The floating point arithmetic instructions use the IEEE 754-1985 single precision format, but do not report IEEE exceptions. Default results are produced for all exception conditions as specified by IEEE for untrapped exceptions. Only IEEE default round-to-nearest rounding mode is provided. No floating-point division or square-root instructions are provided, but instead a reciprocal estimate instruction is provided for division, and a reciprocal square root estimate instruction is provided for square root.

There is also a 32-bit special purpose register that is managed by software to represent a bitmask of vector registers in use. This allows the operating system to optimize vector save and restore algorithms as part of context switch management.

Runtime determination of vector capability

A program can determine whether a system supports the vector extension by reading the **vmx_version** field of the `_system_configuration` structure. If this field is non-zero, then the system processors and operating system contain support for the vector extension. A `__power_vmx()` macro is provided in `/usr/include/sys/systemcfg.h` for performing this test. This can be useful for software that conditionally exploits the vector extension when present, or to use functionally equivalent scalar code paths when not present.

AIX ABI extension

The AIX Application Binary Interface (ABI) has been extended to support the addition of vector register state and conventions. Refer to the *Assembler Language Reference* for a complete description of the ABI extensions.

AIX supports the AltiVec programming interface specification. Below is a table of the C and C++ vector data types. All vector data types are 16 bytes in size, and must be aligned on a 16-byte boundary. Aggregates containing vector types must follow normal conventions of aligning the aggregate to the requirement of its largest member. If an aggregate containing a vector type is packed, then there is no guarantee of 16-byte alignment of the vector type. An AIX compiler supporting the AltiVec programming interface specification is required.

New C and C++ types	Contents
vector unsigned characters	16 unsigned characters
vector signed characters	16 signed characters
vector bool characters	16 unsigned characters
vector unsigned short	8 unsigned short

Table 82. New C and C++ Vector Data Types (continued)

New C and C++ types	Contents
vector signed short	8 signed short
vector bool short	8 unsigned short
vector unsigned integers	4 unsigned integers
vector signed integers	4 signed integers
vector bool integers	4 unsigned integers
vector float	4 float

The following table outlines the vector register usage conventions.

Table 83. Vector Register Conventions

Register type	Register	Status	Use
VRs	VR0	Volatile	Scratch register
	VR1	Volatile	Scratch register
	VR2	Volatile	First vector argument First vector of function return value
	VR3	Volatile	Second vector argument, scratch
	VR4	Volatile	Third vector argument, scratch
	VR5	Volatile	Fourth vector argument, scratch
	VR6	Volatile	Fifth vector argument, scratch
	VR7	Volatile	Sixth vector argument, scratch
	VR8	Volatile	Seventh vector argument, scratch
	VR9	Volatile	Eighth vector argument, scratch
	VR10	Volatile	Ninth vector argument, scratch
	VR11	Volatile	Tenth vector argument, scratch
	VR12	Volatile	Eleventh vector argument, scratch
	VR13	Volatile	Twelfth vector argument, scratch
	VR14:19	Volatile	Scratch

Table 83. Vector Register Conventions (continued)

Register type	Register	Status	Use
	VR20:31	Reserved (default mode) Nonvolatile (extended ABI mode)	When the default Vector enabled mode is used, these registers are reserved and must not be used. In the extended ABI vector enabled mode, these registers are nonvolatile and their values are preserved across function calls
Special Purpose	VRSAVE	Reserved	In the AIX ABI, VRSAVE is not used. An ABI-compliant program must not use or alter VRSAVE.
Special Purpose	VSCR	Volatile	The Vector Status and Control Register contains the saturation status bit and non-Java mode control bit.

The AltiVec Programming Interface Specification defines the VRSAVE register to be used as a bitmask of vector registers in use. AIX requires that an application never modify the VRSAVE register.

The first 12 vector parameters in a function are placed in VR2 through VR13. Unneeded vector parameter registers contain undefined values upon entry to the function. Non-variable length argument list vector parameters are not shadowed in general purpose registers (GPRs). Any additional vector parameters, from 13th and beyond, are passed through memory on the program stack, 16-byte aligned, in their appropriate mapped location within the parameter region corresponding to their position in the parameter list.

For variable length argument lists, **va_list** continues to be a pointer to the memory location of the next parameter. When **va_arg()** accesses a vector type, **va_list** must first be aligned to a 16-byte boundary. The receiver or consumer of a variable-length argument list is responsible for performing this alignment prior to retrieving the vector type parameter.

A non-packed structure or union passed by a value that has a vector member anywhere within it will be aligned to a 16-byte boundary on the stack.

A function that takes a variable-length argument list has all parameters mapped in the argument area ordered and aligned according to their types. The first eight words (32-bit) or doublewords (64-bit) of a variable-length argument list are shadowed in GPRs r3 - r10. This includes vector parameters.

Functions that have a return value declared as a vector data type place the return value in VR2. Any function that returns a vector type or has vector parameters requires a function prototype. This avoids the compiler shadowing the VRs in GPRs for the general case.

Legacy ABI compatibility and interoperability

Due to the nature of interfaces (such as `setjmp()`, `longjmp()`, `sigsetjmp()`, `siglongjmp()`, `_setjmp()`, `_longjmp()`, `getcontext()`, `setcontext()`, `makecontext()`, and `swapcontext()`) that must save and restore nonvolatile machine state, there is risk introduced when considering dependencies between legacy and vector-extended ABI modules. To complicate matters, the `setjmp` family of functions in `libc` reside in a static member of `libc`, which means that every existing AIX binary has a statically bound copy of the `setjmp` family and others that existed with the version of AIX it was linked against. Furthermore, existing

AIX binaries have `jmpbufs` and `ucontext` data structure definitions that are insufficient to house any additional nonvolatile vector register state.

Any cases where legacy modules and new modules interleave calls, or call-backs, where a legacy module could perform a `longjmp()` or `setcontext()`, bypassing normal linkage convention of a vector extended module, has a risk of compromising the nonvolatile vector register state.

For this reason, while the AIX ABI defines nonvolatile vector registers, the default compilation mode when using vectors (AltiVec) in AIX compilers is to not use any of the nonvolatile vector registers. This results in a default compilation environment that safely allows exploitation of vectors (AltiVec) while introducing no risk with respect to interoperability with legacy binaries.

For applications where interoperability and module dependence is completely known, an additional compilation option can be enabled that will allow the use of nonvolatile vector registers. This mode should only be used when all dependent legacy modules and behaviors are fully known and understood as either having no dependence on functions such as `setjmp()`, `sigsetjmp()`, `_setjmp()`, or `getcontext()`, or ensuring that all module transitions are performed using normal subroutine linkage convention, and that no call-backs to an upstream legacy module are used.

The default AltiVec compilation environment predefines `__VEC__`, in accordance with the *AltiVec Technology Programming Interface Manual*.

When the option to use nonvolatile vector registers is enabled, the compilation environment must also predefine `__EXTABI__`. You can compile non-vector enabled modules to be extended ABI-aware by explicitly defining `__AIXEXTABI`. This will ensure that those modules can safely interact with vector-enabled modules that are enabled to use nonvolatile vector registers.

Extended context

In order to support the additional machine state required by the vector extension as well as other extensions such as user keys, AIX 5.3 introduced support for extended context structures. The primary application-visible use of machine-context information is its presence in the `sigcontext` structure provided to signal handlers, and the resulting activation of the machine context in the `sigcontext` upon return from the signal handler. The `sigcontext` structure is actually a subset of the larger `ucontext` structure. The two structures are identical for up to `sizeof(struct sigcontext)`. When AIX builds a signal context to be passed to a signal handler, it actually builds a `ucontext` structure on the signal handler's stack. The machine-context portion of a signal context must contain all of the active machine state, volatile and nonvolatile, for the involuntarily interrupted context. To accomplish this without affecting binary compatibility with existing signal handlers, space previously reserved in the `ucontext` structure now serves as an indication of whether extended context information is available.

A newly defined field in the `ucontext`, `__extctx`, is the address of an extended context structure, `struct __extctx`, as defined in the `sys/context.h` file. A new field, `__extctx_magic`, within the `ucontext` structure indicates whether the extended context information is valid when the value of `__extctx_magic` is equal to `__EXTCTX_MAGIC`. The additional vector machine state for a thread using the vector extension is saved and restored as a member of this new context extension to the `ucontext` structure as a part of signal delivery and return.

The `ucontext` structure is also used on APIs (such as `getcontext()`, `setcontext()`, `swapcontext()`, and `makecontext()`). In these cases, the context needing to be saved is due to a voluntary action, for which calling linkage convention requires only that nonvolatile machine state be saved. Because the default mode of vector enablement on AIX, as described in the ABI section, is to not use nonvolatile vector registers, there are no extensions of the `ucontext` structure required for the majority of applications. If an application chooses to explicitly enable the use of nonvolatile vector registers, it will pick up an extended sized `ucontext` structure that already has space for the `__extctx` field that is included by the implicit definition of `__EXTABI__` by the compiler. The extended `ucontext` can also be picked up by an explicit definition of `__AIXEXTABI`.

Similarly, the `jmp_buf` for use with `setjmp()` or `longjmp()` requires no change for default-mode vector-enabled applications, since nonvolatile vector registers are not used. The explicit enablement of nonvolatile vector registers results in larger `jmp_buf` allocations, due to the implicit definition of

`__EXTABI__` by the compiler. The extended jump buffers can also be activated by explicit definition of `__AIXEXTABI`.

See the `sys/context.h` header file for a more detailed layout of the extended context information.

Vector memory allocation and alignment

Vector data types introduce a data type requiring 16-byte alignment. In accordance with the AltiVec programming interface specification, a set of malloc subroutines (`vec_malloc`, `vec_free`, `vec_realloc`, `vec_calloc`) are provided by AIX that give 16-byte aligned allocations.

Vector-enabled compilation, with `_VEC_` implicitly defined by the compiler, will result in any calls to legacy malloc and calloc being redirected to their vector-safe counterparts, `vec_malloc` and `vec_calloc`, respectively. Non-vector code can also be explicitly compiled to pick up these same malloc and calloc redirections by explicitly defining `__AIXVEC`. The alignment of the default `malloc()`, `realloc()`, and `calloc()` allocations can also be controlled at runtime.

First, externally to any program, a new environment variable, `MALLOCALIGN`, can be set to the default alignment desired for every `malloc()` allocation. An example is

```
MALLOCALIGN=16; export MALLOCALIGN
```

The `MALLOCALIGN` environment variable can be set to any power of 2 greater than or equal to the size of a pointer in the corresponding execution mode (4 bytes for 32-bit mode, 8 bytes for 64-bit mode). If `MALLOCALIGN` is set to an invalid value, then the value is rounded up to the next power of 2, and all subsequent `malloc()` allocations will be aligned to that value.

Also, internally to a program, the program can use a new command option to the `mallopt()` interface to specify the desired alignment for future allocations. An example is

```
rc = mallopt(M_MALIGN, 16);
```

Refer to `mallopt` and `MALLOCALIGN` for more information.

printf and scanf of vector data types

In accordance with the AltiVec programming interface specification, support is added to the AIX versions of `scanf`, `fscanf`, `sscanf`, `wscanf`, `printf`, `fprintf`, `sprintf`, `snprintf`, `wsprintf`, `vprintf`, `vfprintf`, `vsprintf`, and `vwsprintf` for the new vector conversion format strings. The new size formatters are as follows:

- `vl` or `lv` consumes one argument and modifies an existing integer conversion, resulting in vector signed int, vector unsigned int, or vector bool for output conversions or vector signed int * or vector unsigned int * for input conversions. The data is then treated as a series of four 4-byte components, with the subsequent conversion format applied to each.
- `vh` or `hv` consumes one argument and modifies an existing short integer conversion, resulting in vector signed short, or vector unsigned short for output conversions or vector signed short * or vector unsigned short * for input conversions. The data is treated as a series of eight 2-byte components, with the subsequent conversion format applied to each.
- `v` consumes one argument and modifies a 1-byte integer, 1-byte character, or 4-byte floating point conversion. If the conversion is a floating point conversion the result is vector float for output conversion or vector float * for input conversion. The data is treated as a series of four 4-byte floating point components with the subsequent conversion format applied to each. If the conversion is an integer or character conversion, the result is either vector signed char, vector unsigned char, or vector bool char for output conversion or vector signed char * or vector unsigned char * for input conversions. The data is treated as a series of sixteen 1-byte components, with the subsequent conversion format applied to each.

Any conversion format that can be applied to the singular form of a vector-data type can be used with a vector form. The `%d`, `%x`, `%X`, `%u`, `%i`, and `%o` integer conversions can be applied with the `%lv`, `%vl`, `%hv`, `%vh`, and `%v` vector-length qualifiers. The `%c` character conversion can be applied with the `%v` vector

length qualifier. The %a, %A, %e, %E, %f, %F, %g, and %G float conversions can be applied with the %v vector length qualifier.

For input conversions, an optional separator character can be specified excluding whitespace preceding the separator. If no separator is specified, the default separator is a space including whitespace characters preceding the separator, unless the conversion is c, and then the default conversion is null.

For output conversions, an optional separator character can be specified immediately preceding the vector size conversion. If no separator is specified, the default separator is a space, unless the conversion is c, and then the default separator is null.

Threaded applications

Multithreaded applications exploiting the vector extension are also supported. These applications are supported in both system scope (1:1 threading model) and process scope (M:N threading model). If a multithreaded application is compiled with nonvolatile vector registers enabled, the pthreads for that application will be flagged as extended ABI pthreads. The result will be larger context-save buffer allocations within the pthread library for those threads. The dbx AIX debugger also provides full support for machine-level debugging of vector-enabled multithreaded programs.

Compilers

An AIX compiler supporting the vector extension must conform to the AIX ABI Vector Extension. As previously described, the default vector-enabled compilation mode on AIX should be with the use of nonvolatile vector registers disabled. An explicit option to enable the use of nonvolatile vector registers can be provided and enabled at your discretion, after understanding the issues and risks regarding new and old module interoperability.

When enabling the use of nonvolatile vector registers, a C or C++ compiler must predefine **__EXTABI__**. Also, when enabled for any form of vector compilation, a C or C++ compiler is expected to predefine **__VEC__**. If compiling non-vector enabled C or C++ modules for linkage with vector-enabled Fortran modules, it is best that the C or C++ modules be explicitly compiled with at least **__AIXVEC** defined (explicit definition analogous to **__VEC__**), and also **__AIXEXTABI** (explicit definition analogous to **__EXTABI**) if nonvolatile vector registers are enabled in the Fortran modules.

In addition to the Altivec programming interface specification, which provides an explicit extension to the C and C++ languages for vector programming, some compilers will probably allow the exploitation of the vector extension in some optimization settings when targeting a processor that supports the vector extension.

Refer to your compiler documentation for more details.

Assembler

The AIX assembler, in the `/usr/ccs/bin/as` directory, now supports the additional instruction set defined by the vector extension, and specifically as implemented by the PowerPC 970 processor. You can use the new `-m970` assembly mode or the `.machine 970` pseudo op within the source file to enable the assembly of the new vector instructions. Refer to the *Assembler Language Reference* for more information.

Debugger

The dbx AIX debugger, in `/usr/ccs/bin/dbx`, supports machine-level debugging of vector-enabled programs. This support includes the ability to disassemble the new vector instructions, and to display and set vector registers. A new `$instructionset` value of 970 has been defined for enabling disassembly of the PowerPC 970-specific instructions, including the vector instructions, when not running dbx on a PowerPC 970 system. Note that if running dbx on a PowerPC 970, the default `$instructionset` will be 970.

To view vector registers, the subcommand `unset $novregs` must be used, as vector registers are not displayed by default. Also, if the processor does not support the vector extension, or the process or thread being examined is not using the vector extension, then no vector register state will be displayed.

Otherwise, the registers subcommand will print all of the vector registers and their contents in raw hexadecimal.

You can also display the vector registers individually, formatted according to a fundamental type. For instance, print \$vr0 will display the contents of register VR0 as an array of 4 integers. print \$vr0c will display the contents of register VR0 as an array of 16 characters. print \$vr0s will display the contents of register VR0 as an array of 8 shorts, and print \$vr0f will display the contents of register VR0 as an array of 4 floats.

You can assign entire vector registers, for example assign \$vr0 = \$vr1, or assign individual vector elements of the vector register as if assigning an element of an array. For example, assign \$vr0[3] = 0x11223344 sets only the 4th integer member of VR0. assign \$vr0f[0] = 1.123 results in only the first float member of VR0 being set to the value 1.123.

You can trace vector registers throughout the execution of a function or program, for example tracei \$vr0 in main will display the contents of VR0 each time it is modified in main(). Likewise, by specifying one of the format registers (\$vr0f, \$vr0c, \$vr0s) to tracei, each display of the contents will be formatted accordingly.

As long as compilers represent vector data types as arrays of their fundamental types, dbx should also be able to display a vector data type formatted as an array.

Refer to the **dbx** command documentation for more information.

Enablement for third-party debuggers is also provided in the form of PTT_READ_VEC and PTT_WRITE_VEC new ptrace operations, for reading or writing vector register state for a thread. Refer to the ptrace documentation for details.

The **/proc** filesystem is also enhanced to support a **/proc**-based debugger. The status and lwpstatus files for a vector-enabled process and thread, respectively, are extended to include vector register state. A new control message, PCSVREG, is supported on the write of a process or thread control file for setting vector register state. Refer to the **/proc** File Reference for more details.

Core files

AIX also supports the inclusion of vector machine state as part of the core file for a vector-enabled process or thread. Only if a process or thread is using the vector extension will the vector machine state be included in the core image for that thread. Note that if you select the pre-AIX 4.3 core file format, vector state will not be included. The vector state is only supported in the current core file formats. You can use the dbx command to read and display the vector machine state of a vector-enabled core file.

System memory allocation using the malloc subsystem

Memory is allocated to applications using the **malloc** subsystem.

The **malloc** subsystem is a memory management API that consists of the following subroutines:

- **malloc**
- **calloc**
- **realloc**
- **free**
- **mallopt**
- **mallinfo**
- **alloca**
- **valloc**
- **posix_memalign**

The **malloc** subsystem manages a logical memory object called a *heap*. The heap is a region of memory that resides in the application's address space between the last byte of data allocated by the compiler

and the end of the data region. The heap is the memory object from which memory is allocated and to which memory is returned by the **malloc** subsystem API.

The **malloc** subsystem performs the following fundamental memory operations:

- Allocation:
Performed by the **malloc**, **calloc**, **valloc**, **alloca**, and **posix_memalign** subroutines.
- Deallocation:
Performed by the **free** subroutine.
- Reallocation:
Performed by the **realloc** subroutine.

The **mallopt** and **mallinfo** subroutines are supported for System V compatibility. The **mallinfo** subroutine can be used during program development to obtain information about the heap managed by the **malloc** subroutine. The **mallopt** subroutine can be used to disclaim page-aligned, page-sized free memory, and to enable and disable the default allocator. Similar to the **malloc** subroutine, the **valloc** subroutine is provided for compatibility with the Berkeley Compatibility Library.

For additional information, see the following sections:

Working with the process heap

_edata is a symbol whose address is the first byte following the last byte of initialized program data. The **_edata** symbol refers to the start of the process heap, which is enlarged by the **malloc** subsystem when the first block of data is allocated. The **malloc** subsystem enlarges the process heap by increasing the process **brk** value, which denotes the end of the process heap. This is done by calling the **sbrk** subroutine. The **malloc** subsystem expands the process heap as the needs of the application dictate.

The process heap is divided into *allocated* and *freed* memory blocks. The free pool consists of the memory available for subsequent allocation. An allocation is completed by first removing a memory block from the free pool and then returning to the calling function a pointer to this block. A reallocation is completed by allocating a memory block of the new size, moving the data in the original block to the new block, and freeing the original block. The allocated memory blocks consist of the pieces of the process heap being used by the application. Because the memory blocks are not physically removed from the heap (they change state from free to allocated), the size of the process heap does not decrease when memory is freed by the application.

Process address space in 32-bit applications

A 32-bit application program running on the system has an address space that is divided into the following segments:

Segment	Description
0x00000000 to 0x0fffffff	Contains the kernel.
0x10000000 to 0x1fffffff	Contains the application program text.
0x20000000 to 0x2fffffff	Contains the application program data, the process heap, and the application stack.
0x30000000 to 0xcfffffff	Available for use by shared memory or mmap services.
0xd0000000 to 0xdfffffff	Contains shared library text.
0xe0000000 to 0xffffffff	Available for use by shared memory or mmap services.
0xf0000000 to 0xffffffff	Contains the application shared library data.

Process address space in 64-bit applications

A 64-bit application program running on the system has an address space that is divided into the following segments:

Segment	Description
0x0000 0000 0000 0000 to 0x0000 0000 0fff ffff	Contains the kernel.
0x0000 0000 f000 0000 to 0x0000 0000 ffff ffff	Reserved.
0x0000 0001 0000 0000 to 0x07ff ffff ffff ffff	Contains the application program text, application program data, the process heap, and shared memory or mmap services.
0x0800 0000 0000 0000 to 0x08ff ffff ffff ffff	Privately loaded objects.
0x0900 0000 0000 0000 to 0x09ff ffff ffff ffff	Shared library text and data.
0x0f00 0000 0000 0000 to 0x0fff ffff ffff ffff	Application stack.

Note: AIX uses a delayed paging slot allocation technique for storage allocated to applications. When storage is allocated to an application with a subroutine, such as **malloc**, no paging space is assigned to that storage until the storage is referenced. This technique is useful for applications that allocate large sparse memory segments. However, this technique can affect portability of applications that allocate very large amounts of memory. If the application expects that calls to **malloc** will fail when there is not enough backing storage to support the memory request, the application might allocate too much memory. When this memory is referenced later, the machine quickly runs out of paging space and the operating system kills processes so that the system is not completely exhausted of virtual memory. The application that allocates memory must ensure that backing storage exists for the storage being allocated. Setting the **PSALLOC** environment variable to **PSALLOC=early** changes the paging space allocation technique to an early allocation algorithm. In early allocation, paging space is assigned once the memory is requested. For more information, see [Paging space and virtual memory in Operating system and device management](#).

Understanding system allocation policy

The *allocation policy* refers to the set of data structures and algorithms employed to represent the heap and to implement allocation, deallocation, and reallocation. The **malloc** subsystem supports several different allocation policies, including the default allocation policy, the watson allocation policy, the malloc 3.1 allocation policy, and the user-defined allocation policy. The API for accessing the **malloc** subsystem is identical for all allocation policies; only the underlying implementation is different.

You can use the following environment variables to specify the allocation policy and any normal or debug options for that policy:

- **MALLOCTYPE** specifies the allocation policy.
- **MALLOCOPTIONS** specifies normal options to the chosen allocation policy.
- **MALLOCDEBUG** specifies debug options to the chosen allocation policy.
- **MALLOCALIGN** specifies the default **malloc** alignment external to a program.

The default allocation policy is generally more efficient and is the preferred choice for the majority of applications. The other allocation policies have some unique behavioral characteristics that can be beneficial in specific circumstances, as described in [Comparing the Different Allocation Policies](#).

Some options to the various allocation policies are compatible with each other and can be used in tandem. When you are using options in tandem, use a comma (,) to separate options specified by the **MALLOCOPTIONS** and **MALLOCDEBUG** environment variables.

The **MALLOCALIGN** environment variable can be set to the default alignment desired for every **malloc()** allocation. An example is

```
MALLOCALIGN=16; export MALLOCALIGN
```

The **MALLOCALIGN** environment variable can be set to any power of 2 value greater than or equal to the size of a pointer in the corresponding run mode (4 bytes for 32-bit mode, 8 bytes for 64-bit mode). For 32-bit vector-enabled programs, this environment variable can be set to 16, so all **malloc()**s will be suitably aligned for vector data types if necessary. Note that 64-bit vector programs will already receive 16-byte-aligned allocations.

Also, internally to a program, the program can use the **mallopt(M_MALIGN, 16)** routine to change the default **malloc()** to provide 16-byte aligned allocations. The **mallopt(M_MALIGN)** routine allows a program to control the default malloc alignment dynamically at runtime.

Understanding the default allocation policy

The default allocation policy maintains the free space in the heap as nodes in a *cartesian* binary search tree in which the nodes are ordered left-to-right by address (increasing address to the right) and top-to-bottom by length (such that no child is larger than its parent). This data structure imposes no limitation on the number of block sizes supported by the tree, allowing a wide range of potential block sizes. Tree-reorganization techniques optimize access times for node location, insertion, and deletion, and also protect against fragmentation.

The default allocation policy provides support for the following optional capabilities:

Allocation

A small amount of overhead is required to service an allocation request. This is due to the need for a metadata prefix and the need for suitable alignment of each block of memory. The size of the metadata prefix for all allocations is 8 and 16 bytes for 32-bit and 64-bit programs respectively. Each block must be aligned on a 16 or 32 byte boundary, thus the total amount of memory required for an allocation of size *n* is:

```
size = roundup(n + prefix_size, alignment requirement)
```

For example, an allocation of size 37 in a 32-bit process would require `roundup(37 + 8, 16)`, which is equal to 48 bytes.

The node of the tree with the lowest address that is greater than or equal to the size required is removed from the tree. If the block found is larger than the needed size, the block is divided into two blocks: one of the needed size, and the second a remainder. The second block, called the *runt*, is returned to the free tree for future allocation. The first block is returned to the caller.

If a block of sufficient size is not found in the free tree, the heap is expanded, a block the size of the acquired extension is added to the free tree, and allocation continues as previously described.

Deallocation

Memory blocks deallocated with the **free** subroutine are returned to the tree, at the root. Each node along the path to the insertion point for the new node is examined to see if it adjoins the node being inserted. If it does, the two nodes are merged and the newly merged node is relocated in the tree. If no adjoining block is found, the node is simply inserted at the appropriate place in the tree. Merging adjacent blocks can significantly reduce heap fragmentation.

Reallocation

If the size of the reallocated block will be larger than the original block, the original block is returned to the free tree with the **free** subroutine so that any possible coalescence can occur. A new block of the requested size is then allocated, the data is moved from the original block to the new block, and the new block is returned to the caller.

If the size of the reallocated block is smaller than the original block, the block is split and the smaller one is returned to the free tree.

Limitations

The default allocation policy supports the following options:

- [Malloc Multiheap](#)

- [Malloc Buckets](#)
- [Malloc Disclaim](#)
- [Malloc Thread Cache](#)
- [Understanding the no_overwrite Option](#)

Understanding the watson allocation policy

The Watson allocation policy maintains the free space in the heap as nodes in two separate red-black trees: one sorted by address, the other by size. Red-black trees provide simpler and more efficient tree operations than the cartesian trees of the default allocator, thus the watson allocation policy is often faster than the default.

Allocation

The Watson allocation policy has the same overhead requirements as the default allocation policy.

The size tree is searched for the smallest possible block that is greater than or equal to the size required. This block is then removed from the size tree. If the block found is larger than the needed size, the block is divided into two blocks: one block of the remaining size, and the second of the required size. The first block, called the *runt*, is returned to the size tree for future allocation. The second block is returned to the caller. If the block found in the size tree was exactly the required size, the block is removed from both the size and the address tree, and then returned to the caller.

If a block of sufficient size is not found in the free tree, the process heap is expanded, a block the size of this expansion is added to the size and address trees, and allocation continues as previously described.

Deallocation

Memory blocks deallocated with the **free** subroutine are returned to the address tree at the root. Each node along the path to the insertion point for the new node is examined to see if it adjoins the node being inserted. If it does, the two nodes are merged and the newly merged node is relocated in the size tree. If no adjoining block is found, the node is simply inserted at the appropriate place in the both the address and size tree.

After insertion, both red-black trees must be checked for correct balancing.

Reallocation

If the size of the reallocated block will be larger than the original block, the original block is returned to the free trees with the **free** subroutine so that any possible coalescence can occur. A new block of the requested size is then allocated, the data is moved from the original block to the new block, and the new block is returned to the caller.

If the size of the reallocated block is smaller than the original block, the block is split and the remainder is returned to the free tree.

Limitations

The Watson allocation policy supports the following options:

- [Malloc Multiheap](#)
- [Malloc Disclaim](#)
- [Malloc Thread Cache](#)
- [Understanding the no_overwrite Option](#)

Understanding the malloc 3.1 allocation policy

The malloc 3.1 allocation policy can be selected by setting `MALLOCTYPE=3.1` prior to process startup. Thereafter, all 32-bit programs run by the shell will use the malloc 3.1 allocation policy (64-bit programs will continue to use the default allocation policy).

The malloc 3.1 allocation policy maintains the heap as a set of 28 hash buckets, each of which points to a linked list. Each linked list contains blocks of a particular size. The index into the hash buckets indicates the size of the blocks in the linked list. The size of the block is calculated using the following formula:

$$\text{size} = 2^{i+4}$$

where i identifies the bucket. This means that the blocks in the list anchored by bucket zero are $20+4 = 16$ bytes long. Therefore, given that a prefix is 8 bytes in size, these blocks can satisfy requests for blocks between 0 and 8 bytes long. The following table illustrates how requested sizes are distributed among the buckets.

Note: This algorithm can use as much as twice the amount of memory actually requested by the application. An extra page is required for buckets larger than 4096 bytes because objects of a page in size or larger are page-aligned. Because the prefix immediately precedes the block, an entire page is required solely for the prefix.

Bucket	Block Size	Sizes Mapped	Pages Used
0	16	0 ... 8	
1	32	9 ... 24	
2	64	25 ... 56	
3	128	57 ... 120	
4	256	121 ... 248	
5	512	249 ... 504	
6	1K	505 ... 1K-8	
7	2K	1K-7 ... 2K-8	
8	4K	2K-7 ... 4K-8	2
9	8K	4K-7 ... 8K-8	3
10	16K	8K-7 ... 16K-8	5
11	32K	16K-7 ... 32K-8	9
12	64K	32K-7 ... 64K-8	17
13	128K	64K-7 ... 128K-8	33
14	256K	128K-7 ... 256K-8	65
15	512K	256K-7 ... 512K-8	129
16	1M	256K-7 ... 1M-8	257
17	2M	1M-7 ... 2M-8	513
18	4M	2M-7 ... 4M-8	1K + 1
19	8M	4M-7 ... 8M-8	2K + 1
20	16M	8M-7 ... 16M-8	4K + 1
21	32M	16M-7 ... 32M-8	8K + 1
22	64M	32M-7 ... 64M-8	16K + 1
23	128M	64M-7 ... 128M-8	32K + 1
24	256M	128M-7 ... 256M-8	64K + 1
25	512M	256M-7 ... 512M-8	128K + 1
26	1024M	512M-7 ... 1024M-8	256K + 1
27	2048M	1024M-7 ... 2048M-8	512K + 1

Allocation

A block is allocated from the free pool by first converting the requested bytes to an index in the bucket array, using the following equation:

```
needed = requested + 8

If needed <= 16,
then
bucket = 0

If needed > 16,
then
bucket = (log(needed)/log(2) rounded down to the nearest integer) - 3
```

The size of each block in the list anchored by the bucket is $\text{block size} = 2^{\text{bucket} + 4}$. If the list in the bucket is null, memory is allocated using the **sbrk** subroutine to add blocks to the list. If the block size is less than a page, then a page is allocated using the **sbrk** subroutine, and the number of blocks arrived at by dividing the block size into the page size are added to the list. If the block size is equal to or greater than a page, needed memory is allocated using the **sbrk** subroutine, and a single block is added to the free list for the bucket. If the free list is not empty, the block at the head of the list is returned to the caller. The next block on the list then becomes the new head.

Deallocation

When a block of memory is returned to the free pool, the bucket index is calculated as with allocation. The block to be freed is then added to the head of the free list for the bucket.

Reallocation

When a block of memory is reallocated, the needed size is compared against the existing size of the block. Because of the wide variance in sizes handled by a single bucket, the new block size often maps to the same bucket as the original block size. In these cases, the length of the prefix is updated to reflect the new size and the same block is returned. If the needed size is greater than the existing block, the block is freed, a new block is allocated from the new bucket, and the data is moved from the old block to the new block.

Limitations

Setting `MALLOCTYPE=3.1` will only enable the malloc 3.1 policy for 32-bit programs. For 64-bit programs to use the malloc 3.1 policy, the **MALLOCTYPE** environment variable must be explicitly set to `MALLOCTYPE=3.1_64BIT`. This allocation policy is less efficient than the default and is not recommended for use in most cases.

The malloc 3.1 allocation policy supports the following options:

- [Malloc Disclaim](#)
- [Understanding the no_overwrite Option](#)

Understanding the pool allocation policy

Malloc pool is a high performance front end to the libc functions **malloc**, **calloc**, **free**, **posix_memalign** and **realloc** for managing storage objects smaller than 513 bytes. The performance benefits derive from dramatically shorter path lengths and better data cache utilization. For multithreaded applications, there is the further benefit that thread local pool anchors are used to avoid atomic operations. This front end can be used in conjunction with any of the storage management schemes currently provided in libc (*yorktown*, and *watson*).

To use **malloc** pool, run the following command:

```
export MALLOCOPTIONS=pool<:max_size>
```

When this option is specified, a collection of pools is created during **malloc** initialization where each pool is a linked list of fixed sized objects. The smallest pool can hold objects of pointer-size (such as 8 bytes for 32-bit applications or 16 bytes for 64-bit applications). Each successive pool can accommodate objects whose size is pointer-size larger than the previous pool. This means there are 128 pools for 32-bit

applications and 64 pools for 64-bit applications. The collection of pools is represented as an array of pointers that "anchor" the linked lists.

Malloc pool uses its own memory, the pool heap, which is not shared with standard **malloc**. When specified, the **max_size** option is rounded up to the next higher 2 MB value and is used to control the size of the pool heap. The **max_size** option can be specified as a decimal number or a hexadecimal number preceded by 0x or 0X (for example, `export MALLOCOPTIONS=pool1:0x17000000` will set **max_size** to 24 MB after rounding up).

For 32-bit applications, the size of the pool heap starts at 2 MB. If more storage is required and the total pool heap storage is less than **max_size**, an additional 2 MB is acquired. Each 2 MB area will be on a 2 MB boundary, but need not be contiguous to any of the other 2 MB areas. For 64-bit applications, a single contiguous pool heap of **max_size** is allocated during **malloc** initialization and never extended. If **max_size** is not specified, it defaults to 512 MB for 32-bit applications and 32 MB for 64-bit applications. For both 32- and 64-bit modes, **max_size** will be set to 512 MB if a larger size is specified. For 32-bit mode, the **max_size** is set to 512MB, and for 64-bit mode, the **max_size** is set to 3.7 GB if a larger size is specified.

Storage Utilization

All pool anchors are initially set to NULL or empty. When **malloc** pool services a request and the corresponding pool is empty, a routine is called that allocates storage from the pool heap in contiguous 1024-byte chunks on 1024-byte boundaries. Multiple objects of the requested size are "created". The address of the first is returned to satisfy the request, while the remaining objects are linked together and placed on the pool anchor. For each 1024-byte chunk, there is a 2-byte entry in an auxiliary table that is used by free to determine the size of a returned object.

When an object is freed by **malloc** pool, it is merely "pushed" onto the appropriate pool anchor. No attempt is made to coalesce blocks to create larger sized objects.

Because of this behavior, **malloc** pool may use more storage than other forms of **malloc**.

Alignment

The default alignment for the **malloc()**, **calloc()**, and **realloc()** subroutines must be specified by setting the **MALLOCALIGN** environment variable appropriately. The **posix_memalign()** subroutine continues working even if the **MALLOCALIGN** environment variable is not set. If **MALLOCALIGN** is greater than 512, **malloc** pool is not used.

Cache efficiency

The memory objects allocated with **malloc pool** have no prefixes or suffixes. Data cache lines are therefore more densely packed with application usable data. As all memory objects that are a power of 2 in size are aligned on a boundary equal to that size, each object is contained within the minimal number of cache lines. The **malloc** and **free** subroutines do not scan trees or linked lists and therefore do not "pollute" the cache.

Multithreaded support

Malloc pools can improve performance significantly in a multithreaded scenario because it reduces lock contention and the need of atomic operations.

Load balancing support

In some multithreaded scenarios, one-thread's free pool might grow very large due to repetitive freeing of dynamically allocated memory. However, other threads may not be able to use this memory.

Load-balancing support causes a thread to release half the memory in each pool to a global pool after the pool reaches a threshold value so that other threads can use it. You can tune the threshold values at which a thread's pool will be readjusted.

To turn on the load-balancing support, the following options must be exported:

```
export MALLOCOPTIONS=pool1:0x80000000,pool_balanced
```

```
export MALLOCFREEPOOL=min_size<-max_size>:threshold_value<,-min_size<-max_size>:  
threshold_value, ... >,default:threshold
```

The following example sets the threshold value for the pools that provide a memory of 0 -16 bytes and 256 chunks, and the threshold value of the pool that serves 32-byte chunks to 512-byte chunks. For the rest of the pools, 128-byte chunks is the threshold value.

```
export MALLOCFREEPOOL=0-16:256,32:512,default:128
```

Debugging support

There is no debug version of this high-performance front end. If the **MALLOCDEBUG** environment variable is set, the pool option is ignored. It is expected that applications will be debugged using "normal" **malloc** prior to activating pooling.

Understanding the user-defined allocation policy

The malloc subsystem provides a mechanism through which users may develop their own algorithms for managing the system heap and allocating memory.

Understanding the no_overwrite option

An additional option available to all allocation policies is `no_overwrite`. To reduce the overhead of glink code within the malloc subsystem, the function descriptor for the **malloc** subsystem APIs are overwritten with the function descriptor for the actual underlying implementation. Because some programs, such as third-party debuggers, might not work when function pointers are modified in this manner, the `no_overwrite` option can be used to disable this optimization.

To disable this optimization, set `MALLOCOPTIONS=no_overwrite` prior to process startup.

Comparing the various allocation policies

The various malloc allocation policies elaborated above provide flexibility to application developers when used separately or combined in supported ways. It is the responsibility of the developer to recognize the unique needs of an application and to tune the various allocation policy parameters in a beneficial way.

Comparing the default and malloc 3.1 allocation policies

Because the malloc 3.1 allocation policy rounds up the size of each allocation request to the next power of 2, it can produce considerable virtual- and real-memory fragmentation and poor locality of reference. The default allocation policy is generally a better choice because it allocates exactly the amount of space requested and is more efficient about reclaiming previously used blocks of memory.

Unfortunately, some application programs may depend inadvertently on side effects of the malloc 3.1 allocation policy for acceptable performance or even for correct functioning. For example, a program that overruns the end of an array may function correctly when using the malloc 3.1 allocator only because of the additional space provided by the rounding-up process. The same program is likely to experience erratic behavior or even fail when used with default allocator because the default allocator allocates only the number of bytes requested.

As another example, because of the inefficient space reclamation of the malloc 3.1 allocation algorithm, the application program almost always receives space that has been set to zeros (when a process touches a given page in its working segment for the first time, that page is set to zeros). Applications may depend on this side effect for correct execution. In fact, zeroing out of the allocated space is not a specified function of the **malloc** subroutine and would result in an unnecessary performance penalty for programs that initialize only as required and possibly not to zeros. Because the default allocator is more aggressive about reusing space, programs that are dependent on receiving zeroed storage from **malloc** will probably fail when the default allocator is used.

Similarly, if a program continually reallocates a structure to a slightly greater size, the malloc 3.1 allocator may not need to move the structure very often. In many cases, the **realloc** subroutine can make use of the extra space provided by the rounding implicit in the malloc 3.1 allocation algorithm. The default allocator will usually have to move the structure to a slightly larger area because of the likelihood that something else has been called by the **malloc** subroutine just above it. This may present the appearance of a degradation in **realloc** subroutine performance when the default allocator is used instead of the

malloc 3.1 allocator. In reality, it is the surfacing of a cost that is implicit in the application program's structure.

Debugging application mismanagement of the system heap

The **malloc** subsystem offers a collection of debugging tools intended to help the application developer debug and correct errors in a program's heap management. These debugging tools are controlled through the **MALLOCDEBUG** environment variable.

Synopsis of malloc environment variable and options

The following table shows the compatibility between the **MALLOCTYPE** and **MALLOCOPTIONS** environment variables.

*Table 84. Compatibility between **MALLOCTYPE** and **MALLOCOPTIONS** Environment Variables*

	multiheap (and sub- options)	buckets (and suboptio ns)	Thread Cache	disclaim	no_overwrite
Default Allocator	yes	yes	yes	yes	yes
3.1	no	no	yes	yes	yes
Watson	no	no	no	no	no
Watson2	no	no	no	no	no
user:	no	no	no	no	yes

*Table 85. Compatibility between **MALLOCDEBUG** and **MALLOCTYPE** Environment Variables*

	York Town<Default Allocator>	3.1	Watson	Watson2	user:
catch_overflow (and suboptions)	yes	no	yes	yes	no
report_allocations	yes	no	yes	yes	no
postfree_checking	yes	no	yes	yes	no
validate_ptrs	yes	no	yes	yes	no
trace	yes	no	yes	yes	no
log	yes	no	yes	yes	no
verbose	no	no	no	no	no

All **MALLOCDEBUG** options are compatible and supported with **MALLOCOPTIONS**.

Understanding the Watson2 allocation policy

The Watson2 malloc subsystem adapts to the behavior of the application when it changes from a single thread to multiple threads and from multiple threads to a single thread. It uses a thread-specific mechanism that uses a varying number of heap structures, which depend on the behavior of the program. Therefore no configuration options are required. The Watson2 malloc subsystem has O(logN) amortized the cost per operation for many workloads because vast number of operations can be run at a constant time without synchronization.

Allocation

Allocation is handled through a combination of mechanisms. These mechanisms depend on parameters, such as the number of active threads, size of the request and the deallocation history of the process. The

set of mechanisms reaches from a thread-specific caching and use a variable number of heaps, which has thread affinity to a Double-red-black tree and Page-based coalescing.

Deallocation

Deallocation depends on the same parameters as the allocation behavior. Typically, a returning block is captured in the thread-specific cache. Based on the heap affinity and capacity utilization, memory might be returned to one of the multiple heap structures. At times, the content from the multiple heap structure is consolidated to a common heap structure to improve coalescing and to reduce heap fragmentation. To improve robustness against application errors, the allocator identifies deallocation of invalid pointers or corrupted blocks, to a certain degree and filters these operations.

Reallocation

Large blocks of memory that are adequate are reused. If the current block cannot satisfy the request, it is replaced with a regular deallocation and allocation.

Limitations

The Watson2 malloc subsystem is adaptive to the application and does not require any further options, but the Watson2 malloc subsystem supports the following debug features that are controlled by the MALLOCDEBUG variable: `validate_ptrs`, `report_allocations`, and `trace`. The reports that are related to allocations can be redirected to a file by using the `output:<filename>` option. See “[Debug malloc tool](#)” on [page 623](#) for detailed information on the MALLOCDEBUG variable.

Related concepts

[Program address space overview](#)

The Base Operating System provides a number of services for programming application program memory use.

[User-defined malloc replacement](#)

Users can replace the memory subsystem (**malloc**, **calloc**, **realloc**, **free**, **mallopt** and **mallinfo** subroutines) with one of their own design.

[Malloc multiheap](#)

By default, the malloc subsystem uses a single heap, or free memory pool.

[Malloc buckets](#)

Malloc buckets provides an optional buckets-based extension of the default allocator.

[Malloc trace](#)

Malloc Trace provides an optional extension to the malloc subsystem for use with the trace facility.

[Malloc log](#)

Malloc Log is an optional extension of the **malloc** subsystem, allowing the user to obtain information regarding the active allocations currently held by the calling process. This data can then be used in problem determination and performance analysis.

[Malloc disclaim](#)

Malloc Disclaim is an optional extension of the malloc subsystem, providing the user with a means to enable the automatic disclaiming of memory returned by the **free** subroutine.

Related reference

[Debug malloc tool](#)

Debugging applications that are mismanaging memory allocated by the **malloc** subsystem can be difficult and tedious. This is because there is generally no synchronicity between the insertion of an error and the exposure of its resulting symptom.

User-defined malloc replacement

Users can replace the memory subsystem (**malloc**, **calloc**, **realloc**, **free**, **mallopt** and **mallinfo** subroutines) with one of their own design.

Note: Replacement Memory Subsystems written in C++ are not supported due to the use of the **libc.a** memory subsystem in the C++ library **libC.a**.

The existing memory subsystem works for both threaded and non-threaded applications. The user-defined memory subsystem must be threadsafe so that it works in both threaded and non-threaded processes. Because there are no checks to verify that it is, if a non-threadsafe memory module is loaded in a threaded application, memory and data may be corrupted.

The user defined memory subsystem 32- and 64- bit objects must be placed in an archive with the 32-bit shared object named `mem32.o` and the 64-bit shared object named `mem64.o`.

The user-shared objects must export the following symbols :

- `__malloc__`
- `__free__`
- `__realloc__`
- `__calloc__`
- `__mallinfo__`
- `__mallopt__`
- `__malloc_init__`
- `__malloc_prefork_lock__`
- `__malloc_postfork_unlock__`

The user-shared objects can optionally export the following symbol:

- `__malloc_start__`
- `__posix_memalign__`

Execution does not stop if these symbols do not exist.

The functions are defined as follows:

void *__malloc__(size_t):

This function is the user equivalent of the **malloc** subroutine.

void __free__(void *):

This function is the user equivalent of the **free** subroutine.

void *__realloc__(void *, size_t):

This function is the user equivalent of the **realloc** subroutine.

void *__calloc__(size_t, size_t):

This function is the user equivalent of the **calloc** subroutine.

int __mallopt__(int, int):

This function is the user equivalent of the **mallopt** subroutine.

struct mallinfo __mallinfo__():

This function is the user equivalent of the **mallinfo** subroutine.

void __malloc_start__()

This function will be called once before any other user-defined malloc entry point is called.

void __posix_memalign__()

This function is the user equivalent of the **posix_memalign** subroutine. If this symbol does not exist, the execution will not stop, but a call made to the `posix_memalign` subroutine will cause unexpected results.

The following functions are used by the thread subsystem to manage the user-defined memory subsystem in a multithreaded environment. They are only called if the application and/or the user defined module are bound with **libpthreads.a**. Even if the the user-defined subsystem is not threadsafe and not bound with **libpthreads.a**, these functions must be defined and exported. Otherwise, the object will not be loaded.

void __malloc_init__(void)

Called by the pthread initialization routine. This function is used to initialize the threaded-user memory subsystem. In most cases, this includes creating and initializing some form of locking data.

Even if the user-defined memory subsystem module is bound with **libpthread.a**, the user-defined memory subsystem *must* work before `__malloc_init__()` is called.

void __malloc_prefork_lock__(void)

Called by pthreads when the `fork` subroutine is called. This function is used to insure that the memory subsystem is in a known state before the `fork()` and stays that way until the `fork()` has returned. In most cases this includes acquiring the memory subsystem locks.

void __malloc_postfork_unlock__(void)

Called by pthreads when the `fork` subroutine is called. This function is used to make the memory subsystem available in the parent and child after a `fork`. This should undo the work done by `__malloc_prefork_lock__`. In most cases, this includes releasing the memory subsystem locks.

All of the functions must be exported from a shared module. Separate modules must exist for 32- and 64-bit implementations placed in an archive. For example:

- **mem.exp** module:

```
__malloc__
__free__
__realloc__
__calloc__
__mallopt__
__mallinfo__
__malloc_init__
__malloc_prefork_lock__
__malloc_postfork_unlock__
__malloc_start__
```

- **mem_functions32.o** module:

Contains all of the required 32-bit functions

- **mem_functions64.o** module:

Contains all of the required 64-bit functions

The following examples are for creating the shared objects. The `-lpthreads` parameter is needed only if the object uses pthread functions.

- Creating 32-bit shared object:

```
ld -b32 -m -o mem32.o mem_functions32.o \
-bE:mem.exp \
-bM:SRE -lpthreads -lc
```

- Creating 64-bit shared object:

```
ld -b64 -m -o mem64.o mem_functions64.o \
-bE:mem.exp \
-bM:SRE -lpthreads -lc
```

- Creating the archive (the shared objects name must be `mem32.o` for the 32bit object and `mem64.o` for the 64bit object):

```
ar -X32_64 -r archive_name mem32.o mem64.o
```

Enabling the user-defined memory subsystem

The user-defined memory subsystem can be enabled by using one of the following:

- The **MALLOCTYPE** environment variable
- The `_malloc_user_defined_name` global variable in the user's application

To use the **MALLOCTYPE** environment variable, the archive containing the user defined memory subsystem is specified by setting **MALLOCTYPE** to `user:archive_name` where `archive_name` is in the application's `libpath` or the path is specified in the **LIBPATH** environment variable.

To use the `_malloc_user_defined_name` global variable, the user's application must declare the global variable as:

```
char *_malloc_user_defined_name="archive_name"
```

where `archive_name` must be in the application's libpath or a path specified in the **LIBPATH** environment variable.

Note:

1. When a `setuid` application is run, the **LIBPATH** environment variable is ignored so the archive must be in the application's libpath.
2. `archive_name` cannot contain path information.
3. When both the **MALLOCTYPE** environment variable and the `_malloc_user_defined_name` global variable are used to specify the `archive_name`, the archive specified by **MALLOCTYPE** will override the one specified by `_malloc_user_defined_name`.

32-bit and 64-bit considerations

If the archive does not contain both the 32-bit and 64-bit shared objects and the user-defined memory subsystem was enabled using the **MALLOCTYPE** environment variable, there will be problems executing 64-bit processes from 32-bit applications and 32-bit processes from 64-bit applications. When a new process is created using the `exec` subroutine, the process inherits the environment of the calling application. This means that the **MALLOCTYPE** environment variable will be inherited and the new process will attempt to load the user-defined memory subsystem. If the archive member does not exist for this type of program, the load will fail and the new process will exit.

Thread considerations

All of the provided functions must work in a multithreaded environment. Even if the module is linked with `libpthread.a`, at least `__malloc__()` must work before `__malloc_init__()` is called and `pthread` is initialized. This is required because the `pthread` initialization requires `malloc()` before `__malloc_init__()` is called.

All provided memory functions must work in both threaded and non-threaded environments. The `__malloc__()` function should be able to run to completion without having any dependencies on `__malloc_init__()` (that is, `__malloc__()` should initially assume that `__malloc_init__()` has not yet run.) After `__malloc_init__()` has completed, `__malloc__()` can rely on any work done by `__malloc_init__()`. This is required because the `pthread` initialization uses `malloc()` before `__malloc_init__()` is called.

The following variables are provided to prevent unneeded thread-related routines from being called:

- The `__multi_threaded` variable is zero until a thread is created when it becomes non-zero and will not be reset to zero for that process.
- The `__n_pthreads` variable is -1 until `pthread` has been initialized when it is set to 1. From that point on it is a count of the number of active threads.

Example:

If `__malloc__()` uses `pthread_mutex_lock()`, the code might look similar to the following:

```
if (__multi_threaded)
pthread_mutex_lock(mutexptr);

/* ..... work ..... */

if (__multi_threaded)
pthread_mutex_unlock(mutexptr);
```

In this example, `__malloc__()` is prevented from executing `pthread` functions before `pthread` is fully initialized. Single-threaded applications are also accelerated because locking is not done until a second thread is started.

Limitations

Memory subsystems written in C++ are not supported due to initialization and the dependencies of **libc.a** and the **libc.a** memory subsystem.

Error messages are not translated because the **setlocale** subroutine uses `malloc()` to initialize the locales. If `malloc()` fails then the **setlocale** subroutine cannot finish and the application is still in the POSIX locale. Therefore, only the default English messages will be displayed.

Existing statically built programs cannot use the user-defined memory subsystem without recompiling.

Error reporting

The first time the **malloc** subroutine is called, the 32- or 64-bit object in the archive specified by the **MALLOCTYPE** environment variable is loaded. If the load fails, a message displays and the application exits. If the load is successful, an attempt is made to verify that all of the required symbols are present. If any symbols are missing, the application is terminated and the list of missing symbols displays.

Related concepts

System memory allocation using the malloc subsystem

Memory is allocated to applications using the **malloc** subsystem.

Debug malloc tool

Debugging applications that are mismanaging memory allocated by the **malloc** subsystem can be difficult and tedious. This is because there is generally no synchronicity between the insertion of an error and the exposure of its resulting symptom.

Adding to the difficulty is the inherent complexity of memory allocation, with thousands of allocations being made, undone, and accessed (perhaps) asynchronously and simultaneously, all within a multithreaded context that necessitates robust and efficient synchronization.

It is for these reasons that the focus of our debugging tools is primarily to move the time of symptom detection closer to the time of error insertion. This helps the application developer to pinpoint more precisely which section of code is responsible for committing the error.

Many different debugging tools have been developed for use with malloc. Some can be used in combination with other debugging tools and with all allocation policies; others are more limited in their use. Many of the debugging tools consume resources additional to those required by the process. It is up to the application developer to provide adequate resources when necessary.

Performance Considerations

The debug malloc tools are not appropriate for full-time, constant, or system-wide use. Although they are designed for minimal performance impact upon the application being debugged, significant negative impact on overall system throughput can result if they are used widely throughout a system. In particular, setting `MALLOCDEBUG=catch_overflow` in the **/etc/environment** file is not recommended, and will likely cause significant system problems, such as excessive use of paging space. The debug malloc tools should only be used to debug single applications or small groups of applications at the same time.

Because of the extra work involved in making various run-time checks, **malloc** subsystem performance will degrade variably with debug malloc tools enabled (depending on which tool is being used), but not to the point that applications will become unusable. After the problem is resolved, the debug malloc tools should be turned off to restore **malloc** subsystem performance.

Disk and memory considerations

With the `catch_overflow` or Malloc Log tools enabled, the **malloc** subsystem will consume significantly more memory.

For `catch_overflow`, each malloc request is increased by $4096 + 2$ times the size of unsigned long, then rounded up to the next multiple of the `PAGESIZE` macro. `catch_overflow` might prove to be too memory-intensive to use for very large applications, but for the majority of applications that need memory debugging, the extra use of memory should not cause a problem. For large applications, the use

of the **debug_range** and **functionset** options to `catch_overflow` can significantly lower memory usage, allowing the program to be debugged piecemeal.

For Malloc Log, an allocation record is stored for every active allocation in the process. This memory overhead can be minimized by specifying a low number of saved stack pointers.

If the application being debugged frequently calls **malloc** subsystem allocation routines, it might encounter memory usage problems with debug malloc tools enabled that could prevent the application from executing properly in a single segment. If this occurs, it may be helpful to enable the application to access additional memory by using the **ulimit** command and the **-bmaxdata** option of the **ld** command.

For the purpose of running with debug malloc tools enabled, set the `ulimit` for both the data (`-d`) and stack (`-s`) variables as follows:

```
ulimit -d unlimited
ulimit -s unlimited
```

To reserve the maximum of 8 segments for a 32-bit process, the **-bmaxdata** option should be specified as `-bmaxdata:0x80000000`.

When the debug malloc tools are turned off, the default values for `ulimit` and `-bmaxdata` can be restored.

For more information about the **ulimit** command and the **-bmaxdata** option, see [Large Program Support](#).

The debug malloc tools are not appropriate for use in some debugging situations. Because some of the debug malloc tools require the overhead of a page or more per allocation, programs that issue many small allocation requests will see their memory usage increase dramatically. These programs might encounter new failures as memory allocation requests are denied due to a lack of memory or paging space. These failures are not necessarily errors in the program being debugged, and they are not errors in the debug malloc tool.

One specific example of this is the X server, which issues numerous tiny allocation requests during its initialization and operation. Any attempt to run the X server using the `X` or `xinit` commands with `catch_overflow` enabled will result in the failure of the X server due to a lack of available memory. It is possible, however, to debug X in piecemeal fashion using the `debug_range` or `functionset` options. X clients in general will not encounter functional problems running with `catch_overflow` enabled. To use `catch_overflow` on an X client program, take the following steps:

1. Start the X server with `catch_overflow` turned off.
2. Start a terminal window (for example, `dtterm`, `xterm`, `aixterm`).
3. Set the appropriate environment variables within the terminal window session to enable `catch_overflow`.
4. Invoke the X client program to be debugged from within the same window.

Enabling debug malloc

Debug Malloc is not enabled by default, but is enabled and configured by setting the **MALLOCDEBUG** environment variable to the appropriate option. If more than one option is required, options can be separated by a comma (,). Options requested in tandem must be compatible with each other.

Note: To disable the debug malloc, unset the **MALLOCDEBUG** environment variable by using the **unset MALLOCDEBUG** command.

Malloc debugging tools

The following malloc debugging tools are available:

- [Buffer Overflow Detection](#)
 - [align](#)
 - [override_signal_handling](#)
 - [debug_range](#)

- [functionset](#)
- [allow_overreading](#)
- [postfree_checking](#)
- [Malloc Trace](#)
- [Malloc Log](#)
 - [report_allocations](#)
 - [validate_ptrs](#)
- [Malloc Detect](#)
 - [verbose](#)
 - [checkarena](#)
 - [output](#)
 - [continue](#)
- [Malloc debug Fill](#)

Buffer overflow detection

Memory management errors are sometimes caused by the application program writing past the end of an allocated buffer. Because this often has no immediate consequence, symptoms do not arise until much later when the memory that was overwritten (usually belonging to another allocation) is referenced and no longer contains the data originally stored therein.

The `catch_overflow` debug option exists to allow users to identify memory overwrites, overreads, duplicate frees, and reuse of freed memory allocated by the **malloc** subroutine. Memory problems detected by the `catch_overflow` tool result in an abort call or a segmentation violation (**SIGSEGV**). In most cases, when an error is detected, the application stops immediately and a core file is produced.

The `catch_overflow` option affects the allocations of the following allocation policies and options:

- Default Allocation Policy
- Watson Allocation Policy
- Malloc Multiheap Option
- Malloc Threadcache Option
- Malloc Disclaim Option

The `catch_overflow` debug option is enabled by setting `MALLOCDEBUG=catch_overflow`. This will turn on identification of memory overwrites and overreads.

align

By default, the **malloc** subroutine returns a pointer aligned on a 2-word boundary. This is necessary for standards conformance and for programs which cannot accept unaligned memory accesses (e.g. programs using DCE components). However, due to a quirk in the implementation of the `catch_overflow` option, it is possible for a program to overwrite a buffer by an amount less than the alignment value without being detected by `catch_overflow`. The `align` option can be used to tell the **malloc** subsystem to disregard this default alignment in order to decrease or eliminate the number of bytes by which a buffer can be overwritten without detection. A custom alignment can be specified for any power of two between 0 and 4096 inclusive (e.g. 0,1,2,4,...). The values 0 and 1 are treated as the same, that is, there is no memory alignment; therefore any memory accesses beyond the allocated area will cause a `SEGFault`.

The `align` option is part of the `catch_overflow` option and is only meaningful when `catch_overflow` is enabled. To enable a non-default alignment, set the **MALLOCDEBUG** environment variable as follows:

```
MALLOCDEBUG=catch_overflow,align:n
```

where n is the desired alignment.

To calculate how many bytes of overreads or overwrites the `catch_overflow` option will allow for a given allocation request when n is the requested alignment and `size` is the number of bytes to be allocated, use the following formula:

```
(((size / n) + 1) * n) - size) % n
```

The following example demonstrates the effect of the `align` option on the application's ability to perform overreads or overwrites with the `catch_overflow` option enabled. In this example, the `align` option is specified with a value of 2:

```
MALLOCDEBUG=align:2,catch_overflow
```

The `catch_overflow` option handles overreads and overwrites as follows:

- When an even number of bytes is allocated, `malloc` allocates exactly the number of bytes requested, which will allow for 0 bytes of overreads or overwrites.
- When an odd number of bytes is allocated, `malloc` allocates the number of bytes requested, plus one additional byte to satisfy the required alignment. This allows for 1 byte of possible overreads or overwrites.

override_signal_handling

The `catch_overflow` option reports errors in one of the following ways:

- Memory access errors (such as trying to read or write past the end of allocated memory) cause a segmentation violation (SIGSEGV), resulting in a core dump.
- For other types of errors (such as trying to free space that was already freed), the `catch_overflow` option will output an error message, then call the `abort` function, which will send a SIGIOT signal to end the current process.

If the calling program is blocking or catching the SIGSEGV and the SIGIOT signals, the `catch_overflow` option will be prevented from reporting errors. The `override_signal_handling` option provides a means of bypassing this situation without recoding and rebuilding the application.

If the `override_signal_handling` option is specified, the `catch_overflow` option will perform the following actions upon each call to a **malloc** subsystem routine:

1. Disable any existing signal handlers set up by the application for SIGSEGV or SIGIOT.
2. Set the action for both SIGIOT and SIGSEGV to the default (SIG_DFL).
3. Unblock both SIGIOT and SIGSEGV.

If an application signal handler modifies the action for SIGSEGV between memory allocation routine calls and then attempts an invalid memory access, the `catch_overflow` option will be unable to report the error (the application will not exit and no core file will be produced).

Note:

1. The `override_signal_handling` option can be ineffective in a threaded application environment because the `catch_overflow` option uses the **sigprocmask** subroutine and many threaded processes use the **pthread_sigmask** subroutine.
2. If a thread calls the **sigwait** subroutine without including SIGSEGV and SIGIOT in the signal set and the `catch_overflow` option subsequently detects an error, the thread will hang because the `catch_overflow` option can only generate SIGSEGV or SIGIOT.
3. If a pointer to invalid memory is passed to a kernel routine, the kernel routine will fail and usually return with `errno` set to EFAULT. If the application is not checking the return from the system call, this error might be undetected.

debug_range

By default, if the `catch_overflow` option is enabled, buffer overflow detection is performed for every allocation in the program. If the `debug_range` option is specified, only allocation requests that fall between a user-defined minimum and maximum size will have buffer overflows detected by the `catch_overflow` option. Otherwise, no buffer overflow detection will be performed. This option allows the user to control the amount of extra memory resources consumed by the `catch_overflow` option by only using the tool in specific cases.

The `debug_range` option is only meaningful in the context of the `catch_overflow` option. It is enabled as follows:

```
MALLOCDEBUG=catch_overflow,debug_range:min:max
```

where *min* is the lower bound and *max* is the upper bound of the range in which buffer overflow detection is to be performed. If 0 is specified as a minimum value, then anything that is less than the maximum value will have buffer overflow detection performed. If 0 is specified as a maximum value, then anything that is greater than the minimum value will have buffer overflow detection performed.

Limitation: Due to an internal implementation requirement, each allocation will still necessarily be at least a page size in length. Therefore the `debug_range` option merely reduces the overhead of the `catch_overflow` option rather than eliminating it.

If the `realloc` subroutine is called with an allocation request that falls within the user-specified range, buffer overflow detection is performed even if the original allocation was not within the specified range. The reverse of this is also true.

Note: If the `override_signal` option is set in conjunction with the `debug_range` option, the overriding of the SIGIOT and SIGSEGV signal behavior is performed for all allocations.

functionset

Due to an internal implementation requirement, each allocation will still necessarily be at least a page size in length. Therefore the `functionset` option merely reduces the overhead of the `catch_overflow` option rather than eliminating it.

If the `realloc` subroutine is called from a function that is a member of the user-specified function list, buffer overflow detection is performed even if the original allocation was not made from a specified function. The reverse of this is also true.

Note: If the `override_signal` option is set in conjunction with the `functionset` option, the overriding of the SIGIOT and SIGSEGV signal behavior is performed for all allocations.

The `functionset` option does not check the validity of the functions specified in the list.

allow_overreading

By default, when the `catch_overflow` debug option is enabled and the calling program attempts to read past the end of allocated memory, a segmentation violation will occur and the process will core dump. However, the user may not be interested in catching this type of error, and may have enabled `catch_overflow` in order to catch more dangerous overwrites. Specifying the `allow_overreading` option will cause the `catch_overflow` option to ignore overreads so that other types of errors, which may be considered more serious, can be detected first.

The `allow_overreading` option is only meaningful in the context of the `catch_overflow` option. It is enabled as follows:

```
MALLOCDEBUG=catch_overflow,allow_overreading,
```

postfree_checking

The `postfree_checking` option consumes a substantial amount of extra memory. Programs with very large memory requirements may not be able to use the `postfree_checking` option.

Malloc trace

Malloc Trace is a debugging option designed to allow tracing of all calls to the **malloc** subsystem API through the system trace facility.

Malloc log

Malloc Log is a debugging option designed to provide the user with a runtime database of active allocations in the **malloc** subsystem.

report_allocations

The `report_allocations` option is a tool for detecting memory leaks in an application program. The `report_allocations` option uses the database constructed by Malloc Log to report a list of allocations currently held by the user. A record of each successful allocation is made at the time of the request by Malloc Log. When an allocation is deallocated, Malloc Log removes its record from the database. At process exit, the list of allocations still active is printed to `stderr`, giving a list of allocations that were never freed by their callers.

The `report_allocations` option requires the functionality of Malloc Log to work. Thus, Malloc Log is implicitly enabled when `report_allocations` is enabled. The `report_allocations` option is enabled as follows:

```
MALLOCDEBUG=report_allocations
```

validate_ptrs

By default, the **malloc** subsystem APIs do not validate their input pointers to ensure that they actually reference memory previously allocated. If one of these pointers is invalid, severe heap corruption can occur. Specifying the `validate_ptrs` option causes the **malloc** subsystem APIs to perform extensive validation on their input pointers. If a pointer is found to be invalid (that is, it does not reference memory previously allocated by a call to the **malloc** subsystem API), an error message stating why it is invalid is printed, the `abort` function is called, and a core file is produced. The `validate_ptrs` option is similar to the `verbose` sub-option. The `validate_ptrs` option does not take effect if the `postfree_checking` option is enabled.

The `validate_ptrs` option is enabled as follows:

```
MALLOCDEBUG=validate_ptrs
```

Malloc detect

Malloc Detect is a debugging option designed to detect and report corruption of the internal **malloc** subsystem data structures on every call to a **malloc** subsystem API.

verbose

Sub-option of Malloc Detect.

checkarena

Sub-option of Malloc Detect.

output

By default, the **malloc** debugging options send their output to `stderr`. This may not be desired for all programs. The `output` option exists to provide an alternate destination for printed information. Output can be sent to either `stderr`, `stdout`, or to any file on the system.

The `output` option is enabled as follows:

```
MALLOCDEBUG=output:<filename>
```

continue

Many **malloc** debugging options call `abort()` when they detect an error. This is not always the desired behavior for all programs. The `continue` option exists to instruct the **malloc** subsystem to continue

after the detection of a synchronous error rather than to abort the process. Error messages will still be logged to the appropriate channels.

The continue option is enabled as follows:

```
MALLOCDEBUG=continue
```

Malloc debug fill

Malloc debug fill is a debugging option designed to fill up the memory allocated through the malloc() calls with user specified pattern for debug purposes..

The pattern should be specified as a string (for example, export MALLOCDEBUG=fill:"abc" will set the memory allocated through malloc with the pattern "abc") and a maximum of 128 characters is allowed. If the pattern is not specified, fill option is ignored.

The malloc debug fill option can be enabled as follows:

```
MALLOCDEBUG=fill:pattern
```

Pattern can be a octal or hexadecimal numbers specified in the form of a string. i.e the pattern "\101" , is treated as the octal notation for character 'A' and the pattern "\x41", is treated as the hexadecimal notation for character 'A'

If an invalid octal number is specified, for example \777 which cannot be contained within 1 byte, will be stored as \377, the maximum octal value that can be stored as 1 byte.

Related concepts

[System memory allocation using the malloc subsystem](#)

Memory is allocated to applications using the **malloc** subsystem.

Malloc multiheap

By default, the malloc subsystem uses a single heap, or free memory pool.

However, it also provides an optional multiheap capability to allow the use of multiple heaps of free memory, rather than just one.

The purpose of providing multiple-heap capability in the **malloc** subsystem is to improve the performance of threaded applications running on multiprocessor systems. When the **malloc** subsystem is limited to using a single heap, simultaneous memory-allocation requests received from threads running on separate processors are serialized. The **malloc** subsystem can therefore only service one thread at a time, resulting in a serious impact on multiprocessor system performance.

With malloc multiheap capability enabled, the **malloc** subsystem creates a fixed number of heaps for its use. It will begin to use multiple heaps after the second thread is started (process becomes multithreaded). Each memory-allocation request will be serviced using one of the available heaps. The **malloc** subsystem can then process memory allocation requests in parallel, as long as the number of threads simultaneously requesting service is less than or equal to the number of heaps.

If the number of threads simultaneously requesting service exceeds the number of heaps, additional simultaneous requests will be serialized. Unless this occurs on an ongoing basis, the overall performance of the **malloc** subsystem should be significantly improved when multiple threads are making calls to the **malloc** subroutine in a multiprocessor environment.

Enabling malloc multiheap

Malloc multiheap is not enabled by default. It is enabled and configured by setting the **MALLOPTIOIONS** environment variable. To enable malloc multiheap with the default settings, set MALLOPTIOIONS=multiheap prior to process startup. Setting **MALLOPTIOIONS** in this manner will enable malloc multiheap in its default configuration, with 32 heaps and the fast heap selection algorithm.

Malloc multiheap options

The Malloc Multiheap options are as follows:

- multiheap:n
- considersize

Each of these options is described in detail later in this document.

To set any of these options, use the following syntax:

```
MALLOCOPTIONS=[multiheap:n] | [considersize]
```

One or both options can be specified in any order, as long as options are comma-separated, as in the following example:

```
MALLOCOPTIONS=multiheap:3,considersize
```

In the preceding example, malloc multiheap would be enabled with three heaps and a somewhat slower heap selection algorithm that tries to minimize process size.

Each configuration option should only be specified once when setting **MALLOCOPTIONS**. If a configuration option is specified more than once per setting, only the final instance will apply.

The Malloc Multiheap options are described as follows:

multiheap:n

By default, the maximum number of heaps available to malloc multiheap is 32. The *multiheap:n* option can be used to change the maximum number of heaps to any value from 1 through 32, where *n* is the number of heaps. If *n* is set to a value outside the given range, the default value of 32 is used. Only enable as many heaps as are necessary for the process' requirements. Unnecessarily enabled heaps can increase the amount of fragmentation and waste.

considersize

By default, malloc multiheap selects the next available heap. If the *considersize* option is specified, malloc multiheap will use an alternate heap-selection algorithm that tries to select an available heap that has enough free space to handle the request. This may minimize the working set size of the process by reducing the number of **sbrk** subroutine calls. However, because of the additional processing required, the *considersize* heap-selection algorithm is somewhat slower than the default heap selection algorithm.

If the heaps are unable to allocate space, the **malloc** subroutine will return NULL and set errno to ENOMEM. If there is no available memory in the current heap, the **malloc** subsystem will check the other heaps for available space.

Related concepts

[System memory allocation using the malloc subsystem](#)

Memory is allocated to applications using the **malloc** subsystem.

Malloc buckets

Malloc buckets provides an optional buckets-based extension of the default allocator.

It is intended to improve malloc performance for applications that issue large numbers of small allocation requests. When malloc buckets is enabled, allocation requests that fall within a predefined range of block sizes are processed by malloc buckets. All other requests are processed in the usual manner by the default allocator.

Malloc buckets is not enabled by default. It is enabled and configured prior to process startup by setting the **MALLOCOPTIONS** environment variable.

Bucket composition and sizing

A bucket consists of a block of memory that is subdivided into a predetermined number of smaller blocks of uniform size, each of which is an allocatable unit of memory. Each bucket is identified using a bucket number. The first bucket is bucket 0, the second bucket is bucket 1, the third bucket is bucket 2, and so on. The first bucket is the smallest, and each succeeding bucket is larger in size than the preceding bucket, using a formula described later in this section. A maximum of 128 buckets are available per heap.

The block size for each bucket is a multiple of a bucket-sizing factor. The bucket-sizing factor equals the block size of the first bucket. Each block in the second bucket is twice this size, each block in the third bucket is three times this size, and so on. Therefore, a given bucket's block size is determined as follows:

```
block size = (bucket number + 1) * bucket sizing factor
```

For example, a bucket-sizing factor of 16 would result in a block size of 16 bytes for the first bucket (bucket 0), 32 bytes for the second bucket (bucket 1), 48 bytes for the third bucket (bucket 2), and so on.

The bucket-sizing factor must be a multiple of 8 for 32-bit implementations and a multiple of 16 for 64-bit implementations in order to guarantee that addresses returned from malloc subsystem functions are properly aligned for all data types.

The bucket size for a given bucket is determined as follows:

```
bucket size = number of blocks per bucket * (malloc overhead +  
((bucket number + 1) * bucket sizing factor))
```

The preceding formula can be used to determine the actual number of bytes required for each bucket. In this formula, malloc overhead refers to the size of an internal malloc construct that is required for each block in the bucket. This internal construct is 8 bytes long for 32-bit applications and 16 bytes long for 64-bit applications. It is not part of the allocatable space available to the user, but is part of the total size of each bucket.

The number of blocks per bucket, number of buckets, and bucket-sizing factor are all set with the **MALLOCOPTIONS** environment variable.

Processing allocations from the buckets

A block will be allocated from one of the buckets whenever malloc buckets is enabled and an allocation request falls within the range of block sizes defined by the buckets. Each allocation request is serviced from the smallest possible bucket to conserve space.

If an allocation request is received for a bucket and all of its blocks are already allocated, malloc buckets will automatically enlarge the bucket to service the request. The number of new blocks added to enlarge a bucket is always equal to the number of blocks initially contained in the bucket, which is configured by setting the **MALLOCOPTIONS** environment variable.

Support for multiheap processing

The malloc multiheap capability provides a means to enable multiple malloc heaps to improve the performance of threaded applications running on multiprocessor systems. Malloc buckets supports up to 128 buckets per heap. This allows the **malloc** subsystem to support concurrent enablement of malloc buckets and malloc multiheap so that threaded processes running on multiprocessor systems can benefit from the buckets algorithm.

Enabling malloc buckets

Malloc buckets is not enabled by default, but is enabled and configured by setting the following environment variables:

- **MALLOCTYPE**
- **MALLOCOPTIONS**

The **MALLOCTYPE** environment variable must be set to the default allocator when using Malloc Buckets. To enable malloc buckets with default settings, set **MALLOCOPTIONS=buckets** prior to process startup. To enable malloc buckets with user-specified configuration options, set **MALLOCOPTIONS=buckets, options** prior to process startup, where *options* is a comma-separated list of one or more predefined configuration options.

Malloc buckets configuration options

The **MALLOCOPTIONS** environment variable can be used to provide malloc buckets with one or more of the following predefined configuration options:

```
number_of_buckets:n
bucket_sizing_factor:n
blocks_per_bucket:n
bucket_statistics:[stdout|stderr|pathname]
no_mallinfo
```

Each of these options is described in detail in **MALLOCOPTIONS**.

To set the the **MALLOCOPTIONS** environment variable, use the following syntax:

```
MALLOCOPTIONS=[buckets,[ number_of_buckets:n | bucket_sizing_factor:n | blocks_per_bucket:n |
bucket_statistics:[stdout|stderr|pathname] | no_mallinfo],...]
```

More than one option can be specified (and in any order), as long as options are comma-separated, for example:

```
MALLOCOPTIONS=buckets,number_of_buckets:128,bucket_sizing_factor:8,bucket_statistics:stderr
MALLOCOPTIONS=buckets,bucket_statistics:stdout,blocks_per_bucket:512
```

Commas are the only valid delimiters for separating configuration options in this syntax. The use of other delimiters (such as blanks) between options will cause configuration options to be parsed incorrectly.

Each configuration option should only be specified once when setting the **MALLOCOPTIONS** environment variable. If a configuration option is specified more than once per setting, only the final instance will apply.

If a configuration option is specified with an invalid value, malloc buckets writes a warning message to standard error and then continues execution using a documented default value.

The Malloc Buckets configuration options will be recognized by the **malloc** subsystem only if the buckets option is set, as in the following example:

```
MALLOCOPTIONS=number_of_buckets:8,buckets,bucket_statistics:stderr
```

Malloc buckets options

number_of_buckets:n

The `number_of_buckets:n` option can be used to specify the number of buckets available per heap, where *n* is the number of buckets. The value specified for *n* will apply to all available heaps.

The default value for `number_of_buckets` is 16. The minimum value allowed is 1. The maximum value allowed is 128.

bucket_sizing_factor:n

The `bucket_sizing_factor:n` option can be used to specify the bucket-sizing factor, where *n* is the bucket-sizing factor in bytes.

The value specified for `bucket_sizing_factor` must be a multiple of 8 for 32-bit implementations and a multiple of 16 for 64-bit implementations. The default value for `bucket_sizing_factor` is 32 for 32-bit implementations and 64 for 64-bit implementations.

blocks_per_bucket:n

The `blocks_per_bucket:n` option can be used to specify the number of blocks initially contained in each bucket, where *n* is the number of blocks. This value is applied to all of the buckets. The value of *n* is also used to determine how many blocks to add when a bucket is automatically enlarged because all of its blocks have been allocated.

The default value for `blocks_per_bucket` is 1024.

bucket_statistics:[stdout|stderr|pathname]

The `bucket_statistics` option will cause the **malloc** subsystem to output a statistical summary for malloc buckets upon normal termination of each process that calls the malloc subsystem while malloc buckets is enabled. This summary shows buckets-configuration information and the number of allocation requests processed for each bucket. If multiple heaps have been enabled by way of malloc multiheap, the number of allocation requests shown for each bucket will be the sum of all allocation requests processed for that bucket for all heaps.

The buckets statistical summary will be written to one of the following output destinations, as specified with the `bucket_statistics` option.

- **stdout** - standard output
- **stderr** - standard error
- *pathname* - a user-specified path name

If a user-specified path name is provided, statistical output is appended to the existing contents of the file (if any).

Standard output should not be used as the output destination for a process whose output is piped as input into another process.

The `bucket_statistics` option is disabled by default.

Note: One additional allocation request is always shown in the first bucket for the **atexit** subroutine that prints the statistical summary. For threaded processes, additional allocation requests are shown for some of the buckets due to **malloc** subsystem calls issued by the **pthread**s library.

no_mallinfo

If you specify `MALLOCOPTIONS=no_mallinfo`, the **mallinfo** setting is disabled and the information about the heap that is managed by the **malloc** subsystem is not logged.

Malloc buckets default configuration

The following table summarizes the **malloc** buckets default configuration.

Configuration Option	Default Value (32-bit)	Default Value (64-bit)
number of buckets per heap	16	16
bucket sizing factor	32 bytes	64 bytes
allocation range	1 to 512 bytes (inclusive)	1 to 1024 bytes (inclusive)
number of blocks initially contained in each bucket	1024	1024
bucket statistical summary	disabled	disabled

The default configuration for **malloc** buckets should be sufficient to provide a performance improvement for many applications that issue large numbers of small allocation requests. However, it may be possible to achieve additional gains by setting the **MALLOCOPTIONS** environment variable to modify the default configuration. Before modifying the default configuration, become familiar with the application's memory requirements and usage. **Malloc** buckets can then be enabled with the `bucket_statistics` option to fine-tune the buckets configuration.

Limitations

Because of variations in memory requirements and usage, some applications may not benefit from the memory allocation scheme used by **malloc** buckets. Therefore, it is not advisable to enable **malloc** buckets for system-wide use. For optimal performance, **malloc** buckets should be enabled and configured on a per-application basis.

Related concepts

System memory allocation using the **malloc** subsystem

Memory is allocated to applications using the **malloc** subsystem.

Malloc trace

Malloc Trace provides an optional extension to the **malloc** subsystem for use with the trace facility.

Traces of the **malloc**, **realloc**, and **free** subroutines are recorded for use in problem determination and performance analysis.

Malloc Trace is not enabled by default, but can be enabled and configured prior to process startup through the **MALLOCDEBUG** environment variable.

Events recorded by malloc trace

The tracehook IDs used for Malloc Trace are as follows:

- HKWD_LIBC_MALL_SUBSYSTEM
- HKWD_LIBC_MALL_INTERNAL

When tracing is enabled for HKWD_LIBC_MALL_SUBSYSTEM, the input parameters, as well as return values for each call to **malloc**, **realloc**, and **free** subroutines are recorded in the trace subsystem.

When tracing is enabled for HKWD_LIBC_MALL_INTERNAL and the Malloc Detect debugging tools are enabled, Malloc Trace will log a trace event whenever Malloc Detect detects an error in the internal data structures of the malloc subsystem.

Enabling malloc trace

Malloc Trace is not enabled by default. It is enabled and configured by setting the **MALLOCDEBUG** environment variable. To enable Malloc Trace, set the **MALLOCDEBUG** environment variable by typing the following on the command line:

```
MALLOCDEBUG=trace
```

To enable other Malloc Debug features, set the **MALLOCDEBUG** environment variable as follows:

```
MALLOCDEBUG=[trace, other_option]
```

Limitations

The Malloc Trace debugging feature is compatible with the following malloc policies and options:

- [Default Allocation Policy](#)
- [3.1 Allocation Policy](#)
- [Watson Allocation Policy](#)
- [Watson2 Allocation Policy](#)
- [Malloc Buckets](#)
- [Malloc Multiheap](#)
- [Malloc Thread Cache](#)
- [Malloc Disclaim Option](#)

Related concepts

[System memory allocation using the malloc subsystem](#)

Memory is allocated to applications using the **malloc** subsystem.

Malloc log

Malloc Log is an optional extension of the **malloc** subsystem, allowing the user to obtain information regarding the active allocations currently held by the calling process. This data can then be used in problem determination and performance analysis.

Data recorded in the malloc log

Malloc Log records the following data for each active allocation:

- The address returned to the caller.
- The size of the allocation.
- The heap which the allocation was serviced from.
- The stack traceback of the calling function. The depth of the traceback that is recorded is a configurable option.

In addition, the following information can optionally be logged:

- The process ID of the calling process.
- The thread ID of the calling thread.
- The sequence number of the allocation with respect to process startup.
- The real time in which the allocation was made.

With Malloc Log enabled, each successful allocation requires an additional amount of overhead used to store the metadata. This overhead is around 50-100 bytes for a 32-bit application; twice that for a 64-bit application. The amount of overhead changes depending on the options that are set.

The Malloc Log information can be accessed in the following ways:

- Using the DBX malloc subcommand.
- Using the `report_allocations` malloc debug option.

Enabling malloc log

Malloc Log is not enabled by default. To enable Malloc Log with the default settings, set the **MALLOCDEBUG** environment variable as follows:

```
MALLOCDEBUG=log
```

To enable Malloc Log with user-specified configuration options, set the **MALLOCDEBUG** environment variable as follows:

```
MALLOCDEBUG=log:extended,stack_depth:6
```

Note: The default Malloc Log values are as follows:

extended

Defaults to off. Specifying this option will enable the logging of the optional allocation metadata specified above. This parameter will affect the amount of overhead required for each allocation. This option does not have any effect if `MALLOCTYPE` is set to `watson2`.

stack_depth

Used to specify the depth of the function-call stack that is recorded for each allocation. This parameter affects the amount of overhead required for each allocation. The default value is 4, the maximum is 64.

Limitations

The performance of all programs can degrade when Malloc Log is enabled, due to the cost of storing data to memory. Memory usage will also increase.

The Malloc Log debugging feature is compatible with the following malloc policies and options:

- [Default Allocation Policy](#)
- [Watson Allocation Policy](#)
- [Watson2 Allocation Policy](#)
- [Malloc Buckets](#)
- [Malloc Multiheap](#)
- [Malloc Thread Cache](#)
- [Malloc Disclaim](#)

Related concepts

[System memory allocation using the malloc subsystem](#)

Memory is allocated to applications using the **malloc** subsystem.

Malloc disclaim

Malloc Disclaim is an optional extension of the malloc subsystem, providing the user with a means to enable the automatic disclaiming of memory returned by the **free** subroutine.

This is useful in instances where a process has a high paging-space usage, but is not actually using the memory.

Malloc Disclaim is not enabled by default. It can be enabled and configured prior to process startup through the **MALLOCOPTIONS** environment variable, as follows:

```
MALLOCOPTIONS=disclaim
```

Related concepts

System memory allocation using the malloc subsystem

Memory is allocated to applications using the **malloc** subsystem.

Malloc detect

Malloc Detect provides an optional error reporting and detection extension to the **malloc** subsystem. Information on errors that occurred in the malloc environment will be reported and actions can be performed, if specified.

Malloc Detect can be divided into three distinct capabilities:

- Detecting an error
- Reporting the error using stderr
- Reporting the error using an application-defined function

Detecting an error

Some errors in the **malloc** subsystem are easy to detect. Errors such as freeing a pointer which is not a valid heap address are detected synchronously in the codepath of free. However, errors caused by asynchronous events, such as heap corruption, are much more difficult to pinpoint. The `check_arena` option of Malloc Detect is designed to check for this type of corruption on a synchronous basis. On every call to a **malloc** subsystem API, a check of the internal data structures is performed. If corruption is found, it is then reported to the application. This provides another point in time with which to debug these difficult problems. The `check_arena` option can be enabled by setting the **MALLOCDEBUG** environment variable as follows:

```
MALLOCDEBUG=checkarena
```

Reporting an error using stderr

The normal way to report errors in the **malloc** subsystem is through the return value and the `errno` environment variable. The `verbose` option of Malloc Detect allows these errors to be printed to the standard error stream of the application program. This provides a more visible, more detailed way of reporting errors in malloc. The `verbose` option can be enabled by setting the **MALLOCDEBUG** environment variable as follows:

```
MALLOCDEBUG=verbose
```

Reporting an error using an application-defined function

Malloc Detect also allows the user to provide a function that the **malloc** subsystem will call whenever it encounters an error. Malloc Report will call the application-provided function, wait for a return, and then return as usual. This facility allows the application to perform whatever necessary debug tasks are required before allowing the program to continue. To enable this facility, the application must set the

value of a global function pointer **malloc_err_function** equal to the value of the application's error routine. For example:

```
extern void (*malloc_err_function)(int, ...)
malloc_err_function = &application_malloc_err_hdl
```

Limitations

The Malloc Detect debugging feature is compatible with the following malloc policies and options:

- [Default Allocation Policy](#)
- [Malloc Buckets Option](#)
- [Watson Allocation Policy](#)
- [Malloc Multiheap](#)
- [Malloc Threadcache](#)
- [Malloc Disclaim Option](#)

Configuring and using the malloc thread cache

The Malloc Thread Cache maintains a per-thread pool of unallocated memory for the purpose of reducing contention for the global heap structures.

This cache attempts to preallocate memory pieces for future use according to the pattern of allocations already performed by the thread. If an allocation request can be serviced using one of the unallocated pieces in the Thread Cache, it is removed from the cache and returned to the caller. If an allocation request cannot be serviced using an unallocated piece in the cache, the request is passed on to the global heap structure.

Thread cache allocation strategy

The first time a thread requests memory less than 4096 bytes, the thread cache preallocates multiple chunks of memory of the same size from the global heap structure. It also reserves a larger piece of memory to be used to service future requests. If the thread frees a piece of memory, that piece will be kept in the thread cache for future allocations. However, if during a free the size of the thread cache goes over a certain threshold, half of the elements in the cache will be returned to the backing allocator. Essentially, the behavior of the thread cache can be described as a "batch processor", grouping up individual allocation/deallocation calls to be ran at one single time. This leads to less contention on the global heaps and, in many cases, efficiency gains.

Enabling malloc thread cache

Malloc Thread Cache is enabled by default for the Watson allocator. It can be disabled prior to process startup by setting the **MALLOCOPTIONS** environment variable as follows:

```
$ MALLOCOPTIONS=threadcache:off
```

Malloc Thread Cache is not enabled by default for the default allocator. It can be enabled prior to process startup by setting the **MALLOCOPTIONS** environment variable as follows.

```
$ MALLOCOPTIONS=threadcache
```

Writing reentrant and threadsafe code

In single-threaded processes, only one flow of control exists. The code executed by these processes thus need not be reentrant or threadsafe. In multithreaded programs, the same functions and the same resources may be accessed concurrently by several flows of control.

To protect resource integrity, code written for multithreaded programs must be reentrant and threadsafe.

Reentrance and thread safety are both related to the way that functions handle resources. Reentrance and thread safety are separate concepts: a function can be either reentrant, threadsafe, both, or neither.

This section provides information about writing reentrant and threadsafe programs. It does not cover the topic of writing thread-efficient programs. Thread-efficient programs are efficiently parallelized programs. You must consider thread efficiency during the design of the program. Existing single-threaded programs can be made thread-efficient, but this requires that they be completely redesigned and rewritten.

Reentrance

A reentrant function does not hold static data over successive calls, nor does it return a pointer to static data. All data is provided by the caller of the function. A reentrant function must not call non-reentrant functions.

A non-reentrant function can often, but not always, be identified by its external interface and its usage. For example, the **strtok** subroutine is not reentrant, because it holds the string to be broken into tokens. The **ctime** subroutine is also not reentrant; it returns a pointer to static data that is overwritten by each call.

Thread safety

A threadsafe function protects shared resources from concurrent access by locks. Thread safety concerns only the implementation of a function and does not affect its external interface.

In C language, local variables are dynamically allocated on the stack. Therefore, any function that does not use static data or other shared resources is trivially threadsafe, as in the following example:

```
/* threadsafe function */
int diff(int x, int y)
{
    int delta;

    delta = y - x;
    if (delta < 0)
        delta = -delta;

    return delta;
}
```

The use of global data is thread-unsafe. Global data should be maintained per thread or encapsulated, so that its access can be serialized. A thread may read an error code corresponding to an error caused by another thread. In AIX, each thread has its own **errno** value.

Making a function reentrant

In most cases, non-reentrant functions must be replaced by functions with a modified interface to be reentrant. Non-reentrant functions cannot be used by multiple threads. Furthermore, it may be impossible to make a non-reentrant function threadsafe.

Returning data

Many non-reentrant functions return a pointer to static data. This can be avoided in the following ways:

- Returning dynamically allocated data. In this case, it will be the caller's responsibility to free the storage. The benefit is that the interface does not need to be modified. However, backward compatibility is not ensured; existing single-threaded programs using the modified functions without changes would not free the storage, leading to memory leaks.
- Using caller-provided storage. This method is recommended, although the interface must be modified.

For example, a **strtoupper** function, converting a string to uppercase, could be implemented as in the following code fragment:

```
/* non-reentrant function */
char *strtoupper(char *string)
{
    static char buffer[MAX_STRING_SIZE];
    int index;
```



```

    for (index = 0; string[index]; index++)
        buffer[index] = toupper(string[index]);
    buffer[index] = 0;

    return buffer;
}

```

This function is not reentrant (nor threadsafe). To make the function reentrant by returning dynamically allocated data, the function would be similar to the following code fragment:

```

/* reentrant function (a poor solution) */
char *strtoupper(char *string)
{
    char *buffer;
    int index;

    /* error-checking should be performed! */
    buffer = malloc(MAX_STRING_SIZE);

    for (index = 0; string[index]; index++)
        buffer[index] = toupper(string[index]);
    buffer[index] = 0;

    return buffer;
}

```

A better solution consists of modifying the interface. The caller must provide the storage for both input and output strings, as in the following code fragment:

```

/* reentrant function (a better solution) */
char *strtoupper_r(char *in_str, char *out_str)
{
    int index;

    for (index = 0; in_str[index]; index++)
        out_str[index] = toupper(in_str[index]);
    out_str[index] = 0;

    return out_str;
}

```

The non-reentrant standard C library subroutines were made reentrant using caller-provided storage.

Keeping data over successive calls

No data should be kept over successive calls, because different threads may successively call the function. If a function must maintain some data over successive calls, such as a working buffer or a pointer, the caller should provide this data.

Consider the following example. A function returns the successive lowercase characters of a string. The string is provided only on the first call, as with the **strtok** subroutine. The function returns 0 when it reaches the end of the string. The function could be implemented as in the following code fragment:

```

/* non-reentrant function */
char lowercase_c(char *string)
{
    static char *buffer;
    static int index;
    char c = 0;

    /* stores the string on first call */
    if (string != NULL) {
        buffer = string;
        index = 0;
    }

    /* searches a lowercase character */
    for (; c = buffer[index]; index++) {
        if (islower(c)) {
            index++;
            break;
        }
    }

    return c;
}

```

```
    return c;
}
```

h

This function is not reentrant. To make it reentrant, the static data, the **index** variable, must be maintained by the caller. The reentrant version of the function could be implemented as in the following code fragment:

```
/* reentrant function */
char reentrant_lowercase_c(char *string, int *p_index)
{
    char c = 0;

    /* no initialization - the caller should have done it */

    /* searches a lowercase character */
    for (; c = string[*p_index]; (*p_index)++) {
        if (islower(c)) {
            (*p_index)++;
            break;
        }
    }
    return c;
}
```

The interface of the function changed and so did its usage. The caller must provide the string on each call and must initialize the index to 0 before the first call, as in the following code fragment:

```
char *my_string;
char my_char;
int my_index;
...
my_index = 0;
while (my_char = reentrant_lowercase_c(my_string, &my_index)) {
    ...
}
```

Making a function threadsafe

In multithreaded programs, all functions called by multiple threads must be threadsafe. However, a workaround exists for using thread-unsafe subroutines in multithreaded programs. Non-reentrant functions usually are thread-unsafe, but making them reentrant often makes them threadsafe, too.

Locking shared resources

Functions that use static data or any other shared resources, such as files or terminals, must serialize the access to these resources by locks in order to be threadsafe. For example, the following function is thread-unsafe:

```
/* thread-unsafe function */
int increment_counter()
{
    static int counter = 0;

    counter++;
    return counter;
}
```

To be threadsafe, the static variable **counter** must be protected by a static lock, as in the following example:

```
/* pseudo-code threadsafe function */
int increment_counter();
{
    static int counter = 0;
    static lock_type counter_lock = LOCK_INITIALIZER;

    pthread_mutex_lock(counter_lock);
    counter++;
    pthread_mutex_unlock(counter_lock);
}
```

```
    return counter;
}
```

In a multithreaded application program using the threads library, mutexes should be used for serializing shared resources. Independent libraries may need to work outside the context of threads and, thus, use other kinds of locks.

Workarounds for thread-unsafe functions

It is possible to use a workaround to use thread-unsafe functions called by multiple threads. This can be useful, especially when using a thread-unsafe library in a multithreaded program, for testing or while waiting for a threadsafe version of the library to be available. The workaround leads to some overhead, because it consists of serializing the entire function or even a group of functions. The following are possible workarounds:

- Use a global lock for the library, and lock it each time you use the library (calling a library routine or using a library global variable). This solution can create performance bottlenecks because only one thread can access any part of the library at any given time. The solution in the following pseudocode is acceptable only if the library is seldom accessed, or as an initial, quickly implemented workaround.

```
/* this is pseudo code! */
lock(library_lock);
library_call();
unlock(library_lock);

lock(library_lock);
x = library_var;
unlock(library_lock);
```

- Use a lock for each library component (routine or global variable) or group of components. This solution is somewhat more complicated to implement than the previous example, but it can improve performance. Because this workaround should only be used in application programs and not in libraries, mutexes can be used for locking the library.

```
/* this is pseudo-code! */
lock(library_moduleA_lock);
library_moduleA_call();
unlock(library_moduleA_lock);

lock(library_moduleB_lock);
x = library_moduleB_var;
unlock(library_moduleB_lock);
```

Reentrant and threadsafe libraries

Reentrant and threadsafe libraries are useful in a wide range of parallel (and asynchronous) programming environments, not just within threads. It is a good programming practice to always use and write reentrant and threadsafe functions.

Using libraries

Several libraries shipped with the AIX Base Operating System are threadsafe. In the current version of AIX, the following libraries are threadsafe:

- Standard C library (**libc.a**)
- Berkeley compatibility library (**libbsd.a**)

Some of the standard C subroutines are non-reentrant, such as the **ctime** and **strtok** subroutines. The reentrant version of the subroutines have the name of the original subroutine with a suffix **_r** (underscore followed by the letter *r*).

When writing multithreaded programs, use the reentrant versions of subroutines instead of the original version. For example, the following code fragment:

```

token[0] = strtok(string, separators);
i = 0;
do {
    i++;
    token[i] = strtok(NULL, separators);
} while (token[i] != NULL);

```

should be replaced in a multithreaded program by the following code fragment:

```

char *pointer;
...
token[0] = strtok_r(string, separators, &pointer);
i = 0;
do {
    i++;
    token[i] = strtok_r(NULL, separators, &pointer);
} while (token[i] != NULL);

```

Thread-unsafe libraries may be used by only one thread in a program. Ensure the uniqueness of the thread using the library; otherwise, the program will have unexpected behavior, or may even stop.

Converting libraries

Consider the following when converting an existing library to a reentrant and threadsafe library. This information applies only to C language libraries.

- Identify exported global variables. Those variables are usually defined in a header file with the **export** keyword. Exported global variables should be encapsulated. The variable should be made private (defined with the **static** keyword in the library source code), and access (read and write) subroutines should be created.
- Identify static variables and other shared resources. Static variables are usually defined with the **static** keyword. Locks should be associated with any shared resource. The granularity of the locking, thus choosing the number of locks, impacts the performance of the library. To initialize the locks, the one-time initialization facility may be used.
- Identify non-reentrant functions and make them reentrant. For more information, see [Making a Function Reentrant](#).
- Identify thread-unsafe functions and make them threadsafe. For more information, see [Making a Function threadsafe](#).

Related concepts

[One-time initializations](#)

Some C libraries are designed for dynamic initialization, in which the global initialization for the library is performed when the first procedure in the library is called.

Related information

[admin](#)

[cdc](#)

[delta](#)

[get](#)

[prs](#)

[sccsdiff](#)

[sccsfile](#)

Packaging software for installation

This topic provides information about preparing applications to be installed with the **installp** command.

This section describes the format and contents of the software product installation package that must be supplied by the product developer. It gives a description of the required and optional files that is part of a software installation or update package.

A software product installation package is a backup-format file which contains the files of the software product, required installation control files, and optional installation customization files. The **installp** command is used to install and update software products.

An *installation package* contains one or more separately installable, logically grouped units called *filesets*. Each fileset in a package must belong to the same product.

A *fileset update* or *update package* is a package which contains modifications to an existing fileset.

Throughout this topic, the term *standard system* is used to refer to a system that is not configured as a diskless system.

Installation procedure requirements

- Installation must *not* require user interaction. Product configuration that requires user interaction must occur before or after installation.
- All installations of independent filesets or updates to interdependent filesets must be able to be performed during a single installation.
- No system restart should be required for installation. The installation will stop portions of the system related to the installation, and a system restart is required after installation in order for the installation to take full effect.

Package control information requirements

Package control information must:

- Specify all installation requirements the filesets have on other filesets.
- Specify all file system size requirements for the fileset installation.

Format of a software package

An installation or update package must be a single file in backup format that can be restored by the **installp** command during installation. This file can be distributed on tape, diskette, or CD-ROM.

Package partitioning requirements

In order to support diskless or dataless client workstations, machine-specific portions of the package (the *root part*) must be separated from the machine-shareable portions of the package (the *usr part*). The *usr* part of the package contains files that reside in the `/usr` or `/opt` file system.

Installation of the root part of the package must not modify any files in the `/usr` file system. The `/usr` file system is not writable during installation of the root part of a diskless or dataless client system. The machine-specific (root) part should include everything that is not in the `/usr` or `/opt` file systems.

Packaging for workload partitions

Some software products require special consideration when being packaged for workload partitioning (WPAR). For a software product to successfully deploy in a WPAR, it must be packaged such that it does not attempt to write into the `/usr` or `/opt` file systems during the root part of the processing, because a WPAR mounts those file systems in read-only mode. Similarly, any configuration that is to be performed in each system where the product is installed must be performed from the root part of the package.

If a fileset is not intended to ever be installed into a workload partition, it must be designated with the PRIVATE attribute in the **`lpp_name`** file of the package.

If a fileset needs to be configured differently when it is installed within a WPAR, a packaging script checks the INUWPAR environment variable to determine whether the fileset is being installed within a WPAR.

If a fileset is configured differently when installed into a WPAR, it will be reconfigured when a WPAR is created from a system copy, because the fileset was not installed into a WPAR originally. Fileset owners can create programs in the `/usr/lib/wpars/wparconvert.d/usr` and `/usr/lib/wpars/wparconvert.d/root` directories, which are to be run when converting the *usr* and root parts from their fileset to be run within a system copy WPAR. All executable files within those directories are executed in alphabetical order (C locale) when a system copy WPAR is first started.

Software vital product data

Information about a software product and its installable options is maintained in the Software Vital Product Data (SWVDP) database. The SWVDP consists of a set of commands and the Object Data Manager (ODM) object classes for the maintenance of software product information. The SWVDP commands are provided for the user to query (**lspp**) and verify (**lppchk**) installed software products. The ODM object classes define the scope and format of the software product information that is maintained.

The **installp** command uses the Object Data Manager to maintain the following information in the SWVDP database:

- The name of the software product (for example, bos.adt)
- The version of the software product
- The release level of the software product, which indicates changes to the external programming interface of the software product
- The modification level of the software product, which indicates changes that do not affect the software product's external interface
- The fix level of the software product, which indicates small updates that are to be built into a regular modification level at a later time
- The name, checksum, and size of the files that make up the software product or option
- The state of the software product: available, applying, applied, committing, committed, rejecting, or broken
- Technology level and APAR information
- The destination directory and installer for non-installp packaged software, where applicable.

Software product packaging parts

In order to support installation in the client/server environment, the installation packaging is divided in the following parts:

usr

Contains the part of the product that can be shared among several machines with compatible hardware architectures. For a standard system, these files are stored in the **/usr** or **/opt** file tree.

root

Contains the part of the product that cannot be shared among machines. Each client must have its own copy. Most of this software requiring a separate copy for each machine is associated with the configuration of the machine or product. For a standard system, files in the root part are stored in the root (**/**) file tree. The root part of a fileset must be in the same package as the usr part of the fileset. If a fileset contains a root part, it must also contain a usr part.

Sample file system guide for package partitioning

Following is a brief description of file systems and directories. You can use this as a guide for splitting a product package into root, usr, and share parts.

Some root-part directories and their contents:

/dev

Local machine device files

/etc

Machine configuration files such as **hosts** and **passwd**

/sbin

System utilities needed to boot the system

/var

System-specific data files and log files

Some usr-part directories and their contents include:

/usr/bin

Commands and scripts (ordinary executables)

/usr/sbin

System administration commands

/usr/include

Include files

/usr/lib

Libraries, non-user commands, and architecture-dependent data

/opt

Libraries, non-user commands, and scripts normally associated with non-operating system products.

Package and fileset naming conventions

Use the following conventions when naming a software package and its filesets:

- A package name (*PackageName*) should begin with the product name. If a package has only one installable fileset, the fileset name can be the same as the *PackageName*. All package names must be unique.
- A fileset name has the form:

```
ProductName.PackageName.FilesetName.extension
```

where:

- `ProductName` identifies the product or solution group.
- `PackageName` identifies a functional group within the product.
- `FilesetName` (optional) identifies a specific functional set of files and libraries to be installed.
- `Extension` (optional) further describes the contents.
- A fileset name contains more than one character and begins with a letter.
- All characters in the fileset name must be ASCII characters. Valid characters are upper and lowercase letters, digits, underscores (`_`), plus signs (`+`), and minus signs. The period (`.`) is used as a separator in the fileset name.
- A fileset name cannot end with a period or dot.
- The maximum length for a fileset name is 144 bytes.
- All fileset names must be unique within the package.

Extension	Fileset Description
.adt	Application development toolkit
.com	Common code required by similar filesets
.compat	Compatibility code that may be removed in a future release
.diag	Diagnostics support
.fnt	Fonts
.help. <i>Language</i>	Common Desktop Environment (CDE) help files for a particular language

Special naming considerations for device driver packaging

The configuration manager command (**cfgmgr**) automatically installs software support for detectable devices that are available on the installation media and packaged with the following naming convention:

```
devices.BusTypeID.CardID.Extension
```

where:

- *BusTypeID* specifies the type of bus to which the card attaches (for example, **pci** for PCI)
- *CardID* specifies the unique hexadecimal identifier associated with the card type
- *Extension* specifies the part of the driver that is included (for example, *rte* could be the extension for runtime and *diag* is the extension for diagnostics.)

For example, assume that an ethernet device attaches to the PCI bus and is identified by the configuration manager as having a unique card identifier of 1410bb02. The package of filesets associated with this ethernet device would be named **devices.pci.1410bb02**. The runtime environment fileset within this package would be named **devices.pci.1410bb02.rte**.

Special naming considerations for message catalog packaging

A user installing a package can request that the message catalogs be installed automatically. When this request is made, the system automatically installs message filesets for the primary language if the message filesets are available on the installation media and packaged with the following naming convention:

```
Product.msg.Language.SubProduct
```

The optional *SubProduct* suffix is used when a product has multiple message catalog filesets for the same language, each message catalog fileset applying to a different *SubProduct*. You can choose to have one message fileset for an entire product.

For example, the *Super_Widget* product has a *plastic* and a *metal* set of fileset options. All *Super_Widget* English US message catalogs can be packaged in a single fileset named *Super_Widget.msg.en_US*. If separate message catalog filesets are needed for the *plastic* and *metal* options, the English US message catalog filesets would be named *Super_Widget.msg.en_US.plastic* and *Super_Widget.msg.en_US.metal*.

Note: A message fileset that conforms to this naming convention *must* contain an installed-requisite (**instreq**) on another fileset in the product in order to ensure automatic installation of the message fileset.

File names

Files delivered with the software package cannot have names containing commas or colons. Commas and colons are used as delimiters in the control files used during the software installation process. File names can contain non-ASCII characters. The full path for the file name cannot be more than 128 characters.

Fileset revision level identification

The fileset level is referred to as the *level* or alternatively as the *v.r.m.f* or *VRMF* and has the form:

```
Version.Release.Modification.FixLevel
```

where:

- *Version* is a numeric field of 1 to 2 digits that identifies the version number.
- *Release* is a numeric field of 1 to 2 digits that identifies the release number.
- *Modification* is a numeric field of 1 to 4 digits that identifies the modification level.
- *FixLevel* is a numeric field of 1 to 4 digits that identifies the fix level.

A *base fileset installation level* is the full initial installation level of a fileset. This level contains all files in the fileset, as opposed to a fileset update, which may contain a subset of files from the full fileset.

All filesets in a software package should have the same fileset level, though it is not required for AIX 4.1-formatted packages.

For all new levels of a fileset, the fileset level must increase. The **installp** command uses the fileset level to check for a later level of the product on subsequent installations.

fileset level precedence reads from left to right (for example, 5.2.0.0 is a newer level than 4.3.0.0).

Contents of a software package

This section describes the files contained in an installation or update package. File path names are given for installation package types. For update packages, wherever *PackageName* is part of the path name, it is replaced by *PackageName/FilessetName/FilessetLevel*.

The `usr` part of an installation or update package contains the following installation control files:

- **./lpp_name**: This file provides information about the software package to be installed or updated. For performance reasons, the **lpp_name** file should be the first file in the backup-format file that makes up a software installation package.
- **./usr/lpp/PackageName/liblpp.a**: This archive file contains control files used by the installation process for installing or updating the `usr` part of the software package.
- All files, backed up relative to root, that are to be restored for the installation or update of the `usr` part of the software product.

If the installation or update package contains a root part, the root part contains the following files:

- **./usr/lpp/PackageName/inst_root/liblpp.a**: This library file contains control files used by the installation process for installing or updating the root part of the software package.
- All files that are to be restored for the installation or update of the root part of the software package. For a base fileset installation level these files must be backed up relative to **./usr/lpp/PackageName/inst_root**.

Example contents of a software package

The `farm.apps` package contains the `farm.apps.hog 4.1.0.0` fileset. The `farm.apps.hog 4.1.0.0` fileset delivers the following files:

```
/usr/bin/raisehog (in the usr part of the package)
/usr/sbin/sellhog
(in the usr part of the package)

/etc/hog
(in the root part of the package)
```

The `farm.apps` package contains at least the following files:

```
./lpp_name
./usr/lpp/farm.apps/liblpp.a
./usr/lpp/farm.apps/inst_root/liblpp.a
./usr/bin/raisehog
./usr/sbin/sellhog
./usr/lpp/farm.apps/inst_root/etc/hog
```

Fileset update `farm.apps.hog 4.1.0.3` delivers updates to the following files:

```
/usr/sbin/sellhog
/etc/hog
```

The fileset update package contains the following files:

```
./lpp_name
./usr/lpp/farm.apps/farm.apps.hog/4.1.0.3/liblpp.a
./usr/lpp/farm.apps/farm.apps.hog/4.1.0.3/inst_root/liblpp.a
./usr/sbin/sellhog
./usr/lpp/farm.apps/farm.apps.hog/4.1.0.3/inst_root/etc/hog
```

Note: The file from the root part of the package was restored under an **inst_root** directory. Files installed for the machine-dependent root part of a package are restored relative to an **inst_root** directory. This facilitates installation of machine-specific files in the root file systems of multiple systems. The root part files are installed into the root portions of systems by copying the files from the **inst_root** directory. This allows multiple machines to share a common machine-independent `usr` part.

The `lpp_name` package information file

Each software package must contain the **lpp_name** package information file. The **lpp_name** file gives the **installp** command information about the package and each fileset in the package. Refer to the figure for

an example **lpp_name** file for a fileset update package. The numbers and arrows in the figure refer to fields that are described in the table that follows.

Field Name	Format	Separator	Description
1. Format	Integer	White space	Indicates the release level of installp for which this package was built. The values are: <ul style="list-style-type: none"> • 1 - AIX 3.1 • 3 - AIX 3.2 • 4 - AIX 4.1 and later
2. Platform	Character	White space	Indicates the platform for which this package was built. The values are: <ul style="list-style-type: none"> • R - RISC • I - Intel • N - Neutral
3. Package Type	Character	White space	Indicates whether this is an installation or update package and what type. The values are: <ul style="list-style-type: none"> • I - Installation • S - Single update • SR - Single update required • ML - Technology level update
4. PackageName	Character	White space	The name of the software package (<i>PackageName</i>).
	{	New line	Indicates the beginning of the repeatable sections of fileset-specific data.
5. Fileset name	Character	White space	The complete name of the fileset. This field begins the heading information for the fileset or fileset update.
6. Level	Shown in Description column	White space	The level of the fileset to be installed. The format is: <i>Version.Release.Modification.FixLevel</i> Note: The level can be defined additionally with syntax of <, >, and = combinations. For example, *prereq bos.rte v<5 or *prereq bos.rte v=5 r=3.
7. Volume	Integer	White space	Indicates the volume number where the fileset is located if shipped on multivolume media.
8. Bosboot	Character	White space	Indicates whether a bosboot is needed following the installation. The values are: <ul style="list-style-type: none"> • N - Do not start bosboot • b - start bosboot
9. Content	Character	White space	Indicates that the parts included in the fileset or fileset update. The values are: <ul style="list-style-type: none"> • B -usr and root part • U -usr part only
10. Language	Character	White space	Should be set to the language displayed if the C locale is selected. Generally set to en_US.
11. Description	Character	# or new line	Fileset description. The description is limited to 60 characters.
12. Comments	Character	New line	(Optional) Additional comments.
	[New line	Indicates the beginning of the body of the fileset information.
13. Requisite information	Described following table	New line	(Optional) Installation dependencies the fileset has on other filesets and fileset updates. See the section following this table for detailed description.
	%	New line	Indicates the separation between requisite and size information.

Table 87. *lpp_name* file fields (continued)

Field Name	Format	Separator	Description
14. Size and License Agreement information	Described later in this topic	New line	Size requirements by directory and license agreement information. See Size and License Agreement Information Section later in this topic for detailed description.
	%	New line	Indicates the separation between size and licensing information.
	%	New line	Indicates the separation between licensing and supersede information.
15. Supersede Information	Described later in topic	New line	Information regarding a prior fileset that this fileset replaces.
	%	New line	Indicates the separation between licensing and fix information.
16. Fix information	Described later in topic	New line	Information regarding the fixes contained in the fileset update. See Fix Information Section later in this topic for detailed description.
]	New line	Indicates the end of the body of the fileset information.
	}	New line	Indicates the end of the repeatable sections of fileset-specific information.

Table 88. Examples

```

1 23 4
| || |
4 RSfarm.apps { | 6 7 8 9 10 11
5--> farm.apps.hog04.01.0000.0003 1 N U en_US Hog Utilities
12--># ...
[
13--> *ifreq bos.farming.rte (4.2.0.0) 4.2.0.15
%
14--> /usr/sbin 48
14--> /usr/lpp/farm.apps/farm.apps.hog/4.1.0.3 280
14--> /usr/lpp/farm.apps/farm.apps.hog/inst_root/4.1.0.3.96
14--> /usr/lpp/SAVESPACE 48
14--> /lpp/SAVESPACE 32
14--> /usr/lpp/bos.hos/farm.apps.hog/inst_root/4.1.0.3/ etc 32
%
%
15--> ranch.hog 4.1.0.0
%
16--> IX51366 Hogs producing eggs.
16--> IX81360 Piglets have too many ears.
]
}

```

Requisite information section

The requisite information section contains information about installation dependencies on other filesets or fileset updates. Each requisite listed in the requisite section must be satisfied according to the requisite rules in order to apply the fileset or fileset update.

Before any installing or updating occurs, the **installp** command compares the current state of the filesets to be installed with the requirements listed in the **lpp_name** file. If the **-g** flag was specified with the **installp** command, any missing requisites are added to the list of filesets to be installed. The filesets are ordered for installation according to any prerequisites. Immediately before a fileset is installed, the **installp** command again checks requisites for that fileset. This check verifies that all requisites installed earlier in the installation process were installed successfully and that all requisites are satisfied.

In the following descriptions of the different types of requisites, *RequiredFilesetLevel* represents the minimum fileset level that satisfies the requirements. Except when explicitly blocked for reasons described in [Supersede Information Section](#), newer levels of a fileset satisfy requisites on an earlier level. For example, a requisite on the `plum.tree 2.2.0.0` fileset is satisfied by the `plum.tree 3.1.0.0` fileset.

Prerequisite

A prerequisite indicates that the specified fileset must be installed at the specified fileset level or at a higher level before the current fileset can install successfully. If a prerequisite fileset is scheduled to be installed, the **installp** command orders the list of filesets to install to make sure the prerequisite is met. A prerequisite on a fileset within the same package is not guaranteed.

Syntax

```
*prereq Fileset RequiredFilesetLevel
```

Alternate syntax

```
Fileset RequiredFilesetLevel
```

A fileset update contains an implicit prerequisite to its base-level fileset. If this implicit prerequisite is not adequate, you must specify explicitly a different prerequisite. The *Version* and *Release* of the update and the implicit prerequisite are the same. If the *FixLevel* of the update is 0, the *ModificationLevel* and the *FixLevel* of the implicit prerequisite are both 0. Otherwise, the implicit prerequisite has a *ModificationLevel* that is the same as the *ModificationLevel* of the update and a *FixLevel* of 0. For example, a 4.1.3.2 level update requires its 4.1.3.0 level to be installed before the update installation. A 4.1.3.0 level update requires its 4.1.0.0 level to be installed before the update installation.

Corequisite

A corequisite indicates that the specified fileset must be installed for the current fileset to function successfully. At the end of the installation process, the **installp** command issues warning messages for any corequisites that are not met. A corequisite can be used to indicate requisites between filesets within the same package.

Syntax

```
*coreq Fileset RequiredFilesetLevel
```

If requisite

An if requisite indicates that the specified fileset is required to be at *RequiredFilesetLevel* only if the fileset is installed at *InstalledFilesetLevel*. This is most commonly used to coordinate dependencies between fileset updates. The following example shows an if requisite:

```
*ifreq plum.tree (1.1.0.0) 1.1.2.3
```

Syntax

```
*ifreq Fileset [(InstalledFilesetLevel)] RequiredFilesetLevel
```

If the `plum.tree` fileset is not already installed, this example does not cause it to be installed. If the `plum.tree` fileset is already installed at any of the following levels, this example does not cause the 1.1.2.3 level to be installed:

1.1.2.3

This level matches the *RequiredFilesetLevel*.

1.2.0.0

This level is a different base fileset level.

1.1.3.0

This level supersedes the *RequiredFilesetLevel*.

If the `plum.tree` fileset is already installed at any of the following levels, this example causes the 1.1.2.3 level to be installed:

1.1.0.0

This level matches the *InstalledFilesetLevel*.

1.1.2.0

This level is the same base level as the *InstalledFilesetLevel* and a lower level than the *RequiredFilesetLevel*.

The (*InstalledFilesetLevel*) parameter is optional. If it is omitted, the *Version* and *Release* of the *InstalledFilesetLevel* and the *RequiredFilesetLevel* are assumed to be the same. If the *FixLevel* of the *RequiredFilesetLevel* is 0, the *ModificationLevel* and the *FixLevel* of the *InstalledFilesetLevel* are both 0. Otherwise, the *InstalledFilesetLevel* has a *ModificationLevel* that is the same as the *ModificationLevel* of the *RequiredFilesetLevel* and a *FixLevel* of 0. For example, if the *RequiredFilesetLevel* is 4.1.1.1 and no *InstalledFilesetLevel* parameter is supplied, the *InstalledFilesetLevel* is 4.1.1.0. If the *RequiredFilesetLevel* is 4.1.1.0 and no *InstalledFilesetLevel* parameter is supplied, the *InstalledFilesetLevel* is 4.1.0.0.

Installed requisite

An installed requisite indicates that the specified fileset should be installed automatically only if its corresponding fileset is already installed or is on the list of filesets to install. An installed requisite also is installed if the user explicitly requests that it be installed. A fileset update can not have an installed requisite. Because a fileset containing the message files for a particular package should not be installed automatically without some other part of the package being installed, a message fileset should always contain an installed requisite for another fileset in its product.

Syntax

```
*instreq Fileset RequiredFilesetLevel
```

Group requisite

A group requisite indicates that different requisite conditions can satisfy the requisite. A group requisite can contain prerequisites, corequisites, if-requisites, and nested group requisites. The *Number* preceding the { *RequisiteExpressionList* } identifies how many of the items in the *RequisiteExpressionList* are required. For example, >2 states that at least three items in the *RequisiteExpressionList* are required.

Syntax

```
>Number { RequisiteExpressionList }
```

Requisite information section examples

1. The following example illustrates the use of corequisites. The `book.create 12.30.0.0` fileset cannot function without the `layout.text 1.1.0.0` and `index.generate 2.3.0.0` filesets installed, so the requisite section for `book.create 12.30.0.0` contains:

```
*coreq layout.text 1.1.0.0
*coreq index.generate 2.3.0.0
```

The `index.generate 3.1.0.0` fileset satisfies the `index.generate` requisite, because 3.1.0.0 is a newer level than the required 2.3.0.0 level.

2. The following example illustrates the use of the more common requisite types. Fileset `new.fileset.rte 1.1.0.0` contains the following requisites:

```
*prereq database.rte 1.2.0.0
*coreq spreadsheet.rte 1.3.1.0
*ifreq wordprocessorA.rte (4.1.0.0) 4.1.1.1
*ifreq wordprocessorB.rte 4.1.1.1
```

The `database.rte` fileset must be installed at level 1.2.0.0 or higher before the `new.fileset.rte` fileset can be installed. If `database.rte` and `new.fileset.rte` are installed in the same installation session, the installation program installs the `database` fileset before the `new.fileset.rte` fileset.

The `spreadsheet.rte` fileset must be installed at level 1.3.1.0 or higher for the `new.fileset.rte` fileset to function properly. The `spreadsheet.rte` fileset does not need to be installed before the `new.fileset.rte` fileset is installed, provided both are installed in the same

installation session. If an adequate level of the `spreadsheet.rte` fileset is not installed by the end of the installation session, a warning message will be issued stating that the corequisite is not met.

If the `wordprocessorA.rte` fileset is installed (or being installed with `new.fileset.rte`) at level 4.1.0.0, then the `wordprocessorA.rte` fileset update must be installed at level 4.1.1.1 or higher.

If the `wordprocessorB.rte` fileset is installed (or being installed with `new.fileset.rte`) at level 4.1.1.0, then the `wordprocessorB.rte` fileset update must be installed at level 4.1.1.1 or higher.

3. The following example illustrates an installed requisite. Fileset `Super.msg.fr_FR.Widget` at level 2.1.0.0 contains the following install requisite:

```
*instreq Super.Widget 2.1.0.0
```

The `Super.msg.fr_FR.Widget` fileset cannot be installed automatically when the `Super.Widget` fileset is not installed. The `Super.msg.fr_FR.Widget` fileset can be installed when the `Super.Widget` fileset is not installed if the fileset is listed explicitly on the list of filesets to be installed..

4. The following example illustrates a group requisite. At least one of the prerequisite filesets listed must be installed (both can be installed). If installed, the `spreadsheet_1.rte` fileset must be at level 1.2.0.0 or higher or the `spreadsheet_2.rte` fileset must be at level 1.3.0.0 or higher.

```
>0 {  
*prereq spreadsheet_1.rte 1.2.0.0  
*prereq spreadsheet_2.rte 1.3.0.0  
}
```

Size and license agreement information

The size and license agreement information section contains information about the disk space and license agreement requirements for the fileset.

Size Information

This information is used by the installation process to ensure that enough disk space is available for the installation or update to succeed. Size information has the form:

```
Directory PermanentSpace [TemporarySpace]
```

Additionally, the product developer can specify **PAGESPACE** or **INSTWORK** in the full-path name field to indicate disk space requirements for paging space and workspace needed in the package directory during the installation process.

Directory

The full path name of the directory that has size requirements.

PermanentSpace

The size (in 512-byte blocks) of the permanent space needed for the installation or update.

Permanent space is space that is needed after the installation completes. This field has a different meaning in the following cases:

If *Directory* is **PAGESPACE**, *PermanentSpace* represents the size of page space needed (in 512-byte blocks) to perform the installation.

If *Directory* is **INSTWORK**, *PermanentSpace* represents the number of 512-byte blocks needed for extracting control files used during the installation. These control files are the files that are archived to the **liblpp.a** file.

TemporarySpace

The size (in 512-byte blocks) of the temporary space needed for the installation only. Temporary space is released after the installation completes. The *TemporarySpace* value is optional. An example of temporary space is the space needed to re-link an executable object file. Another example is the

space needed to archive an object file into a library. To archive into a library, the **install** command makes a copy of the library, archives the object file into the copied library, and moves the copied library over the original library. The space for the copy of the library is considered temporary space.

When *Directory* is **INSTWORK**, *TemporarySpace* represents the number of 512-byte blocks needed for the unextracted **liblpp.a** file.

The following example shows a size information section:

```
/usr/bin      30
/lib          40  20
PAGESPACE    10
INSTWORK     10   6
```

Because it is difficult to predict how disk file systems are mounted in the file tree, the directory path name entries in the size information section should be as specific as possible. For example, it is better to have an entry for **/usr/bin** and one for **/usr/lib** than to have a combined entry for **/usr**, because **/usr/bin** and **/usr/lib** can exist on different file systems that are both mounted under **/usr**. In general, it is best to include entries for each directory into which files are installed.

For an update package only, the size information must include any old files (to be replaced) that will move into the **save** directories. These old files will be restored if the update is later rejected. In order to indicate these size requirements, an update package must specify the following special directories:

/usr/lpp/SAVESPACE

The **save** directory for usr part files. By default, the usr part files are saved in the **/usr/lpp/PackageName/FilesetName/FilesetLevel.save** directory.

/lpp/SAVESPACE

The **save** directory for root part files. By default, the root part files are saved in the **/lpp/PackageName/FilesetName/FilesetLevel.save** directory.

License agreement information

A recent addition to the AIX install process allows product owners to require that a customer *sign* a license agreement before the install of the product. The **installp** command reads a fileset's *lpp_name* file as usual for any image. If there are LAF or LAR entries in the size section of the *lpp_name* file, the **installp** command calls the **inulag** command which displays the license in a window and records the customer acceptance of the license. If the customer refuses to accept the license, the installation is halted for that product.

Where to place the license file

Placement of the license files is up to the product owner. However, it is strongly suggested that the license be placed in **/usr/swlag/LANG**. The suggested name for the License Agreement file is given as **ProductName_VersionRelease.la**. There is no requirement to use this name or location. The **installp** command and this packaging is simply supplying a method to deliver the information to customers. All tools and file content must be provided by the product.

Translation requirements for the license file

If there is a requirement that the License Agreement files be translated into the supported languages, it is recommended that the files be separated into different files by language. This translation can require that several files be created.

How to ship the license file

A product can ship the license file as part of its main product or in a separate fileset for the product. Many products are choosing to create a separate fileset for shipping just the license files. This enables a product with many different features to create one license file and fileset that can ship on all the media required for each feature instead of having to include the files in several different features. The current suggestions for the fileset name are **lpp.license** or **lpp.loc.license**. Most products are currently using the first suggestion. If you want the license fileset to be hidden from the customer in standard installs, use **lpp.loc.license** because the license fileset does not need to be selected for installation.

How to package the license file

The file itself is never listed in the *fileset.al* or the *fileset.inventory* files for the fileset. The `installp` command finds the license because of the entry in the size section of the *ProductName* file. The types of entries are:

LAF

License **A**greement **F**ile tells the `installp` command that this particular license file is shipped in this fileset.

License agreement files are indicated by a size section entry with:

```
LAF<lang>license_file size
```

Term

Meaning

LAF

stands for **L**icense **A**greement **F**ile

<lang>

language to which the file is translated. These are typically entries like `en_US`, `fr_FR`, `Ja_JP`, and `zh_TW`. If `<lang>` is not specified, then the agreement file is assumed to be non-translated and encoded as ASCII. If the license agreement file is translated, `<lang>` must be part of the path so that the requirement entry can be associated with the file.

license_file

is the fullpath to the license file as it will be found in the image and placed on the system. The suggested path takes the form `/usr/swlag/en_US/ProductName_VersionRelease.la`

size

This is the real size in 512-byte blocks of the license file so as to ensure that `installp` has allowed enough space to place the license file on the system.

LAR

The **L**icense **A**greement **R**equisite tells the `installp` command that this fileset requires the installation of the listed license agreement file. This is not the same as a prerequisite because it is in the file and not in a fileset. Files and filesets have different formats and different purposes. They should not be confused.

Term

Meaning

LAR

stands for **L**icense **A**greement **R**equisite

req_license_file

The full path to the license file required to install this fileset. Typically these entries use `%L` in the path instead of the actual language name to allow for the appropriate file to be viewed without forcing the customer to view all of the languages.

Generally one fileset ships the files, and contains all of the LAF entries. The other filesets in the product which require this license only ship the LAR entry. The fileset which ships the LAF entries also contains the files listed at their full path location in the BFF image, but the files are not listed in the *fileset.al* or the *fileset.inventory* files for the fileset. The design used for electronic licenses does *not* require that the files be registered to the SWVPD. The `installp` command:

1. Finds a requirement on the file.
2. Checks the system to see if it was accepted.
3. If the file was not accepted
 - a. Finds the fileset that ships the file.
 - b. Extracts (restores) only that file from the BFF image.
 - c. Display the file to the customer.

LAF fileset example

The following is an example of a fileset that ships the license file:

```
iced.tea.loc.license 03.01.0000.0000 1 N U en_US IcedTea Recipe License Information
[
%
INSTWORK 16 160
LAF/usr/swlag/de_DE/iced.tea.la 24
LAF/usr/swlag/DE_DE/iced.tea.la 24
LAF/usr/swlag/en_US/iced.tea.la 24
LAF/usr/swlag/EN_US/iced.tea.la 24
LAF/usr/swlag/es_ES/iced.tea.la 24
LAF/usr/swlag/ES_ES/iced.tea.la 24
LAF/usr/swlag/fr_FR/iced.tea.la 24
LAF/usr/swlag/FR_FR/iced.tea.la 24
LAF/usr/swlag/it_IT/iced.tea.la 24
LAF/usr/swlag/IT_IT/iced.tea.la 24
LAF/usr/swlag/ja_JP/iced.tea.la 24
LAF/usr/swlag/JA_JP/iced.tea.la 32
LAF/usr/swlag/Ja_JP/iced.tea.la 24
LAF/usr/swlag/ko_KR/iced.tea.la 24
LAF/usr/swlag/KO_KR/iced.tea.la 24
LAF/usr/swlag/pt_BR/iced.tea.la 24
LAF/usr/swlag/PT_BR/iced.tea.la 24
LAF/usr/swlag/ru_RU/iced.tea.la 24
LAF/usr/swlag/RU_RU/iced.tea.la 48
LAF/usr/swlag/zh_CN/iced.tea.la 16
LAF/usr/swlag/zh_TW/iced.tea.la 16
LAF/usr/swlag/Zh_TW/iced.tea.la 16
LAF/usr/swlag/ZH_TW/iced.tea.la 24
%
%
%
%
%
%
%
%
]
```

LAR Fileset example

The following is an example of a fileset which ships the license requirement to the license file:

```
iced.tea.server 03.01.0000.0010 1 N B en_US Iced Tea Recipe Group
[
*prereq bos.net.tcp.client 5.1.0.10
*coreq iced.tea.tools 5.1.0.10
*coreq Java14.sdk 1.4.0.1
%
/usr/bin 624
/usr/lib/objrepos 24
/usr/include 16
/usr/include/sys 56
/usr/lpp/iced.tea 22
/usr/samples/iced.tea 8
/usr/samples/iced.tea/server 504
/usr/lpp/iced.tea/inst_root/etc/tea 8
/usr/iced.tea 8
/usr/lpp/iced.tea/inst_root/etc/tea/Top 8
INSTWORK 208 96
/lpp/iced.tea 104
/etc/tea 8
/etc/objrepos 8
/etc/tea/Top 8
/tmp 0 6
LAR/usr/swlag/%L/iced.tea.la 0
%
%
%
%
%
%
%
%
]
```

Supersede information section

The supersede information section indicates the levels of a fileset or fileset update for which this fileset or fileset update may (or may not) be used as a replacement. Supersede information is optional and is only applicable to AIX 4.1-formatted fileset base installation packages and AIX 3.2-formatted fileset update packages.

A newer fileset supersedes any older version of that fileset unless the supersedes section of the **lpp_name** file identifies the latest level of that fileset it supersedes. In the rare cases where a fileset

does not supersede all earlier levels of that fileset, the **installp** command does not use the fileset to satisfy requisites on levels older than the level of the fileset listed in the `supersedes` section.

A fileset update supersedes an older update for that fileset only if it contains all of the files, configuration processing, and requisite information contained in the older fileset update. The **installp** command determines that a fileset update supersedes another update for that fileset in the following conditions:

- The version, release, and modification levels for the updates are equal, the fix levels are both non-zero, and the update with the higher fix level does not contain a prerequisite on a level of the fileset greater than or equal to the level of the update with the lower fix level.
- The version and release levels for the updates are equal, and the update with the higher modification level does not contain a prerequisite on a level of the fileset greater than or equal to the level of the update with the lower modification level.

For example, the fileset update `farm.apps.hog 4.1.0.1` delivers an update of `/usr/sbin/sellhog`. Fileset update `farm.apps.hog 4.1.0.3` delivers updates to the `/usr/sbin/sellhog` file and the `/etc/hog` file. Fileset update `farm.apps.hog 4.1.1.2` delivers an update to the `/usr/bin/raisehog` file.

Update `farm.apps.hog 4.1.0.3` supersedes `farm.apps.hog 4.1.0.1` because it delivers the same files and applies to the same level, `farm.apps.hog 4.1.0.0`.

Update `farm.apps.hog 4.1.1.2` does not supersede either `farm.apps.hog 4.1.0.3` or `farm.apps.hog 4.1.0.1` because it does not contain the same files and applies to a different level, `farm.apps.hog 4.1.1.0`. Update `farm.apps.hog 4.1.1.0` supersedes `farm.apps.hog 4.1.0.1` and `farm.apps.hog 4.1.0.3`.

Supersede section for fileset installation levels (base levels)

An AIX 4.1-formatted fileset installation package can contain the following supersede entries:

Barrier Entry

Identifies the fileset level where a major incompatibility was introduced. Such an incompatibility keeps the current fileset from satisfying requisites to levels of the fileset earlier than the specified level.

Compatibility Entry

Indicates the fileset can be used to satisfy the requisites of another fileset. A compatibility entry is used when a fileset has been renamed or made obsolete. Only one fileset can supersede a given fileset. You may specify only one compatibility entry for each fileset.

The `lpp_name` file can contain at most one barrier and one compatibility entry for a fileset.

A barrier entry consists of the fileset name and the fileset level when the incompatibility was introduced. A barrier entry is necessary for a fileset only in the rare case that a level of the fileset has introduced an incompatibility such that functionality required by dependent filesets has been modified or removed to such an extent that requisites on previous levels of the fileset are not met. A barrier entry must exist in all subsequent versions of the fileset indicating the latest level of the fileset that satisfies requisites by dependent filesets.

For example, if a major incompatibility was introduced in fileset `Bad.Idea 6.5.6.0`, the supersede information section for each `Bad.Idea` fileset installation package from fileset level `6.5.6.0` onward would contain a `Bad.Idea 6.5.6.0` barrier entry. This barrier entry would prevent a requisite of `Bad.Idea 6.5.4.0` from being satisfied by any levels of `Bad.Idea` greater than or equal to `6.5.6.0`.

A compatibility entry consists of a fileset name (different from the fileset in the package) and a fileset level. The fileset level identifies the level at which requisites to the specified fileset (and earlier levels of that fileset) are met by the fileset in the installation package. The compatibility is useful when the specified fileset is obsolete or has been renamed, and the function from the specified fileset is contained in the current fileset. The fileset level in the compatibility entry should be higher than any level expected to exist for the specified fileset.

As an example, assume the `Year.Full 19.91.0.0` fileset is no longer delivered as a unit but is instead broken into several smaller, individual filesets. Only one of the smaller resulting filesets, perhaps `Winter`

19.94.0.0, should contain a compatibility entry of Year.Full 19.94.0.0. This compatibility entry allows the Winter 19.94.0.0 fileset to satisfy the requisites of filesets dependent on Year.Full at levels 19.94.0.0 and earlier.

Supersedes processing

The **installp** command provides the following special features for installing filesets and fileset updates which supersede other filesets or fileset updates:

- If the installation media does not contain a fileset or fileset update that the user requested to install, a superseding fileset or fileset update on the installation media can be installed.

For example, assume a user invokes the **installp** command with the **-g** flag (automatically install requisites) to install the `farm.apps.hog 4.1.0.2` fileset. If the installation media contains the `farm.apps.hog 4.1.0.4` fileset only, then the **installp** command will install the `farm.apps.hog 4.1.0.4` fileset because it supersedes the requested level.

- If the system and the installation media do not contain a requisite fileset or fileset update, the requisite can be satisfied by a superseding fileset or fileset update.
- If an update is requested for installation and the **-g** flag is specified, the request is satisfied by the newest superseding update on the installation media.

When the **-g** flag is specified with the **installp** command, any update requested for installation (either explicitly or implicitly) is satisfied by the newest superseding update on the installation media. If the user wants to install a particular level of an update, not necessarily the latest level, the user can invoke the **installp** command without the **-g** flag.

- If an update and a superseding update (both on the installation media) are requested for installation, the **installp** command installs the newer update only.

In this case, if a user wishes to apply a certain update and its superseding update from the installation media, the user must do separate **installp** operations for each update level. Note that this kind of operation is meaningless if the two updates are applied and committed (**-ac**). Committing the second update removes the first update from the system.

Fix information section

The fix information section is optional. The fix information section entries contain a fix keyword and a 60 characters or less description of the problem fixed. A fix keyword is a 16 characters or less identifier corresponding to the fix. Fix keywords beginning with `ix`, `iy`, `IY` and `IX` are reserved for use by the operating system manufacturer.

A technology level is a fix that is a major update level. Periodic preventive maintenance packages are technology levels. A technology level identifier begins with the name of the software product (not the package), followed by a single dot (`.`) and an identifying level, such as `farm.4.1.1.0`.

The **liblpp.a** installation control library file

The **liblpp.a** file is an archive file that contains the files required to control the package installation. You can create a **liblpp.a** file for your package using the **ar** command with the **-g** flag on systems later than AIX 4.3 to ensure that a 32 bit archive is created. This section describes many of the files you can put in a **liblpp.a** archive.

Throughout this section, *Fileset* appears in the names of the control files. *Fileset* represents the name of the separate fileset to be installed within the software package. For example, assume the apply list file is described as *Fileset.al*. The apply list file for the `bos.net.tcp.client` option of the `bos.net` software product would then be `bos.net.tcp.client.al`.

For any files you include in the **liblpp.a** archive file other than the files listed in this section, you should use the following naming conventions:

- If the file is used in the installation of a specific fileset, the file name should begin with the *Fileset*. prefix.
- If the file is used as a common file for several filesets in the same package, the file name should begin with the **lpp.** prefix.

Many files described in this section are optional. An optional file is necessary only if the function the file provides is required for the fileset or fileset update. Unless stated, a file pertains to both full installation packages and fileset update packages.

Data files contained in the liblpp.a file

Fileset.al

Apply list. This file lists all files to be restored for this fileset. Files are listed one per line with a path relative to root, as in `./usr/bin/pickle`. An apply list file is required if any files are delivered with the fileset or fileset update.

Fileset.cfginfo

Special instructions file. This file lists one keyword per line, each keyword indicating special characteristics of the fileset or fileset update. The only currently recognized keyword is **BOOT**, which causes a message to be generated after installation is complete indicating that the system needs to be restarted.

Fileset.cfgfiles

List of user-configurable files and processing instructions for use when applying a newer or equal installation level of a fileset that is already installed. Before restoring the files listed in the *Fileset.al* file, the system saves the files listed in *Fileset.cfgfiles* file. Later, these saved files are processed according to the handling methods specified in the *Fileset.cfgfiles* file.

Fileset.copyright

Required copyright information file for the fileset. This file consists of the full name of the software product followed by the copyright notices.

Fileset.err

Error template file used as input to the **errupdate** command to add or delete entries in the Error Record Template Repository. This file is commonly used by device support software. The **errupdate** command creates a *Fileset.undo.err* file for cleanup purposes. See the **errupdate** command for information about the format of the *Fileset.err* file. This file can only be included in the root part of a fileset.

Fileset.fixdata

Stanza format file. This file contains information about the fixes contained in a fileset or fileset update.

Fileset.inventory

The inventory file. This file contains required software vital product data for the files in the fileset or fileset update. The inventory file is a stanza-format file containing an entry for each file to be installed or updated.

Fileset.namelist

List of obsolete filesets and current fileset (if applicable) that once contained files now existing in the fileset to be installed. This file is used for installation of repackaged software products only.

Fileset.odmadd or Fileset.*.odmadd

Stanzas to be added to ODM (Object Data Manager) databases.

Fileset.rm_inv

Remove inventory file. This file is for installation of repackaged software products only and must exist if the fileset is not a direct replacement for an obsolete fileset. This stanza-format file contains names of files that need to be removed from obsolete filesets.

Fileset.size

This file contains the space requirements for the files included in this fileset as described earlier in this section.

Fileset.trc

Trace report template file. The **trcupdate** command uses this file to add, replace, or delete trace report entries in the `/etc/trcfmt` file. The **trcupdate** command creates a *Fileset.undo.trc* file for cleanup purposes. Only the root part of a package can contain *Fileset.trc* files.

lpp.acf

Archive control file for the entire package. This file is needed only when adding or replacing an archive member file to an archive file that already exists on the system. The archive control file consists of

lines containing pairs of the member file in the temporary directory as listed in the *Fileset.al* file and the archive file that the member belongs to, both listed relative to root as in:

```
./usr/ccs/lib/libc/member.o ./usr/ccs/lib/libc.a
```

lpp.README

Readme file. This file contains information the user should read before using the software. This file is optional and can also be named **README**, **lpp.doc**, **lpp.instr**, or **lpp.lps**.

productid

Product identification file. This file consists of a single line indicating the product name, the product identifier (20-character limit), and the optional feature number (10-character limit).

Optional executable files contained in the liblpp.a file

The product-specific executable files described in this section are called during the installation process. Unless otherwise noted, file names that end in **_i** are used during installation processing only, and file names that end in **_u** are used in fileset update processing only. All files described in this section are optional and can be either shell scripts or executable object modules. Each program should have a return value of 0 (zero), unless the program is intended to cause the installation or update to fail.

Fileset.config or Fileset.config_u

Modifies configuration near the end of the default installation or update process. *Fileset.config* is used during installation processing only.

Fileset.odmdel or Fileset.*.odmdel

Updates ODM database information for the fileset prior to adding new ODM entries for the fileset. The **odmdel** file naming conventions enables a fileset to have multiple **odmdel** files.

Fileset.pre_d

Indicates whether a fileset may be removed during deinstall. The program must return a value of 0 (zero) if the fileset can be removed. Filesets are removable by default. The program should generate error messages indicating why the fileset is not removable.

Fileset.pre_i or Fileset.pre_u

Runs prior to restoring or saving the files from the apply list in the package, but after removing the files from a previously installed version of the fileset.

Fileset.pre_rej

Runs prior to the reject operation or prior to the preview of a reject operation of the fileset. Use the script to determine whether the fileset can be rejected. Do not use this script to run any commands which change anything on the system. If the script exits with a non-zero return code, the reject operation is not allowed.

Fileset.pre_rm

Runs during a fileset installation prior to removing the files from a previously installed version of the fileset.

Fileset.post_i or Fileset.post_u

Runs after restoring the files from the apply list of the fileset installation or update.

Fileset.unconfig Fileset.unconfig_u

Undoes configuration processing performed in the installation or update during a deinstall or a reject of a fileset. *Fileset.unconfig* is used during deinstallation processing only.

Fileset.unodmadd

Deletes entries that were added to ODM databases during the installation or update.

Fileset.unpost_i or Fileset.unpost_u

Undoes processing performed following restoring the files from the apply list in the installation or update during a deinstall or a reject of a fileset.

Fileset.unpre_i or Fileset.unpre_u

Undoes processing performed prior to restoring the files from the apply list in the installation or update during a deinstall or a reject of a fileset.

If any of these executable files runs a command that may change the device configuration on a machine, that executable file should check the INUCLIENTS environment variable before running the command. If the INUCLIENTS environment variable is set, the command should not be run. The Network Installation Management (NIM) environment uses the **installp** command for many purposes, some of which require the **installp** command to bypass some of its normal processing. NIM sets the INUCLIENTS environment variable when this normal processing must be bypassed.

If the default installation processing is insufficient for your package, you can provide the following executable files in the **liblpp.a** file. If these files are provided in your package, the **installp** command uses your package-provided files in place of the system default files. Your package-provided files must contain the same functionality as the default files or unexpected results can occur. You can use the default files as models for creating your own files. Use of the default files in place of package-provided files is strongly recommended.

instal

Used in place of the default installation script **/usr/lib/instl/instal**. The **installp** command calls this executable file if a fileset in an installation package is applied.

lpp.cleanup

Used in place of the default installation cleanup script **/usr/lib/instl/cleanup**. The **installp** command calls this executable file if a fileset in an installation or update package has been partially applied and must be cleaned up to put the fileset back into a consistent state.

lpp.deinstal

Used in place of the default fileset removal script **/usr/lib/instl/deinstal**. This executable file must be placed in the **/usr/lpp/PackageName** directory. The **installp** command calls this executable file if a fileset in an installation package is removed.

lpp.reject

Used in place of the default installation rejection script **/usr/lib/instl/reject**. The **installp** command calls this executable if a fileset update in an update package is rejected. (The default **/usr/lib/instl/reject** script is a link to the **/usr/lib/instl/cleanup** script.)

update

Used in place of the default fileset update script **/usr/lib/instl/update**. The **installp** command calls this executable file if a fileset in an update package is applied. (The default **/usr/lib/instl/update** script is a link to the **/usr/lib/instl/instal** script.)

To ensure compatibility with the **installp** command, the **instal** or **update** executable provided with a software package must:

- Process all of the filesets in the software package. It can either process the installation for all the filesets or invoke other executables for each fileset.
- Use the **inusable** command to save the current level of any files to be installed.
- Use **inurest** command to restore all required files for the usr part from the distribution media.
- Use the **inucp** command to copy all required files for the root part from the **/usr/lpp/Package_Name/inst_root** directory.
- Create an **\$INUTEMPDIR/status** file indicating success or failure for each fileset being installed or updated.
- Return an exit code indicating the status of the installation. If the **instal** or **update** executable file returns a nonzero return code and no **status** file is found, the installation process assumes all filesets failed.

Optional executable file contained in the Fileset.al file

Fileset.unconfig_d

Undoes fileset-specific configuration operations performed during the installation and updates of the fileset. The **Fileset.unconfig_d** file is used when the **-u** flag is specified with the **installp** command.

If this file is not provided and the **-u** flag is specified, the **Fileset.unconfig**, **Fileset.unpost_i**, and **Fileset.unpre_i** operations are performed.

Further description of installation control filesThe Fileset.cfgfiles file

The *Fileset.cfgfiles* file lists configuration files that need to be saved in order to migrate to a new version of the fileset without losing user-configured data. To preserve user-configuration data, a *Fileset.cfgfiles* file must be provided in the proper **liblpp.a** file (usr or root).

The *Fileset.cfgfiles* contains a one-line entry for each file to be saved. Each entry contains the file name (a path name relative to root), a white-space separator, and a keyword describing the handling method for the migration of the file. The handling method keywords are:

preserve

Replaces the installed new version of the file with the saved version from the save directory. After replacing the new file with the saved version, the saved file from the configuration save directory is deleted.

auto_merge

During the *Fileset.post_i* processing, the product-provided executables merge necessary data from the installed new version of the file into the previous version of the file saved in the configuration save directory. After the *Fileset.post_i* processing, the **installp** command replaces the installed new version of the file with the merged version in the configuration save directory (if it exists) and then removes the saved file.

hold_new

Replaces the installed new version of the file with the saved version from the save directory. The new version of the file is placed in the configuration save directory in place of the old version. The user will be able to reference the new version.

user_merge

Leaves the installed new version of the file on the system and keeps the old version of the file in the configuration save directory. The user will be able to reference the old version to perform any merge that may be necessary. This keyword should be avoided if possible.

other

Used in any case where none of the other defined handling methods are sufficient. The **installp** command saves the file in the configuration save directory and provides no further support. Any other manipulation and handling of the configuration file must be done by the product-provided executables. The product developer has the responsibility of documenting the handling of the file.

The *Fileset.post_i* executable can be used to do specific manipulating or merging of configuration data that cannot be done through the default installation processing.

Configuration files listed in the *Fileset.cfgfiles* file are saved in the configuration save directory with the same relative path name given in the *Fileset.cfgfiles* file. The name of the configuration save directory is stored in the **MIGSAVE** environment variable. The save directory corresponds to the part of the package being installed. The following directories are the configuration save directories:

/usr/lpp/save.config

For the usr part

/lpp/save.config

For the root part

If the list of files that you need to save varies depending on the currently installed level of the fileset, the *Fileset.cfgfiles* file must contain the entire list of configuration files that might be found. If necessary, the *Fileset.post_i* executable (or executables provided by other products) must handle the difference.

For example, assume you have a fileset (**change.rte**) that has one file that can be configured. So, in the root **change.rte.cfgfiles**, there is one file listed:

```
/etc/change_user  user_merge
```

When migrating from your old fileset (**change.obj**) to **change.rte**, you cannot preserve this file because the format has changed. However, when migrating from an older level **change.rte** to a newer level **change.rte**, the file can be preserved. In this case, you might want to create a **change.rte.post_i** script that checks to see what fileset you are migrating from and acts appropriately. This way, if a user had made changes to the **/etc/change_user** file, they are saved.

The root **change.bar.post_i** script could be as follows:

```
#!/bin/ksh
rc=0
grep -q change.rte $INSTALLED_LIST
if [ $? = 0 ]
then
mv $MIGSAVE/etc/change_user/ /etc/change_user
rc=1
fi
exit $rc
```

`$INSTALLED_LIST` is created and exported by **installp**. See Installation for Control Files Specifically for Repackaged Products for more information about the `Fileset.installed_list` configuration file. The `$MIGSAVE` variable contains the name of the directory in which the root part configuration files are saved.

The **installp** command does not produce a warning or error message if a file listed in the `Fileset.cfgfiles` file is not found. The **installp** command also does not produce a message for a file that is not found during the phase following `Fileset.post_i` processing when saved configuration files are processed according to their handling methods. If any warning or error messages are desired, the product-provided executables must generate the messages.

As an example of the `Fileset.cfgfiles` file, the `Product_X.option1` fileset must recover user configuration data from three configuration files located in the root part of the fileset. The `Product_X.option1.cfgfiles` is included in the root part of the `liblpp.a` file and contains the following:

```
./etc/cfg_leafpreserve
./etc/cfg_pudding_hold_new
./etc/cfg_newtonpreserve
```

The `Fileset.fixdata` file

Fileset.fixdata

A stanza-format file that describes the fixes contained in the fileset update (or in a fileset installation, if used in place of an update)

The information in this file is added to a fix database. The **instfix** command uses this database to identify fixes installed on the system. If the `Fileset.fixdata` exists in a package, the fix information in the fix database is updated when the package is applied.

Each fix in the fileset should have its own stanza in the `Fileset.fixdata` file. A `Fileset.fixdata` stanza has the following format:

```
fix:
name = FixKeyword
abstract = Abstract
type = {f | p}
filesets = FilesetName FilesetLevel
[FilesetName FilesetLevel ...]
[symptom = [Symptom]]
```

- *FixKeyword* cannot exceed 16 characters.
- *Abstract* describes the fix and can not exceed 60 characters.
- In the **type** field, **f** represents a fix, and **p** represents a preventive maintenance update.
- The **filesets** field contains a new-line separated list of filesets and fileset levels.
- *FilesetLevel* is the initial level in which the fileset delivered all or part of the fix.
- *Symptom* is an optional description of the problem corrected by the fix. *Symptom* does not have a character limit.

The following example shows a `Fileset.fixdata` stanza for problem MS21235. The fix for this problem is contained in two filesets.


```

fix:
name = MS21235
abstract = 82 gigabyte diskette drive unusable on Mars
type = f
filesets = devices.mca.8d77.rte 12.3.6.13
             devices.mca.8efc.rte 12.1.0.2
symptom = The 82 gigabyte subatomic diskettes fail to operate in a Martian environment.

```

The Fileset.inventory file

Fileset.inventory

File that contains specific information about each file that is to be installed or updated for the fileset

sysck

Command that uses the *Fileset.inventory* file to enter the file name, product name, type, checksum, size, link, and symlink information into the software information database

The *Fileset.inventory* file is required for each part of a fileset that installs or update files. If the package has a root part that does not contain files to be installed (it does configuration only), the root part does not require the *Fileset.inventory* file.

Note: The *Fileset.inventory* file does not support files which are greater than 2 GB in size. If you ship a file that is greater than 2 GB, include it in your *fileset.al* file, but not in your *Fileset.inventory* file. **sysck** has not been updated to handle files larger than 2GB, and the **/usr** file system on most machines will not be created with capability for files greater than 2GB (by default).

The inventory file consists of ASCII text in stanza format. The name of a stanza is the full path name of the file to be installed. The stanza name ends with a colon (:) and is followed by a new-line character. The file attributes follow the stanza name and have the format *Attribute=Value*. Each attribute is described on a separate line.

A *Fileset.inventory* stanza has the following format:

```

inventory:
type      = type
class    = inventory,apply,C2_exclude,fileset
owner    = owner_name
group    = group_name
mode     = TCB | SUID | SGID,permissions
target   = fullpath_filename
link     = fullpath_to_hardlink [additional_hardlinks]
size     = <blank> | VOLATILE |size
checksum = <blank> | VOLATILE |"checksum"

```

Table 89. Valid attributes

Attribute	Description
file_name	The full path to the file or link relative to root (./) followed immediately by a colon (:) and then a new line. The maximum length for this full pathname is 255 characters.
type	Type of the file_name entry where a valid type is one of the following: Keywords Meaning FILE Standard file DIRECTORY Directory SYMLINK The fullpath to a symbolic link

Table 89. Valid attributes (continued)

Attribute	Description
class	<p>Specifies how the file_name is to be referenced during installation. This field must have at least two of the following keywords:</p> <p>Type Items</p> <p>Meaning</p> <p>inventory Indicates that the file remains after the installation is complete. The file is added to the inventory SWVPD database. This should not be used with an A type file in the <i>fileset.il</i> file.</p> <p>apply Indicates that the file is to be restored from the install media. The <i>file_name</i> field is listed in the apply list (<i>fileset.al</i>). This should not be used for an I type file in the <i>fileset.il</i> file.</p> <p>C2_exclude Indicates that this file should be excluded from execution on a C2 Secure system. If this flag is used, the file should also be listed in the <i>fileset.tcb</i> file.</p>
owner	Specifies the file owner after installation. Do not use the uid for this field. The attribute value must be the owner name and must be no more than 8 characters.
group	Specifies the file group. Do not use the gid for this field. The attribute value must be the group name and must be no more than 8 characters.
mode	<p>Specifies the file mode. The value must contain the permissions of the file in octal format. Any of the following keywords can precede the permissions value. Items in the list are separated by commas.</p> <p>Mode Items</p> <p>Meaning</p> <p>TCB Part of the "Trusted Computing Base". If the file is SUID root or SGID system, then the file must be TCB.</p> <p>SUID This file has the set user ID bit set. This has no meaning on a DIRECTORY entry.</p> <p>SGID This file has the set group ID bit set. If set on a DIRECTORY entry, this forces all files created in that directory to have the same group as the directory.</p> <p>permissions Must be a three-digit octal number, for example, 644.</p> <p>Note: If the type is SYMLINK, the mode must be 777. No other entries are valid.</p>
link	Lists any hard links to this file. If multiple hard links exist, each fullpath to a hard link is separated by a comma. The hardlinks must reside in the same directory as the source file. If the type of the entry is SYMLINK, link is not valid. The maximum length of the full path name is 255 characters.
target	Valid only for type=SYMLINK . This is the full path filename to the target of the link. If the link being created is from /usr/bin to /bin, the filename would be /bin and the target would be /usr/bin. The maximum length of the full path name is 255 characters.

Table 89. Valid attributes (continued)

Attribute	Description
size	<p>Mode Items Meaning</p> <p>blank If this field is blank, the size of the filename is determined at the time of installation. The drawback to using this option is that if a file was corrupted during installation, the customer is not informed.</p> <p>VOLATILE If the file size is expected to change through normal operation, the value for this attribute must be VOLATILE.</p> <p>size The exact size of the file.</p> <p>Note: Do not include the size field for DIRECTORY or SYMLINK entries.</p>
checksum	<p>Mode Items Meaning</p> <p>blank If this field is blank, the value returned from the sum -r command to be placed in the inventory SWVPD database when the file is installed.</p> <p>VOLATILE Specifies the size of the file in blocks. If the file size is expected to change through normal operation, the value for this attribute must be VOLATILE.</p> <p>checksum The exact sum of the file as returned by the sum -r command. This should be placed in double quotes.</p> <p>Note: Do not include the checksum field for DIRECTORY or SYMLINK entries.</p>
<i>fileset</i>	Indicates the fileset to which the file belongs.

Note: The **sysck** command creates hard links and symbolic links during installation if those links do not exist. The root part symbolic links should be packaged in the root part *Fileset.inventory* file.

fileset.inventory example

The following fileset.inventory example demonstrates the use of type.

```

/usr/bin:
  owner = bin
  group = bin
  mode = 755
  type = directory
  class = apply,inventory,bos.rte

/usr/bin/tcbck:
  owner = root
  group = security
  mode = TCB,SUID,550
  type = file
  class = apply,inventory,bos.rte.security
  size = 99770
  checksum = "17077    98 "

/usr/sbin/tsm:
  owner = root
  group = security
  mode = TCB,SUID,555
  links = /usr/sbin/getty,/usr/sbin/login
  class = apply,inventory,bos.rte,security

```

```
size = 55086
checksum = "57960 54 "
```

Installation control files specifically for repackaged products

Fileset.installed_list

File created by the **installp** command when installing the fileset from a package if it is found that the fileset (or some form of it) is already installed on the system at some level

The software information database is searched to determine if either *Fileset* or any filesets listed in the file *Fileset.namelist* (if it exists) are already installed on the system. If so, the fileset and the installation level are written to the *Fileset.installed_list* file.

If it is created, the *Fileset.installed_list* is available at the time of the calls to the **rminstal** and **instal** executables. The *Fileset.installed_list* file can be located in the following directories, the packaging working directories, or *PackageWorkDirectory*:

/usr/lpp/

PackageName for the usr part

/lpp/

PackageName for the root part

The *Fileset.installed_list* file contains a one-line entry for each fileset that was installed. Each entry contains the fileset name and the fileset level.

As an example, assume that while the `storm.rain 1.2.0.0` fileset is being installed, the **installp** command discovers that `storm.rain 1.1.0.0` is already installed. The **installp** command would then create the *PackageWorkDirectory/storm.rain.installed_list* file with the following contents:

```
storm.rain 1.1.0.0
```

As another example, assume the `Baytown.com` fileset contains a `Baytown.com.namelist` file with the following entries:

```
Pelly.com
GooseCreek.rte
CedarBayou.stream
```

While installing the `Baytown.com 2.3.0.0` fileset, the **installp** command finds that `Pelly.com 1.2.3.0` and `CedarBayou.stream 4.1.3.2` are installed. The **installp** command creates the *PackageWorkDirectory/Baytown.com.installed_list* file with the following contents:

```
Pelly.obj 1.2.3.0
CedarBayou.stream 4.1.3.2
```

Table 90. The *Fileset.namelist* file

Attribute	Description
<i>Fileset.namelist</i>	This file is necessary when the fileset name has changed or the fileset contains files previously packaged in obsolete filesets. It contains names of all filesets that previously contained files currently included in the fileset to be installed. Each fileset name must appear on a separate line.

The *Fileset.namelist* file must be provided in the `usr` or `root` part of the `liblpp.a` file. The *Fileset.namelist* file is only valid for installation packages; it is not valid for update packages.

At the beginning of installation, the **installp** command searches the Software Vital Product Data (SWVPD) to determine if the fileset or any fileset listed in the *Fileset.namelist* file is already installed on the system. The **installp** command writes to the *Fileset.installed_list* file the fileset names and fileset levels that are found installed, making this information available to product-provided executables.

As a simple example of a *Fileset.namelist* file, assume the `small.business` fileset replaces an earlier fileset named `family.business`. The `small.business` product package contains the `small.business.namelist` file in its `usr` part **liblpp.a** file. The `small.business.namelist` file contains the following entry:

```
family.business
```

As a more complex example of a *Fileset.namelist* file would be where a fileset is divided into a client fileset and a server fileset. The `LawPractice.client` and `LawPractice.server` filesets replace the earlier `lawoffice.mgr` fileset. The `LawPractice.server` fileset also contains a few files from the obsolete `BusinessOffice.mgr` fileset. The `LawPractice.client.namelist` file in the **liblpp.a** file for the `LawPractice` package contains the following entry:

```
lawoffice.mgr
```

The `LawPractice.server.namelist` file in the **liblpp.a** file for the `LawPractice` package contains the following entries:

```
lawoffice.mgr
BusinessOffice.mgr
```

If the *Fileset.namelist* file contains only one entry and the current fileset is not a direct replacement for the fileset listed in the *Fileset.namelist* file, you must include a *Fileset.rm_inv* file in the **liblpp.a** file. The installation process uses the *Fileset.namelist* file and the *Fileset.rm_inv* file to determine if a fileset is a direct replacement for another fileset. If the *Fileset.namelist* file contains only one entry and there is no *Fileset.rm_inv* file, the installation process assumes the new fileset is a direct replacement for the old fileset. When the new (replacement) fileset is installed, the installation process removes from the system all files from the old (replaced) fileset, even files not included in the new fileset.

In the previous examples, the `small.business` fileset is a direct replacement for the `family.business` fileset, so a `small.business.rm_inv` file should not exist. The `LawPractice.client` fileset is not a direct replacement for the `lawoffice.mgr` fileset, so a `LawPractice.client.rm_inv` file must exist, even if it is empty.

Example 3:

Filesets `bagel.shop.rte` and `bread.shop.rte` have been shipping separately for years. Now, `bagel.shop.rte` is going to ship as a part of `bread.shop.rte`. For this to happen, the `bread.shop.rte.namelist` file would look like:

```
bread.shop.rte
bagel.shop.rte
```

You must also ship an empty `bread.shop.rte.rm_inv` file to indicate that all files from the `bagel.shop.rte` fileset should be removed from the system.

Table 91. The *Fileset.rm_inv* file

Attribute	Description
<i>Fileset.rm_inv</i>	File that contains a list of files, links, and directories to be removed from the system if they are found installed

This file is used when the current fileset is packaged differently from a previous level of the fileset and the installation process should not remove previously installed files based on the fileset's entries in the inventory database.

A simple name change for a fileset is not sufficient to require a *Fileset.rm_inv* file. The *Fileset.rm_inv* file is necessary when a new fileset is either a subset of a previous fileset or a mixture of parts of previous filesets. If a *Fileset.namelist* file exists and contains entries for more than one fileset, you must use the *Fileset.rm_inv* file to remove previously installed levels of the files from the system.

The *Fileset.rm_inv* file consists of ASCII text in stanza format. The name of a stanza is the full path name of the file or directory to be removed if found on the system. The stanza name ends with a colon (:) and is followed by a new-line character. If attributes follow the stanza name, the attributes have the format *Attribute=Value*. Attributes are used to identify hard links and symbolic links that need to be removed. Each attribute is described on a separate line. The following list describes the valid attributes associated with the listed file:

Table 92. Attributes and descriptions

Attribute	Description
links	One or more hard links to the file. The full path names of the links are delimited by commas.
symlinks	One or more symbolic links to the file. The full path names of the links are delimited by commas.

Note: Links need to be listed twice, once as a standalone stanza and once as an attribute to the file that is linked to.

For example, assume the U.S.S.R 19.91.0.0 fileset contains the following files in the */usr/lib* directory: *moscow*, *leningrad*, *kiev*, *odessa*, and *petrograd* (a symbolic link to *leningrad*). The product developers decide to split the U.S.S.R 19.91.0.0 fileset into two filesets: *Ukraine.lib 19.94.0.0* and *Russia.lib 19.94.0.0*. The *Ukraine.lib* fileset contains the *kiev* and *odessa* files. The *Russia.lib* fileset contains the *moscow* file. The *leningrad* file no longer exists and is replaced by the *st.petersburg* file in the *Russia.lib* fileset.

The *Ukraine.lib.rm_inv* file must exist because the *Ukraine.lib* fileset is not a direct replacement for the U.S.S.R fileset. The *Ukraine.lib.rm_inv* file should be empty because no files need to be removed when the *Ukraine.lib* fileset is installed to clean up the migrating U.S.S.R fileset.

The *Russia.lib.rm_inv* file must exist because the *Russia.lib* fileset is not a direct replacement for the U.S.S.R fileset. If the *Russia.lib.rm_inv* file is used to remove the *leningrad* file when the *Russia.lib* fileset is installed, the *Russia.lib.rm_inv* file would contain the following stanza:

```
/usr/lib/leningrad:
  symlinks = /usr/lib/petrograd
/usr/lib/petrograd:
```

Installation files for supplemental disk subsystems

A disk subsystem that will not configure with the provided SCSI or bus-attached device driver requires its own device driver and configuration methods. These installation files are provided on a supplemental diskette (which accompanies the device) and must be in backup format with a *./signature* file and a *./startup* file. The signature file must contain the string **target**. The startup file must use restore by name to extract the needed files from the supplemental diskette and to run the commands necessary to bring the device to the available state.

Format of distribution media

The following types of media can be used to distribute software product installation packages.

The following sections describe the formats that must be used to distribute multiple product packages on each of these media.

Tape

In order to stack multiple product package images onto either a single tape or a set of tapes, the files on each tape in the set must conform to the following format:

- File 1 is empty. (Reserved for bootable tapes.)
- File 2 is empty. (Reserved for bootable tapes.)

- File 3 contains a table of contents file that describes product package images on the set of tapes. Therefore, each tape in the set contains a copy of the same table of contents file, except for the difference of the tape volume number in a multi-volume set.
- Files 4 through $(N+3)$ contain the backup-format file images for product packages 1 through N .
- A product package image file cannot span across two tapes.
- Each file is followed by an end-of-file tape mark.

CD-ROM

A CD-ROM that is to contain multiple product package images must be compliant with the Rock Ridge Group Protocol. Product packages should be stored in the installation directory, which must contain the following:

- The backup-format file images of the product packages.
- A table of contents file named **.toc** that describes the product package images in the directory.

A multiple volume CD-ROM is a CD-ROM that has an additional directory structure to define a set of CD-ROMs as a single installable unit.

A multiple volume CD-ROM must conform to the following rules:

- A **/installp/mvCD** directory exists with the following contents:
 1. A table of contents file (**.toc**) that describes the product package images on all of the CD-ROMs of the set. Each volume of the CD-ROM must have the same **.toc** in **/installp/mvCD**.
 2. An ASCII file named **volume_id** in which the first line consists of the decimal volume number of the CD in the set1.
 3. A symbolic link named **vol% n**, where **n** is the decimal volume number of the of the CD in the set. The target of the symbolic link must be a relative path to the directory of product packages on that particular volume of the CD. The standard value for the symbol link is **../ppc**.
- The table of contents file (**.toc**) in the **/installp/mvCD** conforms to the standard table of contents format. The special characteristic of the multiple volume **.toc** is that the location of each product package image begins with the directory entry **vol% n**, where **n** indicates the volume which contains that particular product package.

AIX 5.2 Example:

fileset A is in file **A.bff** on volume 1. fileset B is in file **B.bff** on volume 2. The field in the table of contents file in **/installp/mvCD** containing the location of the product package images for A and B are **vol%1/A.bff** and **vol%2/B.bff**, respectively. The field in the table of contents file in **/installp/ppc** of volume 1 contains the location of A as **A.bff**. The field in the table of contents file in **/installp/ppc** of volume 2 contains the location of B as **B.bff**.

The CD-ROM directory structure for AIX 5.1 and later allows for specification of other installers, as well as multiple platforms.

Diskette

In order to stack multiple product package images onto a set of diskettes, the following files must be written to the set of diskettes:

- A table of contents file that describes product package images to be included in the set.
- Each product package image file that is to be included in the set.

The files are written to the set of diskettes using the following rules:

- Write the data as a stream to the diskette set with a volume ID sector inserted in block 0 of each diskette in the set. The data from the last block of one volume is treated as logically adjacent to the data from block 1 of the next volume (the volume ID sector is verified and discarded when read).
- Each file begins on a 512-byte block boundary.

- Write the table of contents file first. Pad this file to fill its last sector with null characters (x '00'). At least one null character is required to mark the end of the table of contents file. Thus, a full sector of null characters may be required.
- Following the table of contents file, write each of the product package image files to successive sectors. Pad each file to fill its last sector using null characters. A null character is not required if the file ends on the block boundary.
- Block 0 of each diskette in the set contains a volume ID sector. The format of this sector is:

Position	Description
Bytes 0:3	A magic number for identification. This is a binary integer with a value of decimal 3609823513=x 'D7298918'.
Bytes 4:15	A date and time stamp (in ASCII) that serves as the identification for the set of diskettes. The format is <i>MonthDayHourMinuteSecondYear</i> . The <i>Hour</i> should be a value from 00 to 23. All date and time fields contain two digits. Thus, <i>Month</i> should be represented as 03 instead of 3, and <i>Year</i> should be represented as 94 instead of 1994.
Bytes 16:19	A binary integer volume number within this set of diskettes. The first diskette in the set is x '00000001'.
Bytes 20:511	Binary zeroes.

Table 93. The table of contents file

Field Name	Format	Separator	Description
1. Volume	Character	White space	For the tape and diskette table of contents file, this is the number of the volume containing this data. For the fixed disk or CD-ROM table of contents file, the volume number is 0.
2. Date and time stamp	<i>mmddhhMMssyy</i>	White space	For tape or diskette, this is the time stamp when volume 1 was created. For fixed disk or CD-ROM, this is the time stamp when the .toc file was created. See Date and Time Stamp Format later in this article for detailed description.
3. Header format	Character	New line	A number indicating the format of the table of contents file. Valid entries are: <ul style="list-style-type: none"> • 1 -AIX Version 3.1 • 2 -Version 3.2 • 3 -AIX Version 4.1 or later • B -mksysb tape (invalid for use by installp)
4. Location of product package image	Character	White space	For tape or diskette, this is a character string in the form: <i>vvv:bbbb:sssssss</i> See Location Format for Tape and Diskette later in this article for detailed description. For fixed disk or CD-ROM, this is the file name of the product package image file. Note that this is the file name only and must not be preceded by any part of the path name.
5. Package specific information	lpp_name file format	New line	The contents of the lpp_name file contained in this product package image. See The lpp_name Package Information File for detailed description.

Note: Items 4 and 5 described in the preceding table are repeated for each product package image contained on the media.

Date and time stamp format

A date and time stamp format is an ASCII string that has the following format:

MonthDayHourMinuteSecondYear

The *Hour* should be a value from 00 to 23. All date and time fields contain two digits. Thus, *Month* should be represented as 03 instead of 3, and *Year* should be represented as 94 instead of 1994.

Location format for tape and diskette

The location has the format of *vvv:bbbb:sssssss* where each letter represents a digit and has the following meaning:

For tape

vvv

is the volume number of the tape.

bbbb

is the file number on the tape of the product package image.

sssssss

is the size of the file in bytes.

For diskette

vvv

is the volume number of the diskette.

bbbb

is the block number on diskette where the product package image file begins.

sssssss

is the size of the file in bytes (including padding to the end of the block boundary).

AIX Relocatable Install

While AIX relocatable installation is now supported with the native AIX installation utilities (such as **installp**, **instfix**, **lspp**, and **lppchk**), the use of relocation is of particular interest to applications that need to be installed within a Workload Partition. This is because default System WPAR configurations do not include a writeable **/usr** or **/opt** filesystem. Therefore, application installation may need to be retargeted to locations other than the traditional **/usr** or **/opt** location.

In addition to being able to install filesets in the default install location (that is, **/**), an administrator can now install relocatable packages into alternate root install locations. This enables the administrator to do the following:

- Install and maintain multiple installations of the same **installp** package in a single instance of the AIX operating system
- Install and maintain multiples versions of the same **installp** package in a single instance of the AIX operating system
- Use native **installp** tracking tools (such as **lppchk**, **lspp**, **instfix**, and **inulag**) to verify and report installation data on all relocated installation instances
- Attach and detach pre-installed software locations on a given system (application hosting).

Terminology

root install path

The base directory where an application is installed. The **installp** default install path is **/**.

default install path

The default root install path (that is, **/**).

relocated install path

Any root install path that is not the default install path. The path location may be any valid path that is not **/** and that has a size no greater than 512 characters.

relocatable application

An application that can be installed in a non-default root install path.

USIL (User Specified Install Location)

A relocated install path instance set up by the user.

USIL

A User Specified Install Location, or USIL, is a tracked relocated install path that is created by the administrator. This location is tracked by the system and can be used as an alternate install path for packages that support relocation. Multiple instances and/or versions of the same software package can be installed on a single system by delegating each installation to a separate USIL. An existing USIL instance may be attached or detached from any given system.

Each USIL instance maintains its own set of Software Vital Product Data (SWVPD) in all three current installp parts:

- *<InstallRoot>/etc/objrepos*
- *<InstallRoot>/usr/lib/objrepos*
- *<InstallRoot>/usr/share/lib/objrepos*

Notes:

1. The current SWVPD object classes include the following:
 - product
 - lpp
 - inventory
 - history
 - fix
 - vendor
 - lag
2. Each USIL instance will mirror the default SWVPD structure within the relocated path.

USIL management commands

/usr/sbin/mkusil

Creates or attaches a new USIL instance.

mkusil -R RelocatePath -c Comments [XFa]

-a

Attach an existing installation as a USIL instance

-c

Comments to include in the USIL definition (visible with the **lsusil** command)

-R

Path to new USIL location. Must be a valid directory.

-X

Automatically expand if space needed

/usr/sbin/lsusil

Lists existing USIL instances.

lsusil [-R RelocatePath | "ALL"]

-R

Path to existing USIL location.

/usr/sbin/rmusil

Removes an existing USIL instance.

rmusil -R RelocatePath

-R

Path to existing USIL location.

Note: The **rmusil** command only removes the USIL reference in the SWVPD. No files are removed in the USIL install path.

/usr/sbin/chusil

Changes an attribute of an existing USIL instance.

chusil -R RelocatePath -c NewComments [X]

- c**
New comments to include in the USIL definition (visible with the **lsusil** command).
- R**
Path to existing USIL location.
- X**
Automatically expand if space needed

Relocatable installation utilities

To preserve code isolation, all USIL changes are isolated to separately compiled modules. The relocated installation utilities include the following user level modules:

- **/usr/sbin/mkusil**
- **/usr/sbin/rmusil**
- **/usr/sbin/lsusil**
- **/usr/sbin/chusil**

Note:

1. Each utility takes the **-R RelocatePath** flag.
2. When working with relocatable **installp** packages, the above utilities must be used.

Requirements for relocatable applications packaging

The application packaging must support relocatable installation. The following are the recommended guidelines:

- A relocatable application package may not deliver (write) inventory objects outside of its root install location.
- A relocatable application package may not deliver (write) data using packaging customization outside of its root install location.
- The relocatable application package must contain the **RELOCATABLE** extended packaging attribute for each relocatable fileset. The fileset is the smallest installable unit that may be relocated.
- The relocatable application package may not have requisites that are located in external relocated paths. It may have requisites to filesets installed in the default install path or in its own install path.

Requirements for relocatable application execution

The application design must support execution from an installation environment. The following are the recommended guidelines:

- The application must have a method to determine its root install location or function such that it has no dependency on the install location.
- The application must reference all application specific executable components relative to its root install location.
- The application must reference all application specific data components relative to its root install location or it must be designed to share the data with other application instances.
- The application should not make any persistent changes outside of its root install location.

USIL connector ODM class object

The USIL connector ODM Class Object resides in the **/etc/objrepos/usilc** file and contains data that links the default SWVPD with all USIL instances.

The following is the entry for this object class that will be contained in swvdpd.cre:

```
/* User Install Location Connector */
/* Connects the default install path to all relocated install paths. */
class usilc {
    vchar path[1024]; /* USIL path */
    vchar comments[2048]; /* USIL Comments */
}
```

```
long    flags;                /* USIL flags */
};
```

Listing all install paths with the -R "ALL" or -R "all" option

The **lspp** and **lppchk** commands can perform listing operations on all install locations if the **-R "ALL"** syntax is used.

Attach/Detach Operations

The attach operation allows the user to integrate an existing detached USIL path into the SWVPD.

For example, if the administrator creates a master USIL instance with various relocatable applications installed for the purposes of application hosting. The administrator then copies or NFS mounts this USIL instance to various systems and uses the attach feature to integrate the USIL instance into the SWVPD. The detach operation removes reference to the USIL instance.

installp licensing

A new USIL instance starts out with an empty LAG (**installp** license agreement ODM object class). Any installation of filesets or LPPs that require a license will require a license acceptance according to normal **installp** conventions. The license acceptance does not span USIL instances.

Trusted Computing Base (TCB)

Installing USIL instances is currently not supported on TCB enabled systems.

Relocatable requisites

A new packaging semantic indicates the relocatable requisite location. A packager can specify that a given requisite should be found in the default install path or in the relocated install path.

The following new requisite semantic applies:

- **prereq_r** = prereq in relocated install path
- **ifreq_r** = ifreq in relocated install path
- **coreq_r** = coreq in relocated install path
- **instreq_r** = instreq in relocated install path

The currently defined requisites types (**prereq**, **ifreq**, **coreq**, and **instreq**) are all default requisites (requisites that apply to the default install location).

TOC changes for relocatable packages

The following is a sample of the new requisite sections in the TOC file:

```
sscp.rte.1.0.0.5.U.PRIVATE.bff 4 R S sscp {
sscp.rte 01.00.0000.0005 1 N B En_US Sscp
[
*coreq bos.games 1.1.1.1 <-- default requisite in default requisite section
*prereq bos.rte 1.1.1.1 <-- default requisite in default requisite section
%
/usr/bin 20
/etc 20
INSTWORK 72 40
%
%
%
IY99999 1 APAR text here.
%
RELOCATABLE <-- attribute tag to denote relocatable package
%
*prereq bos.rte 1.1.1.1 <-- default requisite in relocated requisite section
*coreq_r bos.games 1.1.1.1 <-- relocated requisite in relocated requisite section
]
}
```

Notes:

1. If the relocatable requisite section is present during a relocated installation, it is used as the requisite section for the installation.

2. If the relocatable requisite section is not present during a relocated installation, the default requisite section is used. This means all requisites will be default requisites.
3. A default installation (non-relocated) does not use the relocatable requisite section.

Table 94. The *installp* processing of product packages

Command	Description
Apply	<p>When a fileset in a product installation package is applied, it is installed on the system and it overwrites any preexisting version of that fileset, therefore <i>committing</i> that version of the fileset on the system. The fileset may be removed if the user decides the fileset is no longer required.</p> <p>When a fileset update is applied, the update is installed and information is saved (unless otherwise requested) so that the update can be removed later. Filesset updates that were applied can be committed or rejected later.</p>
Commit	When a fileset update is committed, the information saved during the apply is removed from the system. Committing already applied software does not change the currently active version of a fileset.
Reject	When an update is rejected, information saved during the apply is used to change the active version of the fileset to the version previous to the rejected update. The saved information is then removed from the system. The reject operation is valid only for updates. Many of the same steps in the reject operation are performed in a cleanup operation when a fileset or fileset update fails to complete installation.
Remove	When a fileset is removed, the fileset and its updates are removed from the system independent of their state (applied, committed, or broken). The remove operation is valid only for the installation level of a fileset.

Executables provided within a product package can tailor processing for the apply, reject, and remove operations.

Reinstalling a fileset does not perform the same actions that removing and installing the same fileset do. The reinstall action (see **/usr/lib/instl/rminstal**) cleans up the current files from the previous or the same version, but does not run any of the **unconfig** or **unpre*** scripts. Therefore, do not assume that the **unconfig** script was run. The **.config** script should check the environment before assuming that the **unconfig** was completed.

For example, for a **ras.berry.rte** fileset, the config script adds a line to root's **crontab** file. Reinstalling the **ras.berry.rte** fileset results in two **crontab** entries, because the **unconfig** script was not run on reinstall (which removed the **crontab** entry). The config script should always remove the entry and then add it again.

Processing for the Apply Operation

This section describes the steps taken by the **installp** command when a fileset or fileset update is applied.

1. Restore the **lpp_name** product package information file for the package from the specified device.
2. Verify that the requested filesets exist on the installation media.
3. Check the level of the requested filesets to ensure that they may be installed on the system.
4. Restore control files from the **liblpp.a** archive library file into the *package directory* (**/usr/lpp/Package_Name** for **usr** or **usr/root** packages. The control files specifically for the **root** portion of a **usr/root** package reside in **/usr/lpp/Package_Name/inst_root/liblpp.a**).
5. Check disk space requirements.
6. Check that necessary *requisites* (filesets required to be at certain levels to use or install another fileset) are already installed or are on the list to be installed.
7. Determine if there are license agreement requirements which must be satisfied in order to proceed with the installation.

8. If this is an installation package rather than a fileset update package, search the software vital product data (SWVPD) to see if *Fileset* (the fileset being installed) or any filesets listed in the *Fileset.namelist* file are already installed on the system at any level. If *Fileset* is already installed, write the fileset name and installed level to the *Work_Directory/Fileset.installed_list* file.

If no level of *Fileset* is installed, then if any filesets listed in the *Fileset.namelist* file are installed, list all those filesets and levels in the *Work_Directory/Fileset.installed_list* file. *Work_Directory* is the same as the package directory with the exception of **root** parts, which use **/lpp/Package_Name**.

9. If this is an installation package rather than a fileset update package, run the **/usr/lib/instl/rminstal** script to do the following for each fileset being installed.

Note: Unless otherwise specified, files checked for existence must have been restored from the **liblpp.a** control file library.

- a. If *Fileset.pre_rm* exists, run *Fileset.pre_rm* to perform required steps before removing any files from this version or an existing version of *Fileset*.
- b. If *Work_Directory/Fileset.installed_list* exists, move the existing files listed in *Fileset.cfg* files to the configuration file save directory (indicated by the **MIGSAVE** environment variable).
- c. If a version of *Fileset* is already installed, remove the files and SWVPD information (except history) for *Fileset*.
- d. If *Work_Directory/Fileset.installed_list* exists, and *Fileset.rm_inv* exists or *Fileset.namelist* contains more than one fileset or the only fileset listed in *Fileset.namelist* is *bos.obj*, then do the following:
 - i) Remove files and SWVPD inventory information for files listed in the file *Fileset.rm_inv*.
 - ii) Remove files and SWVPD inventory information for files listed in the file *Fileset.inventory*.
 - iii) Remove other SWVPD information for any filesets listed in **Fileset.namelist** which no longer have any SWVPD inventory information.
- e. If *Work_Directory/Fileset.installed_list* exists and contains only one fileset and *Fileset.namelist* contained only one fileset, Remove files and SWVPD information (except history) for that fileset.
- f. For each part of a product package (**usr** part only or **usr** followed by **root**)
 - i) Set INUTREE (*U* for **usr** and *M* for **root**) and INUTEMPDIR (name of created temporary working directory environment variables).
 - ii) If an **instal** control program exists in the package directory (not recommended), run **./instal**, otherwise, run the default script **/usr/lib/instl/instal**. If an **instal** control program does not exist in the package directory, set **INUSAVEDIR** environment variable.
 - iii) If an **update** control program exists in the package directory (not recommended), run **./update**. If an **update** control program does not exist in the package directory, run the default script **/usr/lib/instl/update**.
 - iv) If a **status** file has been successfully created by **instal** or **update**, Use **status** file to determine the success or failure of each fileset. If a **status** file has not been created, assume all requested filesets in package failed to apply.
 - v) If the apply operation for a fileset was successful, update the Software Vital Product Data (SWVPD), then register any associated license agreement requirements. If the apply operation for a fileset was not successful, run **/usr/lib/instl/cleanup** or the package-supplied **lpp.cleanup** from package directory to clean up the failed filesets.

Processing of the default install or update script

The **instal** or **update** executable is invoked from **installp** with the first parameter being the device being used for the installation or update. The second parameter is the full path name to the file containing the list of filesets to be installed or updated, referred to below as **\$FILESETLIST**. The default **instal** and **update** scripts are linked together; processing varies based on whether it is invoked as **instal** or **update**. The current directory is the package directory. A temporary directory INUTEMPDIR is created in **/tmp** to hold working files.

The flow within the default **instal** and **update** script is as follows:

1. Do the following for each fileset listed in the **\$FILESETLIST**:
 - a. If the fileset is an update, Execute **Fileset.pre_u** (pre_update) if it exists. If the fileset is not an update, execute **Fileset.pre_i** (pre_installation) if it exists.
 - b. Build a master list of files to be restored from the package by appending **Fileset.al** to the new file **INUTEMPDIR/master.al**.
 - c. If this is an update, the files are specified to be saved, and the **lpp.acf** archive control file exists.
Save off the library archive members being updated.
 - d. If the processing is successful, append this fileset to the list to be installed in the **\$FILESETLIST.new** file.
2. If this is an update and file saving is specified, run **inusave** to save current versions of the files.
3. If you are processing the root part, run **inuclp** to copy files from the apply list to root part. If you are not processing root part, run **inurest** to restore files from apply list for the usr part.

4. Do the following for each fileset listed in **\$FILESETLIST.new** file:

Note: Failure in any step is recorded in the **status** file and processing for that fileset ends

- a. Determine if this fileset is installed at the same or older level, or if filesets listed in the **Fileset.namelist** are installed. If so, export the **INSTALLED_LIST** and **MIGSAVE** environment variables. This is called a *migration*.
 - b. If you are processing an update, invoke **Fileset.post_u** if it exists. If **Fileset.post_u** does not exist, invoke **Fileset.post_i** if it exists.
 - c. If **Fileset.cfgfiles** exists, run **/usr/lib/instl/migrate_cfg** to handle processing of configuration files according to their specified handling method.
 - d. Invoke **sysck** to add the information in the **Fileset.inventory** file to the Software Vital Product Database (SWVPD).
 - e. Invoke the **tcbck** command to add the trusted computing base information to the system if the **Fileset.tcb** file exists and the trusted computing base **tcb_enabled** attribute is set in the **/usr/lib/objrepos/PdAt** ODM database.
 - f. Invoke **errupdate** to add error templates if **Fileset.err** exists.
 - g. Invoke **trcupdate** to add trace report format templates if **Fileset.trc** exists.
 - h. If update or if **Work_Directory/Fileset.installed_list** exists, invoke each **Fileset.odmdel** and **Fileset.*.odmdel** script to process ODM database deletion commands.
 - i. Invoke **odmadd** on each existing **Fileset.odmadd** and **Fileset.*.odmadd** to add information to ODM databases.
 - j. If this is an update, invoke **Fileset.config_u** (fileset configuration update) if it exists. Otherwise, invoke **Fileset.config** (fileset configuration) if it exists.
 - k. Update the **status** file indicating successful processing for the fileset.
5. Link control files for needed for fileset removal into the package's **deinstl** directory for future use. These files include the following files that might be present in the package directory:

- **lpp.deinstal**
- **Fileset.al**
- **Fileset.inventory**
- **Fileset.pre_d**
- **Fileset.unpre_i**
- **Fileset.unpre_u**
- **Fileset.unpost_i**
- **Fileset.unpost_u**
- **Fileset.unodmadd**

- *Fileset.unconfig*
- *Fileset.unconfig_u*
- *\$SAVEDIR/Fileset.*.rodmadd*
- *SAVEDIR/Fileset.*.unodmadd*

Processing for the reject and cleanup operations

This section describes the steps taken by the **installp** command when a fileset update is rejected or when a fileset or fileset update fails to complete installation. The default **cleanup** and **reject** scripts located in **/usr/lib/instl** are linked together. Their logic differs slightly depending on whether the script was invoked as **reject** or **cleanup**. For **usr/root** filesets or fileset updates, the **root** part is processed before the **usr** part.

1. If rejecting, check requisites to ensure that all dependent product updates are also rejected.
2. For each part of a package (for example, **usr** and **root**):
 - a. Set **INUTREE** (*U* for **usr** and *M* for **root**.) and **INUTEMPDIR** environment variables.
 - b. If **reject** control file exists in current directory (**INULIBDIR**), invoke **./lpp.reject**. Otherwise, invoke the default script **/usr/lib/instl/reject**.
3. Update the Software Vital Product Data.

The **reject** executable is invoked from **installp** with the first parameter being undefined and the second parameter being the full path name to the file containing the list of filesets (referred to below as **\$FILESETLIST**) to be rejected for the update.

The following files are referenced by the default **cleanup** and **reject** script.

The flow within the default **cleanup** and **reject** script is as follows:

1. Do the following for each fileset listed in **\$FILESETLIST**:
 - a. If invoked as **cleanup**, then read the line in the *Package_Name.s* status file to determine which step the installation failed on and skip ahead to the undo action for that step. A **cleanup** operation will only begin at the step where the installation failed. For example, if the installation of a fileset failed in the *Fileset.post_i* script, then the cleanup operation for that fileset would begin at *i*, because there are no actions to undo from subsequent steps in the installation.
 - b. Undo any configuration processing performed during the installation:

If rejecting an update, invoke *Fileset.unconfig_u* if it exists. Otherwise, invoke *Fileset.unconfig* if it exists.
 - c. Run any *Fileset.*.unodmadd* and/or *Fileset.unodmadd* files to remove Object Data Manager (ODM) entries added during the installation.
 - d. Run any *Fileset.*.rodmadd* and/or *Fileset.rodmadd* exist to replace ODM entries deleted during the installation.
 - e. Invoke **trcupdate** if *Fileset.undo.trc* exists to undo any trace format template changes made during the installation.
 - f. Invoke **errupdate** if *Fileset.undo.err* exists to undo any error format template changes made during the installation.
 - g. Invoke **tcback** to delete the trusted computing base information to the system if the *Fileset.tcb* file exists and the trusted computing base attribute **tcb_enabled** is set in the **/usr/lib/objrepos/PdAt** ODM database.
 - h. Invoke **sysck** if *Fileset.inventory* exists to undo changes to the software information database.
 - i. Undo any post_installation processing performed during the installation:

If this is an update, invoke *Fileset.unpost_u* if it exists. Otherwise, invoke *Fileset.unpost_i* if it exists.
 - j. Build a master apply list (called **master.al**) from *Fileset.al* files.
 - k. Add *Fileset* to **\$FILESETLIST.new**.

2. Do the following if **\$INUTEMPDIR/master.al** exists.
 - a. Change directories to / (**root**).
 - b. Remove all files in **master.al**.
3. Do the following while reading **\$FILESETLIST.new**.
 - a. Call **inurecv** to recover all saved files.
 - b. If this is an update, invoke *Fileset.unpre_u* if it exists. Otherwise, invoke *Fileset.unpre_i* if it exists.
 - c. Delete the install/update control files.
4. Remove the *Package_Name.s* status file.

Processing for the remove operation

This section describes the steps taken by the **installp** command when a fileset is removed. For **usr/root** filesets or fileset updates, the **root** part is processed before the **usr** part.

1. Check requisites to ensure that all dependent filesets are also removed.
2. For each part of a product package (for example, **usr** or **root**):
 - a. Set INUTREE (U for **usr**, M for **root**, and S for **share**) and INUTEMPDIR (**installp** working directory generated in **/tmp**) environment variables.
 - b. Change directory to **INULIBDIR**.
 - c. If the **deinstal** control file exists in current directory, run the **./lpp.deinstalscript**. If the **deinstal** control file does not exist in current directory, run the **/usr/lib/instl/deinstal** default script.
3. Remove files belonging to the fileset from the file system.
4. Remove fileset entries from the SWVPD except for history data.
5. Deactivate license agreement requirement registration for the fileset.

The **deinstal** executable is invoked from **installp** with the first parameter being the full path name to the file containing the list of filesets to be removed, referred to below as **\$FILESETLIST**.

The flow within the default **deinstal** script is as follows:

1. Do the following for each fileset listed in input file **\$FILESETLIST**:
2. If *Fileset.unconfig_d* exists

Execute *Fileset.unconfig_d* to remove all configuration changes, Object Data Manager (ODM) changes, and error and trace format changes, and to undo all operations performed in the post-installation and pre-installation scripts for all updates and the base level installation. Use of this file is *not* recommended.
3. If *Fileset.unconfig_d* does not exist,
 - a. For each update for that fileset, do the following:
 - Run all *Fileset.unconfig_u* scripts to undo any update configuration processing.
 - Run all *Fileset.*.unodmadd* and *Fileset.unodmadd* to delete Object Data Manager (ODM) entries added during the update.
 - Run all *Fileset.*.rodmadd* and *Fileset.rodmadd* to add Object Data Manager (ODM) entries deleted during the update.
 - Run **errupdate** if *Fileset.undo.err* exists to undo error log template changes.
 - Run **trcupdate** if *Fileset.undo.trc* exists to undo trace report template changes.
 - Run any *Fileset.unpost_u* to undo any post-installation customization.
 - b. For the fileset base installation level, do the following:
 - Run any *Fileset.*.unodmadd* and/or *Fileset.unodmadd* to delete Object Data Manager (ODM) entries added during the installation.

- Run any *Fileset.*.rodmadd* and/or *Fileset.rodmadd* to add Object Data Manager (ODM) entries deleted during the installation.
 - Run **errupdate** if *Fileset.undo.err* exists to undo error log template changes.
 - Run **trcupdate** if *Fileset.undo.trc* exists to undo trace report template changes.
 - Run *Fileset.unconfig_i* to undo any installation configuration processing.
 - Run *Fileset.unpost_i* to undo any post-file installation customization.
4. Remove the files and software data information installed with the fileset.
 5. If *Fileset.unconfig_d* does not exist,
 - a. For each update for that fileset, run any *Fileset.unpre_u* to undo any pre-file installation customization.
 - b. For the fileset base installation level, run any *Fileset.unpre_i* to undo any pre-file installation customization.
 6. Delete any empty directories associated with the fileset.

Note: If an error is returned from some call during the execution of the **deinstal** executable, the error will be logged, but execution will continue. This is different from the other scripts because execution for that fileset is normally canceled once an error is encountered. However, once the removal of a fileset has begun, there is no recovery; therefore, removal becomes a best effort once an error is found.

The installation status file

\$INUTEMPDIR/status

File that contains a one-line entry for each fileset that was to be installed or updated

The **installp** command uses this **status** file to determine appropriate processing. If you create installation scripts, your scripts should produce a **status** file that has the correct format. Each line in the **status** file has the format:

Table 95. StatusCode Fileset

Status code	Meaning
s	Success, update SWVPD
f	Failure, perform cleanup procedure
b	Bypass, failed, cleanup not needed
i	Requisite failure, cleanup not needed
v	sysck verification failed

The following example of a **status** file indicates to the **installp** command that the installations for the **tcp.client** and **tcp.server** filesets of **bos.net** package were successful and the installation for the **nfs.client** fileset was not successful.

```
s bos.net.tcp.client
s bos.net.tcp.server
f bos.net.nfs.client
```

Installation commands used during installation and update processing

inucp

Copies files from the **/usr/lpp/Package_Name/inst_root** directory to the **/** (root) file tree when installing the root part.

inulag

Acts as the front end to the subroutines to manage license agreements.

inurecv

Recovers saved files for installation failure or software rejection (**installp -r**).

inurest

Restores files from the distribution medium onto the system using an apply list as input.

inuseve

Saves all files specified by an apply list into the save directory belonging to the software product.

inuumsg

Issues messages from the `inuumsg.cat` message catalog file for the software product being installed.

ckprereq

Verifies compatibility of the software product with any dependencies using requisite information supplied in the `lpp_name` file and information about already installed products found in the SWVDP.

sysck

Checks the inventory information during installation and update procedures.

The **sysck** command is in the `/usr/bin` directory. Other commands listed previously are in the `/usr/sbin` directory.

For examples of their use, refer to the default installation script, `/usr/lib/instl/instal`.

Source code control system

The source code control system (SCCS) is a complete system of commands that allows specified users to control and track changes made to an SCCS file. SCCS files allow several versions of the same file to exist simultaneously, which can be helpful when developing a project requiring many versions of large files.

The SCCS commands support Multibyte Character Set (MBCS) characters.

Introduction to SCCS

The SCCS commands form a complete system for creating, editing, converting, or changing the controls on SCCS files. An SCCS file is any text file controlled with SCCS commands. All SCCS files have the prefix **s.**, which sets them apart from regular text files.

Attention: Using non-SCCS commands to edit SCCS files can damage the SCCS files.

Use SCCS commands on an SCCS file. If you wish to look at the structure of an SCCS file, use the **pg** command or a similar command to view its contents. However, do not use an editor to directly change the file.

To change text in an SCCS file, use an SCCS command (such as the **get** command) to obtain a version of the file for editing, and then use any editor to modify the text. After changing the file, use the **delta** command to save the changes. To store the separate versions of a file, and control access to its contents, SCCS files have a unique structure.

An SCCS file is made up of three parts:

- Delta table
- Access and tracking flags
- Body of the text

Delta table in SCCS files

Instead of creating a separate file for each version of a file, the SCCS file system only stores the changes for each version of a file. These changes are referred to as *deltas*. The changes are tracked by the delta table in every SCCS file.

Each entry in the delta table contains information about who created the delta, when they created it, and why they created it. Each delta has a specific SID (SCCS IDentification number) of up to four digits. The first digit is the release, the second digit the level, the third digit the branch, and the fourth digit the sequence.

An example of an SID number is:

SID = 1.2.1.4

that is, release 1, level 2, branch 1, sequence 4.

No SID digit can be 0, so there cannot be an SID of 2.0 or 2.1.2.0, for example.

Each time a new delta is created, it is given the next higher SID number by default. That version of the file is built using all the previous deltas. Typically, an SCCS file grows sequentially, so each delta is identified only by its release and level. However, a file may branch and create a new subset of deltas. The file then has a trunk, with deltas identified by release and level, and one or more branches, which have deltas containing all four parts of an SID. On a branch, the release and level numbers are fixed, and new deltas are identified by changing sequence numbers.

Note: A file version built from a branch does not use any deltas placed on the trunk after the point of separation.

Control and tracking flags in SCCS files

After the delta table in an SCCS file, a list of flags starting with the @ (at sign) define the various access and tracking options of the SCCS file. Some of the SCCS flag functions include:

- Designating users who may edit the files
- Locking certain releases of a file from editing
- Allowing joint editing of the file
- Cross-referencing changes to a file

Body of an SCCS file

The SCCS file body contains the text for all the different versions of the file. Consequently, the body of the file does not look like a standard text file. Control characters bracket each portion of the text and specify which delta created or deleted it. When the SCCS system builds a specific version of a file, the control characters indicate the portions of text that correspond to each delta. The selected pieces of text are then used to build that specific version.

SCCS flag and parameter conventions

This section lists flags for the SCCS commands.

In most cases, SCCS commands accept two types of parameters:

Parameter	Description
flags	Flags consist of a - (minus sign), followed by a lowercase character, which is sometimes followed by a value. Flags control how the command operates.
<i>File or Directory</i>	These parameters specify the file or files with which the command operates. Using a directory name as an argument specifies all SCCS files in that directory.

File or directory names cannot begin with a - (minus sign). If you specify this sign by itself, the command reads standard input or keyboard input until it reaches an end-of-file character. This is useful when using pipes that allow processes to communicate.

Any flags specified for a command apply to all files on the command line and are processed before any other parameters to that command. Flag placement in the command line is not important. Other parameters are processed left to right. Some SCCS files contain flags that determine how certain commands operate on the file. See the [admin](#) command description of SCCS header flags for more information.

Creating, editing, and updating an SCCS file

You can create, edit, and update an SCCS file using the **admin**, **get**, and **delta** commands.

Creating an SCCS File

admin

Creates an SCCS file or changes an existing SCCS file.

- To create an empty SCCS file named `s.test.c`, enter:

```
admin -n s.test.c
```

Using the **admin** command with the **-n** flag creates an empty SCCS file.

- To convert an existing text file into an SCCS file, enter:

```
admin -itest.c s.test.c
There are no SCCS identification keywords in the file (cm7)

ls
s.test.c test.c
```

If you use the **-i** flag, the **admin** command creates delta 1.1 from the specified file. Once delta 1.1 is created, rename the original text file so it does not interfere with SCCS commands (it will act as a backup):

```
mv test.c back.c
```

The message `There are no SCCS identification keywords in the file (cm7)` does not indicate an error.

- To start the `test.c` file with a release number of 3.1, use the **-r** flag with the **admin** command, as follows:

```
admin -itest.c -r3 s.test.c
```

Editing an SCCS file

Attention: Do not edit SCCS files directly with non-SCCS commands, or you can damage the SCCS files.

get

Gets a specified version of an SCCS file for editing or compiling.

1. To edit an SCCS file, enter the **get** command with the **-e** flag to produce an editable version of the file:

```
get -e s.test.c
1.3
new delta 1.4
67 lines

ls
p.test.c s.test.c test.c
```

The **get** command produces two new files, `p.test.c` and `test.c`. The editable file is `test.c`. The `p.test.c` file is a temporary, uneditable file used by SCCS to keep track of file versions. It will disappear when you update your changes to the SCCS file. Notice also that the **get** command prints the SID of the version built for editing, the SID assigned to the new delta when you update your changes, and the number of lines in the file.

2. Use any editor to edit `test.c`, for example:

```
ed test.c
```

You can now work on your actual file. Edit this file as often as you wish. Your changes will not affect the SCCS file until you choose to update it.

3. To edit a specific version of an SCCS file with multiple versions, enter the **get** command with the **-r** flag :

```
get -r1.3 s.test.c
1.3
67 lines

get -r1.3.1.4 s.test.c
1.3.1.4
50 lines
```

Updating an SCCS File

delta

Adds a set of changes (deltas) to the text of an SCCS file.

1. To update the SCCS file and create a new delta with the changes you made while editing, use the **delta** command:

```
$delta s.test.c
Type comments, terminated with EOF or a blank line:
```

2. The **delta** command prompts you for comments to be associated with the changes you made. For example, enter your comments, and then press the Enter key twice:

```
No id keywords (cm7)
1.2
5 lines inserted
6 lines deleted
12 lines unchanged
```

The **delta** command updates the s.prog.c file with the changes you made to the test.c file. The **delta** command tells you that the SID of the new version is 1.2, and that the edited file inserted 5 lines, deleted 6 lines, and left 12 lines unchanged from the previous version.

Controlling and tracking SCCS file changes

The SCCS command and file system are primarily used to control access to a file and to track who altered a file, why it was altered, and what was altered.

Controlling access to SCCS files

The following types of access can be controlled in an SCCS file system:

File access controls

Directories containing SCCS files should be created with permission code 755 (read, write, and execute permissions for owner; read and execute permissions for group members and others). The SCCS files themselves should be created as read-only files (444). With these permissions, only the owner can use non-SCCS commands to modify SCCS files. If a group can access and modify the SCCS files, the directories should have group write permission.

User access controls

The **admin** command with the **-a** flag can designate a group of users that can make changes to the SCCS file. A group name or number can also be specified with this flag.

Version access controls

The **admin** command can lock, or prevent, various versions of a file from being accessed by the **get** command by using header flags.

-fc

Sets a ceiling on the highest release number that can be retrieved

-ff

Sets a floor on the lowest release number that can be retrieved

-fl

Locks a particular release against being retrieved

Tracking changes to an SCCS file

There are three ways to track changes to an SCCS file:

- Comments associated with each delta
- Modification Request (MR) numbers
- The SCCS commands

Tracking changes with delta comments

After an SCCS file is updated and a new delta created, the system prompts for comments to be associated with that delta. These comments can be up to 512 characters long and can be modified with the **cdc** command.

cdc

Changes the comments associated with a delta

The **get** command with the **-l** flag prints out the delta table and all the delta comments for any version of a file. In addition to storing the comments associated with a delta, the delta table automatically stores the time and date of the last modification, the real user ID at the time of the modification, the serial numbers of the delta and its predecessor, and any MR numbers associated with the delta.

Tracking changes with modification request numbers

The **admin** command with the **-fv** flag prompts for MR numbers each time a delta is created. A program can be specified with the **-fv** flag to check the validity of the MR numbers when an attempt is made to create a new delta in the SCCS file. If the MR validity-checking program returns a nonzero exit value, the update will be unsuccessful.

The MR validity-checking program is created by the user. It can be written to track changes made to the SCCS file and index them against any other database or tracking system.

Tracking changes with SCCS commands

sccsdiff

Compares two SCCS files and prints their differences to standard output

The **delta** command with the **-p** flag acts the same as the **sccsdiff** command when the file is updated. Both of these commands allow you to see what changes have been made between versions.

prs

Formats and prints specified portions of an SCCS file to standard output

This command allows you to find the differences in two versions of a file.

Detecting and repairing damaged SCCS files

You can detect and repair damaged SCCS files using the **admin** command.

Procedure

1. Check SCCS files on a regular basis for possible damage. Any time an SCCS file is changed without properly using SCCS commands, damage may result to the file. The SCCS file system detects this damage by calculating the checksum and comparing it with the one stored in the delta table. Check for damage by running the **admin** command with the **-h** flag on all SCCS files or SCCS directories as shown:

```
admin -h s.file1 s.file2 ...
```

OR

```
admin -h directory1 directory2 ...
```

If the **admin** command finds a file where the computed checksum is not equal to the checksum listed in the SCCS file header, it displays this message:

```
ERROR [s.filename]:
1255-057 The file is damaged. (co6)
```

2. If a file was damaged, try to edit the file again or read a backup copy. Once the checksum has been recalculated, any remaining damage will be undetectable by the **admin** command.

Note: Using the **admin** command with the **-z** flag on a damaged file can prevent future detection of the damage.

3. After fixing the file, run the **admin** command with the **-z** flag and the repaired file name:

```
admin -z s.file1
```

List of additional SCCS commands

The following SCCS commands complete the system for handling SCCS files:

Attention: Using non-SCCS commands with SCCS files can damage the SCCS files.

Command	Description
<u>rm</u>del	Removes the most recent delta on a branch from an SCCS file.
<u>s</u>act	Displays current SCCS file editing status.
<u>s</u>ccs	Administration program for the SCCS system. The sccs command contains a set of pseudo-commands that perform most SCCS services.
<u>s</u>ccshelp	Explains an SCCS error message or command.
<u>u</u>nget	Cancels the effect of a previous use of the get -e command.
<u>v</u>al	Checks an SCCS file to see if its computed checksum matches the checksum listed in the header.
<u>v</u>c	Substitutes assigned values in place of identification keywords.
<u>w</u>hat	Searches a system file for a pattern and displays text that follows it.

Subroutines, example programs, and libraries

This topic provides information about what subroutines are, how to use them, and where they are stored.

Subroutines are stored in libraries to conserve storage space and to make the program linkage process more efficient. A *library* is a data file that contains copies of a number of individual files and control information that allows them to be accessed individually. The libraries are located in the **/usr/ccs/lib** and **/usr/lib** directories. By convention, most of them have names of the form **libname.a** where *name* identifies the specific library.

All include statements should be near the beginning of the first file being compiled, usually in the declarations section before **main()**, and must occur before using any library functions. For example, use the following statement to include the **stdio.h** file:

```
#include <stdio.h>
```

You do not need to do anything special to use subroutines from the Standard C library (**libc.a**). The **cc** command automatically searches this library for subroutines that a program needs. However, if you use subroutines from another library, you must tell the compiler to search that library. If your program uses subroutines from the library **libname.a**, compile your program with the flag **-lname** (lowercase L). The following example compiles the program `myprog.c`, which uses subroutines from the **libdbm.a** library:


```
cc myprog.c -ldb
```

You can specify more than one **-l** (lowercase L) flag. Each flag is processed in the order specified.

If you are using a subroutine that is stored in the Berkeley Compatibility Library, bind to the **libbsd.a** library *before* binding to the **libc.a** library, as shown in the following example:

```
cc myprog.c -lbsd
```

When an error occurs, many subroutines return a value of -1 and set an external variable named **errno** to identify the error. The **sys/errno.h** file declares the **errno** variable and defines a constant for each of the possible error conditions.

In this documentation, all system calls are described as *subroutines* and are resolved from the **libc.a** library. The programming interface to system calls is identical to that of subroutines. As far as a C Language program is concerned, a system call is merely a subroutine call. The real difference between a system call and a subroutine is the type of operation it performs. When a program invokes a system call, a protection domain switch takes place so that the called routine has access to the operating system kernel's privileged information. The routine then operates in kernel mode to perform a task on behalf of the program. In this way, access to the privileged system information is restricted to a predefined set of routines whose actions can be controlled.

Note:

1. The following list represents the wString routines that are obsolete for the 64 bit **libc.a**. Their corresponding 64 bit **libc.a** equivalents are included. The routines for the 32 bit **libc.a** can be found in the [wstring](#) subroutine.

32 Bit only	64 Bit Equivalent
wstrcat	wcscat
wstrchr	wcschr
wstrcmp	wcscoll
wstrncpy	wcscpy
wstrcspn	wcscspn
wstrdup	Not available and has no equivalents in the 64 bit libc.a
wstrlen	wcslen
wstrncat	wcsncat
wstrncpy	wcsncpy
wstrpbrk	wcspbrk
wstrrchr	wcsrchr
wstrspn	wcsspn
wstrtok	wcstok

2. All programs that handle multibyte characters, wide characters, or locale-specific information must call the **setlocale** subroutine at the beginning of the program.
3. Programming in a multithreaded environment requires reentrant subroutines to ensure data integrity.

Related concepts

[Tools and utilities](#)

This section provides an overview of the tools and utilities that you can use to develop C compiled language programs.

Related information

[List of time data manipulation services](#)

[Header Files Overview](#)

[itrunc](#)

[printf](#)

[scanf](#)

[setlocale](#)

[sqrt](#)

128-bit long double floating-point data type

The AIX operating system supports a 128-bit long double data type that provides greater precision than the default 64-bit long double data type. The 128-bit data type can handle up to 31 significant digits (compared to 17 handled by the 64-bit long double). However, while this data type can store numbers with more precision than the 64-bit data type, it does not store numbers of greater magnitude.

The following special issues apply to the use of the 128-bit long double data type:

- Compiling programs that use the 128-bit long double data type
- Compliance with the IEEE 754 standard
- Implementing the 128-bit long double format
- Values of numeric macros

Compiling programs that use the 128-bit long double data type

To compile C programs that use the 128-bit long double data type, use the **xlc128** command. This command is an alias to the **xlc** command with support for the 128-bit data type. The **xlc** command supports only the 64-bit long double data type.

The standard C library, **libc.a**, provides replacements for **libc.a** routines which are implicitly sensitive to the size of a long double. Link with the **libc.a** library when compiling applications that use the 64-bit long double data type. Link applications that use 128-bit long double values with both the **libc128.a** and **libc.a** libraries. When linking, be sure to specify the **libc128.a** library before the **libc.a** library in the library search order.

Compliance with IEEE 754 standard

The 64-bit implementation of the long double data type is fully compliant with the IEEE 754 standard, but the 128-bit implementation is not. Use the 64-bit implementation in applications that must conform to the IEEE 754 standard.

The 128-bit implementation differs from the IEEE standard for long double in the following ways:

- Supports only round-to-nearest mode. If the application changes the rounding mode, results are undefined.
- Does not fully support the IEEE special numbers NaN and INF.
- Does not support IEEE status flags for overflow, underflow, and other conditions. These flags have no meaning for the 128-bit long double implementation.
- The 128-bit long double data type does not support the following math APIs: **atanhl**, **cbrtl**, **copysignl**, **exp2l**, **expm1l**, **fdiml**, **fmal**, **fmaxl**, **fminl**, **hypotl**, **ilogbl**, **llrintl**, **llroundl**, **log1pl**, **log2l**, **logbl**, **lrintl**, **lroundl**, **nanl**, **nearbyintl**, **nextafterl**, **nexttoward**, **nexttowardf**, **nexttowardl**, **remainderl**, **remquol**, **rintl**, **roundl**, **scalblnl**, **scalbnl**, **tgammal**, and **truncl**.

Implementing the 128-bit long double format

A 128-bit long double number consists of an ordered pair of 64-bit double-precision numbers. The first member of the ordered pair contains the high-order part of the number, and the second member contains the low-order part. The value of the long double quantity is the sum of the two 64-bit numbers.

Each of the two 64-bit numbers is itself a double-precision floating-point number with a sign, exponent, and significand. The low-order member has a magnitude that is less than 1 unit in the last place of the high part, so the values of the two 64-bit numbers do not overlap and the entire significand of the low-order number adds precision beyond the high-order number.

This representation results in several issues that must be considered in the use of these numbers:

- The precision of 128-bit long double data type is greater than the precision of the double data type, but the exponent range is the same. Therefore, the magnitude of numbers that are represented by using 128-bit long double data type is slightly greater than the magnitude of 64-bit double precision data type.

- As the absolute value of the magnitude decreases (near the denormal range), the additional precision available in the low-order part also decreases. When the value to be represented is in the denormal range, this representation provides no more precision than the 64-bit double-precision data type.
- The actual number of bits of precision can vary. If the low-order part is much less than 1 ULP of the high-order part, significant bits (either all 0's or all 1's) are implied between the significant of the high-order and low-order numbers. Certain algorithms that rely on having a fixed number of bits in the significand can fail when using 128-bit long double numbers.

Values of numeric macros

Because of the storage method for the long double data type, more than one number can satisfy certain values that are available as macros. The representation of 128-bit long double numbers means that the following macros required by standard C in the **values.h** file do not have clear meaning:

- Number of bits in the mantissa (**LDBL_MANT_DIG**)
- Epsilon (**LDBL_EPSILON**)
- Maximum representable finite value (**LDBL_MAX**)

Number of bits in the mantissa

The number of bits in the significand is not fixed, but for a correctly formatted number (except in the denormal range) the minimum number available is 106. Therefore, the value of the **LDBL_MANT_DIG** macro is 106.

Epsilon

The ANSI C standard defines the value of epsilon as the difference between 1.0 and the least representable value greater than 1.0, that is, $b^{-(1-p)}$, where b is the radix (2) and p is the number of base b digits in the number. This definition requires that the number of base b digits is fixed, which is not true for 128-bit long double numbers.

The smallest representable value greater than 1.0 is this number:

```
0x3FF0000000000000, 0x0000000000000001
```

The difference between this value and 1.0 is this number:

```
0x0000000000000001, 0x0000000000000000
0.4940656458412465441765687928682213E-323
```

Because 128-bit numbers usually provide at least 106 bits of precision, an appropriate minimum value for p is 106. Thus, $b^{-(1-p)}$ and $2^{-(105)}$ yield this value:

```
0x3960000000000000, 0x0000000000000000
0.24651903288156618919116517665087070E-31
```

Both values satisfy the definition of epsilon according to standard C. The long double subroutines use the second value because it better characterizes the accuracy provided by the 128-bit implementation.

Maximum long double value

The value of the **LDBL_MAX** macro is the largest 128-bit long double number that can be multiplied by 1.0 and yield the original number. This value is also the largest finite value that can be generated by primitive operations, such as multiplication and division:

```
0x7FEFFFFFFFFFFFFFFF, 0x7C9FFFFFFFFFFFFFFF
0.1797693134862315907729305190789002575e+309
```

List of character manipulation subroutines

The character manipulation functions and macros test and translate ASCII characters.

These functions and macros are of three kinds:

- Character testing
- Character translation
- Miscellaneous character manipulation

The [Character Testing](#) illustrates some of the character manipulation routines.

Character testing

Use the following functions and macros to determine character type. Punctuation, alphabetic, and case-querying functions values depend on the current collation table.

The [ctype](#) subroutines contain the following functions:

isalpha

Is character alphabetic?

isalnum

Is character alphanumeric?

isupper

Is character uppercase?

islower

Is character lowercase?

isdigit

Is character a digit?

isxdigit

Is character a hex digit?

isspace

Is character a blank-space character?

ispunct

Is character a punctuation character?

isprint

Is character a printing character, including space?

isgraph

Is character a printing character, excluding space?

iscntrl

Is character a control character?

isascii

Is character an integer ASCII character?

Character translation

The [conv](#) subroutines contain the following functions:

toupper

Converts a lowercase letter to uppercase

_toupper

(Macro) Converts a lowercase letter to uppercase

tolower

Converts an uppercase letter to lowercase

_tolower

(Macro) Converts an uppercase letter to lowercase

toascii

Converts an integer to an ASCII character

Miscellaneous character manipulation

getc, fgetc, getchar, getw

Get a character or word from an input stream

putc, putchar, fputc, putw

Write a character or word to a stream

List of executable program creation subroutines

The list of executable program creation services consists of subroutines that support a group of commands.

These commands and subroutines allow you to create, compile, and work with files in order to make your programs run.

Subroutine	Description
<u>_end, _text, _edata</u>	Define the last location of a program
<u>confstr</u>	Determines the current value of a specified system variable defined as a string
<u>getopt</u>	Gets flag letters from the argument vector
<u>ldopen, ldaopen</u>	Open a common object file
<u>ldclose, ldaclose</u>	Close a common object file
<u>ldahread</u>	Reads the archive header of a member of an archive file
<u>ldfhread</u>	Reads the file header of a common object file
<u>ldlread, ldlininit, ldllitem</u>	Read and manipulate line number entries of a common object file function
<u>ldshread, ldnsbread</u>	Read a section header of a common object file
<u>ldtbread</u>	Reads a symbol table entry of a common object file
<u>ldgetname</u>	Retrieves a symbol name from a symbol table entry or from the string table
<u>ldlseek, ldnsseek</u>	Seek to line number entries of a section of a common object file
<u>ldohseek</u>	Seek to the optional file header of a common object file
<u>ldrseek, ldnrseek</u>	Seek to the relocation information for a section of a common object file
<u>ldsseek, ldnsseek</u>	Seek to a section of a common object file
<u>ldtbseek</u>	Seeks to the symbol table of a common object file
<u>ldtbindex</u>	Returns the index of a particular common object file symbol table entry
<u>load</u>	Loads and binds an object module into the current process
<u>unload</u>	Unloads an object file
<u>loadbind</u>	Provides specific runtime resolution of a module's deferred symbols
<u>loadquery</u>	Returns error information from the load subroutine or the exec subroutine. Also provides a list of object files loaded for the current process
<u>monitor</u>	Starts and stops execution profiling
<u>nlist</u>	Gets entries from a name list
<u>regcmp, regex</u>	Compile and match regular-expression patterns
<u>setjmp, longjmp</u>	Store a location

Subroutine	Description
<u>sgetl, sputl</u>	Accesses long numeric data in a machine-independent fashion
<u>sysconf</u>	Determines the current value of a specified system limit or option

List of files and directories subroutines

The system provides services to create files, move data into and out of files, and describe restrictions and structures of the file system.

Many of these subroutines are the base for the system commands that have similar names. You can, however, use these subroutines to write new commands or utilities to help in the program development process, or to include in an application program.

The system provides subroutines for:

Controlling files

access, accessx, or faccessx

Determine accessibility of a file

fclear

Clears space in a file

fcntl, dup, or dup2

Control open file descriptors

fsync

Writes changes in a file to permanent storage

getenv

Returns the value of an environment variable

getutent, getutid, getutline, pututline, setutent, endutent, or utmpname

Access utmp file entries

getutid_r, getutline_r, pututline_r, setutent_r, endutent_r, or utmpname_r

Access utmp file entries

lseek or llseek

Move the read-write pointer in an open file

lockfx, lockf, or flock

Controls open file descriptor locks

mknod or mkfifo

Create regular, FIFO, or special files

mktemp or mkstemp

Construct a unique file name

open, openx, or creat

Return a file descriptor and creates files

pclose

Closes an open pipe

pipe

Creates an interprocess channel

popen

Initiates a pipe to a process

pathconf, fpathconf

Retrieve file implementation characteristics

putenv

Sets an environment variable

read, readx, readv, readvx

Read from a file or device

rename

Renames directory or file within a file system

statx, stat, fstatx, fstat, fullstat, fullstat

Get file status

tmpfile

Creates a temporary file

tmpnam or tmpnam

Construct a name for a temporary file

truncate, ftruncate

Make a file shorter

umask

Gets and sets the value of the file creation mask

utimes or utime

Set file access or modification time

write, writex, writev, writevx

Write to a file or device

Working with directories**chdir**

Changes the current working directory

chroot

Changes the effective root directory

getwd, getcwd

Get the current directory path name

glob

Generates a list of path names to accessible files

globfree

Frees all memory associated with the *pglob* parameter

link

Creates additional directory entry for an existing file

mkdir

Creates a directory

opendir, readdir, telldir, seekdir, rewinddir, closedir

Performs operations on directories

readdir_r

Reads a directory

rmdir

Removes a directory

scandir, alphasort

Scan a directory

readlink

Reads the volume of a symbolic link

remove

Makes a file inaccessible by specified name

symlink

Creates a symbolic link to a file

unlink

Removes a directory entry

Manipulating file systems

confstr

Determines the current value of a specified system variable defined by a string

fscntl

Manipulates file system control operations

getfsent, getfsspec, getfsfile, getfstype, setfsent, or endfsent

Get information about a file system

getvfsent, getvfsbytype, getvfsbyname, getvfsbyflag, setvfsent, endvfsent

Get information about virtual file system entries

mnctl

Returns mount status information

quotactl

Manipulates disk quotas

statfs, fstatfs

Get the status of a file's file system

sysconf

Reports current value of system limits or options

sync

Updates all file systems information to disk

umask

Gets and sets the value of the file creation mask

vmount

Mounts a file system

umount, uvmount

Remove a virtual file system from the file tree

List of numerical manipulation subroutines

These functions perform numerical manipulation:

Function	Description
<u>a64l, l64a</u>	Convert between long integers and base-64 ASCII strings
<u>abs, div, labs, ldiv, imul_dbl, umul_dbl, llabs, lldiv</u>	Compute absolute value, division, and multiplication of integers
<u>asin, asinl, acos, acosl, atan, atanl, atan2, atan2l</u>	Compute inverse trigonometric functions
<u>asinh, acosh, atanh</u>	Compute inverse hyperbolic functions
<u>atof, atoff, strtod, strtold, strtodf</u>	Convert an ASCII string to a floating point number
<u>bessell: j0, j1, jn, y0, y1, yn</u>	Compute bessel functions
<u>class, finite, isnan, unordered</u>	Determine types of floating point functions
<u>copysign, nextafter, scalb, logb, ilogb</u>	Compute certain binary floating-point functions
<u>nrnd48, mrand48, jrand48, srand48, seed48, lcong48</u>	Generate pseudo-random sequences
<u>lrand48_r, mrand48_r, nrnd48_r, seed48_r, or srand48_r</u>	Generate pseudo-random sequences
<u>drem or remainder</u>	Compute an IEEE remainder
<u>ecvt, fcvt, gcvt</u>	Convert a floating-point number to a string

Function

erf, erfl, erfc, erfcl

exp, expl, expm1, log, logl, log10, log10l, log1p, pow, powl

floor, floorl, ceil, ceill, nearest,

trunc, rint, itrunc, uitrunc, fmod, fmodl, fabs, fabsl

fp_any_enable, fp_is_enabled, fp_enable_all,

fp_enable, fp_disable_all, fp_disable

fp_clr_flag, fp_set_flag, fp_read_flag, or fp_swap_flag

fp_invalid_op, fp_divbyzero, fp_overflow,

fp_underflow, fp_inexact, fp_any_xcp

fp_iop_snan, fp_iop_infsinf, fp_iop_infdfinf,

fp_iop_zrdzr, fp_iop_infmzr, fp_iop_invcmp

fp_read_rnd, fp_swap_rnd

frexp, frexpl, ldexp, ldexpl, modf, modfl

l64a_r

lgamma, lgammal, gamma

hypot, cabs

l3tol, ltol3

madd, msub, mult, mdiv, pow, gcd, invert,

rpow, msqrt, mcmp, move, min, omin,

fmin, m_in, mout, omout, fmout, m_out, sdiv, itom

rand, srand

rand_r

random, srandom, initstate, setstate

rsqrt

sin, cos, tan

sinh, sinhl, cosh, coshl, tanh, tanhl

sqrt, sqrtl, cbrt

strtol, strtoll, strtoul, strtoull, atol, atoi

Description

Compute error and complementary error functions

Compute exponential, log, and power functions

Round floating-point numbers

Allow operations on the floating-point exception status

Allow operations on the floating-point exception status

Test to see if a floating-point exception has occurred

Test to see if a floating-point exception has occurred

Read and set the IEEE rounding mode

Manipulate floating point numbers

Converts base-64 long integers to strings

Compute the logarithm of the gamma function

Compute Euclidean distance functions and absolute values

Convert between 3-byte integers and long integers

Provide multiple precision integer arithmetic

Generate random numbers

Generates random numbers

Generate better random numbers

Computes the reciprocal of the square root of a number

Compute trigonometric and inverse trigonometric functions

Computes hyperbolic functions

Compute square root and cube root functions

Convert a string to an integer

List of long long integer numerical manipulation subroutines

The following subroutines perform numerical manipulation of integers stored in the long long integer data format:

Subroutine	Description
<u>llabs</u>	Computes the absolute value of a long long integer
<u>lldiv</u>	Computes the quotient and remainder of the division of two long long integers
<u>strtoll</u>	Converts a string to a signed long long integer
<u>strtoull</u>	Converts a string to an unsigned long long integer
<u>wcstoll</u>	Converts a wide character string to a signed long long integer
<u>wcstoull</u>	Converts a wide character string to an unsigned long long integer

List of 128-bit long double numerical manipulation subroutines

The following subroutines perform numerical manipulation of floating-point numbers stored in the 128-bit long double data type.

These subroutines do not support the 64-bit long double data type. Applications that use the 64-bit long double data type should use the corresponding double-precision subroutines.

Subroutine	Description
<u>acosl</u>	Computes the inverse cosine of a floating-point number in long double format
<u>asinl</u>	Computes the inverse sine of a floating-point number in long double format
<u>atan2l</u>	Computes the principal value of the arc tangent of x/y , whose components are expressed in long double format
<u>atanl</u>	Computes the inverse tangent of a floating-point number in long double format
<u>ceil</u>	Computes the smallest integral value not less than a specified floating-point number in long double format
<u>coshl</u>	Computes the hyperbolic cosine of a floating-point number in long double format
<u>cosl</u>	Computes the cosine of a floating-point number in long double format
<u>erfcl</u>	Computes the value of 1 minus the error function of a floating-point number in long double format
<u>erfl</u>	Computes the error function of a floating-point number in long double format
<u>expl</u>	Computes the exponential function of a floating-point number in long double format
<u>fabsl</u>	Computes the absolute value of a floating-point number in long double format
<u>floorl</u>	Computes the largest integral value not greater than a specified floating-point number in long double format
<u>fmodl</u>	Computes the long double remainder of a fraction x/y , where x and y are floating-point numbers in long double format
<u>frexpl</u>	Expresses a floating-point number in long double format as a normalized fraction and an integral power of 2, storing the integer and returning the fraction
<u>ldexpl</u>	Multiplies a floating-point number in long double format by an integral power of 2
<u>lgammal</u>	Computes the natural logarithm of the absolute value of the gamma function of a floating-point number in long double format
<u>log10l</u>	Computes the base 10 logarithm of a floating-point number in long double format

Subroutine	Description
<u>logl</u>	Computes the natural logarithm of a floating-point number in long double format
<u>modfl</u>	Stores the integral part of a real number in a long double variable and returns the fractional part of the real number
<u>powl</u>	Computes the value of x raised to the power of y , where both numbers are floating-point numbers in long double format
<u>sinhl</u>	Computes the hyperbolic sine of a floating-point number in long double format
<u>sinl</u>	Computes the sine of a floating-point number in long double format
<u>sqrtl</u>	Computes the square root of a floating-point number in long double format
<u>strtold</u>	Converts a string to a floating-point number in long double format
<u>tanl</u>	Computes the tangent of a floating-point number in long double format
<u>tanh</u>	Computes the hyperbolic tangent of a floating-point number in long double format

List of processes subroutines

With the introduction of threads, several process subroutines have been extended and other subroutines have been added. Threads, not processes, are now the schedulable entity.

For signals, the handler exists at the process level, but each thread can define a signal mask. Some examples of changed or new subroutines are: **getprocs**, **getthrds**, **ptrace**, **getpri**, **setpri**, **yield** and **sigprocmask**.

The subroutines are listed in the following categories:

Process initiation

exec:, execl, execv, execle, execve, execlp, execvp, or exec

Execute new programs in the calling process

fork or vfork

Create a new process

reboot

Restarts the system

siginterrupt

Sets subroutines to restart when they are interrupted by specific signals

Process suspension

pause

Suspends a process until that process receives a signal

wait, wait3, waitpid

Suspend a process until a child process stops or terminates

Process termination

abort

Terminates current process and produces a memory dump by sending a **SIGOT** signal

exit, atexit, or _exit

Terminate a process

, unatexit,

Unregisters functions that were previously registered by the `atexit` subroutine. If the referenced function is found, it is removed from the list of functions that are called at normal program termination.

kill or killpg

Terminate current process or group of processes with a signal

Process and thread identification

ctermid

Gets the path name for the terminal that controls the current process

cuserid

Gets the alphanumeric user name associated with the current process

getpid, getpgid, or getppid

Get the process ID, process group ID, or the parent process ID, respectively

getprocs

Gets process table entries

getthrds

Gets thread table entries

setpgid or setpgrp

Set the process group ID

setsid

Creates a session and sets process group IDs

uname or unamex

Gets the names of the current operating system

Process accounting

acct

Enables and disables process accounting

ptrace

Traces the execution of a process

Process resource allocation

brk or sbrk

Change data segment space allocation

getdtablesize

Gets the descriptor table size

getrlimit, setrlimit, or vlimit

Limit the use of system resources by current process

getrusage, times, or vtimes

Display information about resource use

plock

Locks processes, text, and data into memory

profil

Starts and stops program address sampling for execution profiling

ulimit

Sets user process limits

Process prioritization

getpri

Returns the scheduling priority of a process

getpriority, setpriority, or nice

Get or set the priority value of a process

setpri

Sets a process scheduling priority to a constant value

yield

Yields the processor to processes with higher priorities

Process and thread synchronization

compare_and_swap

Conditionally updates or returns a single word variable atomically

fetch_and_add

Updates a single word variable atomically

fetch_and_and and fetch_and_or

Sets or clears bits in a single word variable atomically

semctl

Controls semaphore operations

semget

Gets a set of semaphores

semop

Performs semaphore operations

Process signals and masks**raise**

Sends a signal to an executing program

sigaction, sigvec, or signal

Specifies the action to take upon delivery of a signal

sigemptyset, sigfillset, sigaddset, sigdelset, or sigismember

Create and manipulate signal masks

sigpending

Determines the set of signals that are blocked from delivery

sigprocmask, sigsetmask, or sigblock

Set signal masks

sigset, sighold, sigrelse, or sigignore

Enhance the signal facility and provide signal management

sigsetjmp or siglongjmp

Save and restore stack context and signal masks

sigstack

Sets signal stack context

sigsuspend

Changes the set of blocked signals

ssignal or gsignal

Implement a software signal facility

Process messages**msgctl**

Provides message control operations

msgget

Displays a message queue identifier

msgrcv

Reads messages from a queue

msgsnd

Sends messages to the message queue

msgxrcv

Receives an extended message

psignal

Printing system signal messages

List of multithreaded programming subroutines

Programming in a multithreaded environment requires reentrant subroutines to ensure data integrity.

Use the following subroutines rather than the non-reentrant version:

Subroutine	Description
<u>asctime_r</u>	Converts a time value into a character array
<u>getgrnam_r</u>	Returns the next group entry in the user database that matches a specific name
<u>getpwuid_r</u>	Returns the next entry that matches a specific user ID in the use database

The following lists the non-reentrant subroutines in libc.

Subroutine	Description		
asctime	getgrent	gsignal	setkey
auditread	getgrgid	hcreate	setlogmask
closelog	getgrnam	hdestroy	setnetent
crypt	getgroupsbyuser	hsearch	setnetgrent
ctime	getgrset	inet_ntoa	setprotoent
dirname	gethostbyaddr	initstate	setpwent
drand48	gethostbyname	inetgr	setpwfile
ecvt	gethostent	iso_addr	setrpcent
endttyent	getlogin	iso_ntoa	setservent
encrypt	getnetbyaddr	jrand48	setstate
asctime	getnetbyname	l64a	setttyent
endfsent	getnetent	lcong48	setutent
endfsent	getnetgrent	link_ntoa	setutxent
endgrent	getopt	localtime	srand48
endhostent	getprotobyname	lrand48	srandom
endnetent	getprotobynumber	mrand48	ssignal
endnetgrent	getprotoent	mtime	strerror
endprotoent	getpwent	ndutent	strtok
endpwent	getpwnam	nrand48	syslog
endrpcent	getpwuid	ns_ntoa	ttyname
endservent	getrpcbyname	openlog	utmpname
endttyent	getrpcbynumber	pututline	wcstok
endutxent	getrpcent	pututxline	
erand48	getservbyname	rand	
ether_aton	getservbyport	random	
ether_ntoa	getservent	rcmd	
fcvt fgetgrent	getttyent	rcmd2	

Subroutine	Description		
fgetpwent	getttynam	readdir	
getdate	getuinfo	rexec	
getfsent	getutent	re_comp	
getfsent	getutid	re_exec	
getfsfile	getutline	seed48	
getfsfile	getutxent	setfsent	
getfsspec	getutxid	setgrent	
getfstype	getutxline	sethostent	

List of programmer's workbench library subroutines

The Programmers Workbench Library (**libPW.a**) contains routines that are provided only for compatibility with existing programs.

Their use in new programs is not recommended. These interfaces are from AT&T PWB Toolchest.

Routine	Description
any (<i>Character, String</i>)	Determines whether <i>String</i> contains <i>Character</i>
anystr (<i>String1, String2</i>)	Determines the offset in <i>String1</i> of the first character that also occurs in <i>String2</i>
balbrk (<i>String, Open, Close, End</i>)	Determines the offset in <i>String</i> of the first character in the string <i>End</i> that occurs outside of a balanced string as defined by <i>Open</i> and <i>Close</i>
cat (<i>Destination, Source1, Source0</i>)	Concatenates the <i>Source</i> strings and copies them to <i>Destination</i>
clean_up ()	Defaults the cleanup routine
curdir (<i>String</i>)	Puts the full path name of the current directory in <i>String</i>
dname (<i>p</i>)	Determines which directory contains the file <i>p</i>
fatal (<i>Message</i>)	General purpose error handler
fdfopen (<i>fd, Mode</i>)	Same as the stdio fdfopen subroutine
giveup (<i>Dump</i>)	Forces a core dump
imatch (<i>pref, String</i>)	Determines if the string <i>pref</i> is an initial substring of <i>String</i>
lockit (<i>LockFile, Count, pid</i>)	Creates a lock file
move (<i>String1, String2, n</i>)	Copies the first <i>n</i> characters of <i>String1</i> to <i>String2</i>
patoi (<i>String</i>)	Converts <i>String</i> to integer
patol (<i>String</i>)	Converts <i>String</i> to long.
repeat (<i>Destination, String, n</i>)	Sets <i>Destination</i> to <i>String</i> repeated <i>n</i> times

Routine

repl (*String, Old, New*)

satoi (*String, *ip*)

setsig ()

setsig1 (*Signal*)

sname (*String*)

strend (*String*)

trnslat (*s, old, new, Destination*)

unlockit (*lockfile, pid*)

userdir (*uid*)

userexit (*code*)

username (*uid*)

verify (*String1, String2*)

xalloc (*asize*)

xcreat (*name, mode*)

xfree (*aptr*)

xfreeall ()

xlink (*f1, f2*)

xmsg (*file, func*)

xpipe (*t*)

xunlink (*f*)

xwrite (*fd, buffer, n*)

zero (*p, n*)

zeropad (*s*)

Description

Replaces each occurrence of the character *Old* in *String* with the character *New*

Converts *String* to integer and saves it in **ip*

Causes signals to be caught by **setsig1**

General purpose signal handling routine

Gets a pointer to the simple name of full path name *String*

Finds the end of *String*.

Copies string *s* into *Destination* and replace any character in *old* with the corresponding characters in *new*

Deletes the lock file

Gets the user's login directory

Defaults user exit routine

Gets the user's login name

Determines the offset in string *String1* of the first character that is not also in string *String2*

Allocates memory

Creates a file

Frees memory

Frees all memory

Links files

Calls the routine **fatal** with an appropriate error message

Creates a pipe

Removes a directory entry

Writes *n* bytes to the file associated with *fd* from *buffer*

Zeros *n* bytes starting at address *p*

Replaces the initial blanks with the character 0 (zero) in string *s*

File

/usr/lib/libPW.a

Contains routines provided only for compatibility with existing programs

List of security and auditing subroutines

This section lists security and auditing subroutines.

Access control subroutines

Subroutine	Description
<u>acl_chg</u> or <u>acl_fchg</u>	Change the access control information on a file
<u>acl_get</u> or <u>acl_fget</u>	Get the access control information of a file
<u>acl_put</u> or <u>acl_fput</u>	Set the access control information of a file
<u>acl_set</u> or <u>acl_fset</u>	Set the base entries of the access control information of a file
<u>aclx_convert</u>	Convert the access control information from one ACL type to another
<u>aclx_get</u> or <u>aclx_fget</u>	Get the access control information of a file if the ACL associated is of the AIXC type
<u>aclx_gettypeinfo</u>	Retrieve the ACL characteristics given to an ACL type
<u>aclx_gettypes</u>	Retrieve the list of ACL types supported for the file system associated with the path provided
<u>aclx_print</u> or <u>aclx_printStr</u>	Convert the binary access control information into nonbinary, readable format
<u>aclx_put</u> or <u>aclx_fput</u>	Stores the access control information for a file system object
<u>aclx_scan</u> or <u>aclx_scanStr</u>	Convert the access control information that is in nonbinary, readable text format into ACL type-specific native format binary ACL data
<u>chacl</u> or <u>fchac l</u>	Change the permissions on a file
<u>chmod</u> or <u>fchmod</u>	Change file access permissions
<u>chown</u> , <u>fchown</u> , <u>chownx</u> , or <u>fchownx</u>	Change file ownership
<u>frevoke</u>	Revoke access to a file by other processes
<u>revoke</u>	Revoke access to a file
<u>statacl</u> or <u>fstatacl</u>	Retrieve the access control information for a file

Auditing subroutines

Subroutine	Description
<u>audit</u>	Enables and disables system auditing
<u>auditbin</u>	Defines files to contain audit records
<u>auditevents</u>	Gets or sets the status of system event auditing
<u>auditlog</u>	Appends an audit record to an audit bin file
<u>auditobj</u>	Gets or sets the auditing mode of a system data object
<u>auditpack</u>	Compresses and uncompresses audit bins
<u>auditproc</u>	Gets or sets the audit state of a process
<u>auditread</u> or <u>auditread_r</u>	Read an audit record
<u>auditwrite</u>	Writes an audit record

Identification and authentication subroutines

User authentication routines have a potential to store passwords and encrypted passwords in memory. This may expose passwords and encrypted passwords in coredumps.

Subroutine	Description
<u>authenticate</u>	Authenticates the user's name and password

Subroutine

[ckuseracct](#)

[ckuserID](#)

[crypt](#), [encrypt](#), or [setkey](#)

[genpagvalue](#)

[getpagvalue64](#)

[setpagvalue64](#)

[getgrent](#), [getgrgid](#), [getgrnam](#), [setgrent](#), or [endgrent](#)

[getgrgid_r](#)

[getgrnam_r](#)

[getgroupattr](#), [IDtogroup](#), [nextgroup](#), or [putgroupattr](#)

[getlogin](#)

[getlogin_r](#)

[getpass](#)

[getportattr](#) or [putportattr](#)

[getpwent](#), [getpwuid](#), [getpwnam](#), [putpwent](#), [setpwent](#), or [endpwent](#)

[getuinfo](#)

[getuserattr](#), [IDtouser](#), [nextuser](#), or [putuserattr](#)

[getuserpw](#), [putuserpw](#), or [putuserpwhist](#)

[loginfailed](#)

[loginrestrictions](#)

[loginsuccess](#)

[newpass](#)

[passwdexpired](#)

[setpwdb](#) or [endpwdb](#)

[setuserdb](#) or [enduserdb](#)

[system](#)

[tcb](#)

Process subroutines

Subroutine

[getgid](#) or [getegid](#)

Description

Checks the validity of a user account

Authenticates the user

Encrypt or decrypt data

Generates a system-wide unique PAG value for a given PAG name such as afs.

Retrieves 64-bit PAG values for a process.

Stores 64-bit PAG values for a process. .

Accesses the basic group information in the user database

Gets a group database entry for a group ID in a multithreaded environment

Searches a group database for a name in a multithreaded environment

Accesses the group information in the user database

Gets the user's login name

Gets the user's login name in a multithreaded environment

Reads a password

Access the port information in the port database

Access the basic user information in the user database

Finds the value associated with a user

Access the user information in the user database

Access the user authentication data

Records an unsuccessful login attempt

Determines if a user is allowed to access the system

Records a successful login

Generates a new password for a user

Checks the user's password to determine if it has expired

Open or close the authentication database

Open or close the user database

Runs a shell command

Alters the Trusted Computing Base status of a file

Description

Get the real or group ID of the calling process

Subroutine	Description
<u>getgroups</u>	Gets the concurrent group set of the current process
<u>getpcred</u>	Gets the current process security credentials
<u>getpenv</u>	Gets the current process environment
<u>getuid</u> or <u>geteuid</u>	Get the real or effective user ID of the current process
<u>initgroups</u>	Initializes the supplementary group ID of the current process
<u>kleenup</u>	Cleans up the run-time environment of a process
<u>setgid</u>, <u>setrgid</u>, <u>setegid</u>, or <u>setregid</u>	Set the group IDs of the calling process
<u>setgroups</u>	Sets the supplementary group ID of the current process
<u>setpcred</u>	Sets the current process credentials
<u>setpenv</u>	Sets the current process environment
<u>setuid</u>, <u>setruid</u>, <u>setuid</u>, or <u>setreuid</u>	Set the process user IDs
<u>usrinfo</u>	Gets and sets user information about the owner of the current process

List of string manipulation subroutines

The string manipulation functions include:

The string manipulation functions include:

- Locate a character position within a string
- Locate a sequence of characters within a string
- Copy a string
- Concatenate strings
- Compare strings
- Translate a string
- Measure a string

When using these string functions, you do not need to include a header file for them in the program or specify a special flag to the compiler.

The following functions manipulate string data:

bcopy, bcmp, bzero, ffs

Perform bit and byte string operations

gets, fgets

Get a string from a stream

puts, fputs

Write a string to a stream

compile, step, advance

Compile and match regular-expression patterns

strlen, strchr, strrchr, strpbrk, strspn, strcspn, strstr, strtok

Perform operations on strings

jcode

Performs string conversion on 8-bit processing codes.

varargs

Handles a variable-length parameter list

Example: Program for manipulating characters

This section includes a programming example for manipulating characters.

```
/*  
This program is designed to demonstrate the use of "Character  
classification and conversion" subroutines. Since we are dealing  
with characters, it is a natural place to demonstrate the use of  
getchar subroutine and putchar subroutine from the stdio library.
```

The program objectives are:

-Read input from "stdin"

-Verify that all characters are ascii and printable

-Convert all uppercase characters to lowercase

-Discard multiple white spaces

-Report statistics regarding the types of characters

The following routines are demonstrated by this example program:

- getchar
- putchar
- isascii (ctype)
- iscntrl (ctype)
- isspace (ctype)
- isalnum (ctype)
- isdigit (ctype)
- isalpha (ctype)
- isupper (ctype)
- islower (ctype)
- ispunct (ctype)
- tolower (conv)
- toascii (conv)

```
*/
```

```
#include <stdio.h> /* The mandatory include file */  
#include <ctype.h> /* Included for character classification
```

```
subroutines */
```

```

/* The various statistics gathering counters */

int asciicnt, printcnt, punctcnt, uppercnt, lowercnt,
digcnt, alnumcnt, cntrlcnt, spacecnt, totcnt, nonprintcnt, linecnt, tabcnt ;

main()
{

int ch ; /* The input character is read in to this */
char c , class_conv() ;

asciicnt=printcnt=punctcnt=uppercnt=lowercnt=digcnt==0;
cntrlcnt=spacecnt=totcnt=nonprintcnt=linecnt=tabcnt=0;
alnumcnt=0;

while ( (ch =getchar()) != EOF )
{

totcnt++;
c = class_conv(ch) ;
putchar(c);

}
printf("The number lines of of input were %d\n",linecnt);
printf(" The character wise breakdown follows :\n");
printf(" TOTAL  ASCII  CNTRL  PUNCT  ALNUM  DIGITS  UPPER
LOWER  SPACE  TABCNT\n");

printf("%5d %5d %5d %5d %5d %5d %5d %5d %5d %5d\n",totcnt,
asciicnt, cntrlcnt, punctcnt, alnumcnt, digcnt, uppercnt,lowercnt, spacecnt, tabcnt );

}

char class_conv(ch)
char ch;
{

if (isascii(ch)) {

asciicnt++;
if ( iscntrl(ch) && ! isspace(ch)) {

nonprintcnt++ ;
cntrlcnt++ ;
return(' ');
}
}
}

```

```

}
else if ( isalnum(ch)) {

alnumcnt++;
if (isdigit(ch)){

digcnt++;
return(ch);

}
else if (isalpha(ch)){

if ( isupper(ch) ){

uppercnt++ ;
return(tolower(ch));

}
else if ( islower(ch) ){

lowercnt++;
return(ch);

}
else {

/*
We should never be in this situation since an alpha character can only be
either uppercase or lowercase.
*/

fprintf(stderr,"Classification error for %c \n",ch);
return(NULL);

}

}
else if (ispunct(ch) ){

punctcnt++;
return(ch);

}
else if ( isspace(ch) ){

spacecnt++;
if ( ch == '\n' ){
linecnt++;
return(ch);

}
while ( (ch == '\t' ) || ( ch == ' ' ) ) {
if ( ch == '\t' ) tabcnt ++ ;
else if ( ch == ' ' ) spacecnt++ ;
totcnt++;
ch = getchar();

}
ungetc(ch,stdin);
totcnt--;
return(' ');

```

```

}
else {

/*
We should never be in this situation any ASCII character
can only belong to one of the above classifications.
*/
fprintf(stderr,"Classification error for %c \n",ch);
return(NULL);
}

}
else
{

fprintf(stdout,"Non Ascii character encountered \n");
return(toascii(ch));

}

}
}

```

Example: Searching and sorting program

This section includes an example of a searching and sorting program.

```

/**This program demonstrates the use of the following:

-qsort subroutine (a quick sort library routine)

-bsearch subroutine (a binary search library routine)

-fgets, fopen, fprintf, malloc, sscanf, and strcmp subroutines.

The program reads two input files with records in
string format, and prints or displays:

-records from file2, which are excluded in file1

-records from file1, which are excluded in file2

The program reads the input records from both files
into two arrays, which are subsequently sorted in
common order using the qsort subroutine. Each element of
one array is searched for its counterpart entry in the
other array using the bsearch subroutine. If the item is
not found in both arrays, a message indicates the record
was not found. The process is repeated interchanging
the two arrays, to obtain the second list of exclusions.

**/

#include <stdio.h>      /*the library file to be included for
                       /*standard input and output*/

#include <search.h>     /*the file to be included for qsort*/
#include <sys/errno.h>  /*the include file for interpreting

```

```

/*predefined error conditions*/

#define MAXRECS 10000 /*array size limit*/

#define MAXSTR 256 /*maximum input string length*/
#define input1 "file1" /*one input file*/
#define input2 "file2" /*second input file*/
#define out1 "o_file1" /*output file1*/
#define out2 "o_file2" /*output file2*/

main()
{

char *arr1[MAXRECS] , *arr2[MAXRECS] ;/*the arrays to store
input records*/

unsigned int num1 , num2; /*to keep track of the number of
/*input records. Unsigned int
/*declaration ensures
/*compatibility
/*with qsort library routine.*/

int i ;
int compar(); /*the function used by qsort and
/*bsearch*/

extern int errno ; /*to capture system call failures*/
FILE *ifp1 , *ifp2, *ofp1, *ofp2; /*the file pointers for
input and output */

void *bsearch() ; /*the library routine for binary search*/
void qsort(); /*the library routine for quick sort*/
char*malloc() ; /*memory allocation subroutine*/
void exit() ;

num1 = num2 = 0;

/**Open the input and output files for reading or writing
**/

if ( (ifp1 = fopen( input1 , "r" ) ) == NULL )
{

(void) fprintf(stderr,"%s could not be opened\n",input1);
exit(-1);

}

if ( ( ifp2 = fopen( input2 , "r" ) ) == NULL )
{

```



```

(void) fprintf(stderr,"%s could not be opened\n",input2);
exit(-1);

}

if ( ( ofp1 = fopen(out1,"w" ) ) == NULL )
{

(void) fprintf(stderr,"%s could not be opened\n",out1);
exit(-1);

}

if ( ( ofp2 = fopen(out2,"w") ) == NULL )
{

(void) fprintf(stderr,"%s could not be opened\n", out2);
exit(-1);

}

/**Fill the arrays with data from input files. Readline
function returns the number of input records.**/

if ( ( i = readline( arr1 , ifp1 ) ) < 0 )
{

(void) fprintf(stderr,"o data in %s. Exiting\n",input1);
exit(-1);

}
num1 = (unsigned) i;
if ( ( i = readline ( arr2 , ifp2) ) < 0 )
{

(void) fprintf(stderr,"No data in %s. Exiting\n",input2);
exit(-1);

}
num2 = (unsigned) i;

/**
The arrays can now be sorted using qsort subroutine
**/

qsort( (char *)arr1 , num1 , sizeof (char * ) , compar);
qsort( (char *)arr2 , num2 , sizeof (char * ) , compar);

/**When the two arrays are sorted in a common order, the
program builds a list of elements found in one but not
in the other, using bsearch.

Check that each element in array1 is in array2
**/

for ( i= 0 ; i < num1 ; i++ )
{

if ( bsearch((void *)&arr1[i] , (char *)arr2,num2,

sizeof(char * ) , compar) == NULL )

```

```

{
(void) fprintf(ofp1,"%s",arr1[i]);

}

} /**One list of exclusions is complete**/

/**Check that each element in array2 is in array1**/

for ( i = 0 ; i < num2 ; i++ )
{
if ( bsearch((void *)&arr2[i], (char *)arr1, num1
, sizeof(char *) , compar) == NULL )

{
(void) fprintf(ofp2,"%s",arr2[i]);

}

}

/**Task completed, so return**/

return(0);

}

/**The function reads in records from an input
file and fills in the details into the two arrays.**/

readline ( char **aptr, FILE *fp )
{
char str[MAXSTR] , *p ;
int i=0 ;

/**Read the input file line by line**/

while ( fgets(str , sizeof(str) , fp ))
{

/**Allocate sufficient memory. If the malloc subroutine
fails, exit.**/

if ( ( p = (char *)malloc ( sizeof(str))) == NULL )
{

(void) fprintf(stderr,"Insufficient Memory\n");
return(-1);

}
else

{
if ( 0 > strcpy(p, str))

{
(void) fprintf(stderr,"Strcpy failed \n");
return(-1);
}
}
}
}

```

```

}
i++ ; /*increment number of records count*/

}

} /**End of input file reached**/
return(i);/*return the number of records read*/

}

/**We want to sort the arrays based only on the contents of the first field of
the input records. So we get the first field using SSCANF**/

compar( char **s1 , char **s2 )
{

char st1[100] , st2[100] ;
(void) sscanf(*s1,"%s" , st1) ;
(void) sscanf(*s2,"%s" , st2) ;

/**Return the results of string comparison to the calling procedure**/

return(strcmp(st1 , st2));

}

```

List of operating system libraries

This section lists operating system libraries.

Library	Description
/usr/lib/libbsd.a	Berkeley library
/lib/profiled/libbsd.a	Berkeley library profiled
/usr/ccs/lib/libcurses.a	Curses library
/usr/ccs/lib/libc.a	Standard I/O library, standard C library
/lib/profiled/libc.a	Standard I/O library, standard C library profiled
/usr/ccs/lib/libdbm.a	Database Management library
/usr/ccs/lib/libi18n.a	Layout library
/usr/lib/liblvm.a	LVM (Logical Volume Manager) library
/usr/ccs/lib/libm.a	Math library
/usr/ccs/lib/libp/libm.a	Math library profiled
/usr/lib/libodm.a	ODM (Object Data Manager) library
/usr/lib/libPW.a	Programmers Workbench library
/usr/lib/libpthread.a	POSIX compliant Threads library
/usr/lib/libqb.a	Queue Backend library
/usr/lib/librpcsvc.a	RPC (Remote Procedure Calls) library
/usr/lib/librts.a	Run-Time Services library
/usr/lib/libsa.a	Security functions
/usr/lib/libsm.a	System management library

Library	Description
<code>/usr/lib/libsrc.a</code>	SRC (System Resource Controller) library
<code>/usr/lib/libmsaa.a</code>	SVID (System V Interface Definition) math library
<code>/usr/ccs/lib/lib/libmsaa.a</code>	SVID (System V Interface Definition) math library profiled
<code>/usr/ccs/lib/libtermcap.a</code>	Terminal I/O
<code>/usr/lib/liby.a</code>	YP (Yellow Pages) library
<code>/usr/lib/lib300.a</code>	Graphics subroutines for DASI 300 workstations
<code>/usr/lib/lib300s.a</code>	Graphics subroutines for DASI 300s workstations
<code>/usr/lib/lib300S.a</code>	Graphics subroutines for DASI 300S workstations
<code>/usr/lib/lib4014.a</code>	Graphics subroutines for Tektronix 4014 workstations
<code>/usr/lib/lib450.a</code>	Graphics subroutines for DASI 450 workstations
<code>/usr/lib/libcsys.a</code>	Kernel extensions services
<code>/usr/ccs/lib/libdbx.a</code>	Debug program library
<code>/usr/lib/libgsl.a</code>	Graphics Support library
<code>/usr/lib/libieee.a</code>	IEEE floating point library
<code>/usr/lib/libIM.a</code>	Stanza file processing library
<code>/usr/ccs/lib/libl.a</code>	lex library
<code>/usr/lib/libogsl.a</code>	Old graphics support library
<code>/usr/lib/liboldX.a</code>	X10 library
<code>/usr/lib/libplot.a</code>	Plotting subroutines
<code>/usr/lib/librpcsvc.a</code>	RPC services
<code>/usr/lib/librs2.a</code>	Hardware-specific sqrt and itrunc subroutines
<code>/usr/lib/libxgsl.a</code>	Enhanced X-Windows graphics subroutines
<code>/usr/lib/libX11.a</code>	X11 run time library
<code>/usr/lib/libXt.a</code>	X11 toolkit library
<code>/usr/lib/liby.a</code>	yacc run time library

System Management Interface Tool (SMIT)

The System Management Interface Tool (SMIT) is an interactive and extensible screen-oriented command interface.

It prompts users for the information needed to construct command strings and presents appropriate predefined selections or run time defaults where available. This shields users from many sources of extra work or error, including the details of complex command syntax, valid parameter values, system command spelling, or custom shell path names.

You can also build and use alternate databases instead of modifying SMIT's default system database.

The following sections discuss SMIT in detail:

New tasks consisting of one or more commands or inline **ksh** shell scripts can be added to SMIT at any time by adding new instances of predefined screen objects to SMIT's database. These screen objects (described by stanza files) are used by the Object Data Manager (ODM) to update SMIT's database. This database controls SMIT's run-time behavior.

Related information

[dspmsg](#)
[gencat](#)
[ksh](#)
[man](#)
[odmadd](#)
[odmcreate](#)
[odmget](#)
[smit](#)
[spaths](#)

SMIT screen types

There are three main screen types available for the System Management Interface Tool (SMIT). The screens occur in a hierarchy consisting of menu screens, selector screens, and dialog screens.

When performing a task, a user typically traverses one or more menus, then zero or more selectors, and finally one dialog.

The following table shows SMIT screen types, what the user sees on each screen, and what SMIT does internally with each screen:

Screen type	What the user sees on the screen	What SMIT does internally with each screen
Menu	A list of choices	Uses the choice to select the next screen to display.
Selector	Either a list of choices or an entry field	Obtains a data value for subsequent screens. Optionally selects alternative dialogs or selectors.
Dialog	A sequence of entry fields.	Uses data from the entry fields to construct and run the target task command string.

Menus present a list of alternative subtasks; a selection can then lead to another menu screen or to a selector or dialog screen. A selector is generally used to obtain one item of information that is needed by a subsequent screen and which can also be used to select which of several selector or dialog screens to use next. A dialog screen is where any remaining input is requested from the user and where the chosen task is actually run.

A menu is the basic entry point into SMIT and can be followed by another menu, a selector, or a dialog. A selector can be followed by a dialog. A dialog is the final entry panel in a SMIT sequence.

Menu screens

A SMIT menu is a list of user-selectable items. Menu items are typically tasks or classes of tasks that can be performed from SMIT. A user starting with the main SMIT menu selects an item defining a broad range of system tasks. A selection from the next and subsequent menus progressively focuses the user's choice, until finally a dialog is typically displayed to collect information for performance of a particular task.

Design menus to help a user of SMIT narrow the scope of choice to a particular task. Your design can be as simple as a new menu and dialog attached to an existing branch of SMIT, or as complex as an entire new hierarchy of menus, selectors, and dialogs starting at the SMIT applications menu.

At run time, SMIT retrieves all menu objects with a given ID (**id** descriptor value) from the specified object repository. To add an item to a particular SMIT menu, add a menu object having an ID value equal to the value of the **id** descriptor of other non-title objects in the same menu.

Build menus by defining them in a stanza file and then processing the file with the **odmadd** command. A menu definition is compiled into a group of menu objects. Any number of menus, selectors, and dialogs can be defined in one or more files.

Command	Description
odmadd	Adds the menu definitions to the specified object repository.
/usr/lib/objrepos	Default object repository for system information and can be used to store your compiled objects.

At SMIT run time, the objects are automatically retrieved from a SMIT database.

Note: You should always back up the **/usr/lib/objrepos** directory before deleting or adding any objects or object classes. Unanticipated damage to objects or classes needed for system operations can cause system problems.

Selector screens

A SMIT selector prompts a user to specify a particular item, typically a system object (such as a printer) or attribute of an object (such as a serial or parallel printer mode). This information is then generally used by SMIT in the next dialog.

For instance, a selector can prompt a user to enter the name of a logical volume for which to change logical volume characteristics. This could then be used as a parameter in the `sm_cmd_hdr.cmd_to_discover_postfix` field of the next dialog for entry field initialization. Likewise, the selector value could also be used as the value for a subsequent `sm_cmd_opt.cmd_to_list_postfix` field. It can also be used directly as a subsequent initial entry field value. In each case, logical consistency requires that this item either be selected prior to the dialog or be held constant while in the dialog.

Design a selector to request a single piece of information from the user. A selector, when used, falls between menus and dialogs. Selectors can be strung together in a series to gather several pieces of information before a dialog is displayed.

Selectors should usually contain a prompt displayed in user-oriented language and either a response area for user input or a pop-up list from which to select a value; that is, one question field and one answer. Typically the question field is displayed and the SMIT user enters a value in the response area by typing the value or by selecting a value from a list or an option ring.

To give the user a run-time list of choices, the selector object can have an associated command (defined in the `sm_cmd_opt.cmd_to_list` field) that lists the valid choices. The list is not hard-coded, but developed by the command in conjunction with standard output. The user gets this list by selecting the **F4 (Esc+4)=List** function of the SMIT interface.

In a ghost selector (`sm_cmd_hdr.ghost="y"`), the command defined in the `sm_cmd_opt.cmd_to_list` field, if present, is automatically run. The selector screen is not displayed at this time and the user sees only the pop-up list.

The application of a super-ghost selector permits branching following menu selection, where the branch to be taken depends on the system state and not user input. In this case, the **cmd_to_classify** descriptor in the super-ghost selector can be used to get the required information and select the correct screen to present next.

Build selectors by defining them in a stanza file and then processing the file with the **odmadd** command. Several menus, selectors, and dialogs can be defined in a single file. The **odmadd** command adds each selector to the specified object repository. The **/usr/lib/objrepos** directory is the default object repository for system information and is used to store your compiled objects. At SMIT run time, the objects are automatically retrieved from a SMIT database.

Note: Always back up the **/usr/lib/objrepos** directory before deleting or adding any objects or object classes. Unanticipated damage to objects or classes needed for system operations can cause system problems.

Dialog screens

A dialog in SMIT is the interface to a command or task a user performs. Each dialog executes one or more commands, shell functions, and so on. A command can be run from any number of dialogs.

To design a dialog, you need to know the command string you want to build and the command options and operands for which you want user-specified values. In the dialog display, each of these command options and operands is represented by a prompt displayed in user-oriented language and a response area for user input. Each option and operand is represented by a dialog command option object in the Object Data Manager (ODM) database. The entire dialog is held together by the dialog header object.

The SMIT user enters a value in the response area by typing the value, or by selecting a value from a list or an option ring. To give the user a run-time list of choices, each dialog object can have an associated command that lists the valid choices. The associated commands are defined in the `sm_cmd_opt.cmd_to_list` field. The user gets this list by invoking the **F4 (Esc + 4)=List** function of the SMIT interface. This causes SMIT to run the command defined in the associated `cmd_to_list` field and to use its standard output and `stderr` file for developing the list.

Dialog screens can have the following options assigned to the entries:

Option	Function
#	Signifies that a numerical value is expected.
*	Signifies that an entry is mandatory.
+	Signifies that a listing of choices can be obtained using the F4 key.

In a ghost dialog, the dialog screen is not displayed. The dialog runs as if the user had immediately pressed the dialog screen **Enter** key to run the dialog.

Build dialogs by defining them in a stanza file and then processing the file with the **odmadd** command. Several menus, selectors, and dialogs can be defined in a single file. The **odmadd** command adds each dialog definition to the specified object repository. The `/usr/lib/objrepos` directory is the default object repository for system information and can be used to store your compiled objects. At SMIT run time, the objects are automatically retrieved from a SMIT database.

Note: Always back up the `/usr/lib/objrepos` directory before deleting or adding any objects or object classes. Unanticipated damage to objects or classes needed for system operations can cause system problems.

SMIT object classes

A System Management Interface Tool (SMIT) object class created with the Object Data Manager (ODM) defines a common format or record data type for all individual objects that are instances of that object class.

Therefore a SMIT object class is basically a record data type and a SMIT object is a particular record of that type.

SMIT menu, selector, and dialog screens are described by objects that are instances of one of four object classes:

- **sm_menu_opt**
- **sm_name_hdr**
- **sm_cmd_hdr**
- **sm_cmd_opt**

The following table shows the objects used to create each screen type:

Screen Type	Object Class	Object's Use (typical case)
Menu	sm_menu_opt	1 for title of screen
	sm_menu_opt	1 for first item

Screen Type	Object Class	Object's Use (typical case)
	sm_menu_opt	1 for second item

	sm_menu_opt	1 for last item
Selector	sm_name_hdr	1 for title of screen and other attributes
	sm_cmd_opt	1 for entry field or pop-up list
Dialog	sm_cmd_hdr	1 for title of screen and command string
	sm_cmd_opt	1 for first entry field
	sm_cmd_opt	1 for second entry field

	sm_cmd_opt	1 for last entry field

Each object consists of a sequence of named fields and associated values. These are represented in stanza format in ASCII files that can be used by the **odmadd** command to initialize or extend SMIT databases. Stanzas in a file should be separated with one or more blank lines.

Note: Comments in an ODM input file (ASCII stanza file) used by the **odmadd** command must be alone on a line beginning with a # (pound sign) or an * (asterisk) in column one. Only an * (asterisk) comment can be on the same line as a line of the stanza, and must be after the descriptor value.

The following is an example of a stanza for an **sm_menu_opt** object:

```
sm_menu_opt:
  id                = "top_menu"  *name of object class
                    = "050"      *object's (menu screen) name
  id_seq_num        = "050"
  next_id           = "commo"     *id of objects for next menu screen
  text              = "Communications Applications & Services"
  text_msg_file     = ""
  text_msg_set      = 0
  text_msg_id       = 0
  next_type         = "m"         *next_id specified another menu
  alias             = ""
  help_msg_id       = ""
  help_msg_loc      = ""
  help_msg_base     = ""
  help_msg_book     = ""
```

The notation *ObjectClass.Descriptor* is commonly used to describe the value of the fields of an object. For instance, in the preceding **sm_menu_opt** object, the value of **sm_menu_opt.id** is top_menu.

The following is an example of a stanza for an **sm_name_hdr** object:

```
sm_name_hdr:
  id                = ""          *---- used for selector screens
                    = ""          *the name of this selector screen
  next_id           = ""          *next sm_name_hdr or sm_cmd_hdr
  option_id         = ""          *specifies one associated sm_cmd_opt
  has_name_select   = ""
  name              = ""          *title for this screen
  name_msg_file     = ""
  name_msg_id       = 0
  type              = ""
  ghost             = ""
  cmd_to_classify   = ""
  cmd_to_classify_postfix = ""
  raw_field_name    = ""
  cooked_field_name = ""
  next_type         = ""
  help_msg_id       = ""
  help_msg_loc      = ""
  help_msg_base     = ""
  help_msg_book     = ""
```


The following is an example of a stanza for an **sm_cmd_hdr** object:

```
sm_cmd_hdr:
  id          = "" *---- used for dialog screens
  option_id   = "" *the name of this dialog screen
  has_name_select = "" *defines associated set of sm_cmd_opt objects
  name        = "" *title for this screen
  name_msg_file = ""
  name_msg_set = 0
  name_msg_id = 0
  cmd_to_exec  = ""
  ask          = ""
  exec_mode    = ""
  ghost       = ""
  cmd_to_discover = ""
  cmd_to_discover_postfix = ""
  name_size    = 0
  value_size   = 0
  help_msg_id  = ""
  help_msg_loc = ""
  help_msg_base = ""
  help_msg_book = ""
```

The following is an example of a stanza for an **sm_cmd_opt** object:

```
sm_cmd_opt:
  id          = "" *---- used for selector and dialog screens
  id_seq_num  = "" *name of this object
  disc_field_name = "" *"0" if associated with selector screen
  name        = "" *text describing this entry
  name_msg_file = ""
  name_msg_set = 0
  name_msg_id = 0
  op_type     = ""
  entry_type  = ""
  entry_size  = 0
  required    = ""
  prefix      = ""
  cmd_to_list_mode = ""
  cmd_to_list = ""
  cmd_to_list_postfix = ""
  multi_select = ""
  value_index = 0
  disp_values = ""
  values_msg_file = ""
  values_msg_set = 0
  values_msg_id = 0
  aix_values  = ""
  help_msg_id = ""
  help_msg_loc = ""
  help_msg_base = ""
  help_msg_book = ""
```

All SMIT objects have an `id` field that provides a name used for looking up that object. The **sm_menu_opt** objects used for menu titles are also looked up using their `next_id` field. The **sm_menu_opt** and **sm_name_hdr** objects also have `next_id` fields that point to the `id` fields of other objects. These are how the links between screens are represented in the SMIT database. Likewise, there is an `option_id` field in **sm_name_hdr** and **sm_cmd_hdr** objects that points to the `id` fields of their associated **sm_cmd_opt** object(s).

Note: The **sm_cmd_hdr.option_id** object field is equal to each **sm_cmd_opt.id** object field; this defines the link between the **sm_cmd_hdr** object and its associated **sm_cmd_opt** objects.

Two or more dialogs can share common **sm_cmd_opt** objects since SMIT uses the ODM **LIKE** operator to look up objects with the same `sm_cmd_opt.id` field values. SMIT allows up to five IDs (separated by commas) to be specified in a `sm_cmd_hdr.option_id` field, so that **sm_cmd_opt** objects with any of five different `sm_cmd_opt.id` field values can be associated with the **sm_cmd_hdr** object.

The following table shows how the value of an `sm_cmd_hdr.option_id` field relates to the values of the `sm_cmd_opt.id` and `sm_cmd_opt.id_seq_num` fields.

Note: The values in the `sm_cmd_opt.id_seq_num` fields are used to sort the retrieved objects for screen display.

IDs of objects to retrieve (sm_cmd_hdr.option_id)	Objects retrieved (sm_cmd_opt.id)	Display sequence of retrieved objects (sm_cmd_opt.id_seq_num)
"demo.[AB]"	"demo.A"	"10"
	"demo.B"	"20"
	"demo.A"	"30"
	"demo.A"	"40"
"demo.[ACD]"	"demo.A"	"10"
	"demo.C"	"20"
	"demo.A"	"30"
	"demo.A"	"40"
	"demo.D"	"50"
"demo.X,demo.Y,demo.Z"	"demo.Y"	"20"
	"demo.Z"	"40"
	"demo.X"	"60"
	"demo.X"	"80"

The SMIT database

SMIT objects are generated with ODM creation facilities and stored in files in a designated database. The default SMIT database consists of eight files:

- **sm_menu_opt**
- **sm_menu_opt.vc**
- **sm_name_hdr**
- **sm_name_hdr.vc**
- **sm_cmd_hdr**
- **sm_cmd_hdr.vc**
- **sm_cmd_opt**
- **sm_cmd_opt.vc**

The files are stored by default in the **/usr/lib/objrepos** directory. They should always be saved and restored together.

SMIT aliases and fast paths

This section describes SMIT aliases and fast paths

A System Management Interface Tool (SMIT) **sm_menu_opt** object can be used to define a fast path that, when entered with the **smit** command to start SMIT, can get a user directly to a specific menu, selector, or dialog; the alias itself is never displayed. Use of a fast path allows a user to bypass the main SMIT menu and other objects in the SMIT interface path to that menu, selector, or dialog. Any number of fast paths can point to the same menu, selector, or dialog.

An **sm_menu_opt** object is used to define a fast path by setting the `sm_menu_opt.alias` field to "y". In this case, the **sm_menu_opt** object is used exclusively to define a fast path. The new fast path or alias name is specified by the value in the `sm_menu_opt.id` field. The contents of the `sm_menu_opt.next_id` field points to another menu object, selector header object, or dialog header object, depending on whether the value of the `sm_menu_opt.next_type` field is "m" (menu), "n" (selector), or "d" (dialog).

Every non alias **sm_menu_opt** object for a menu title (`next_type="m"`) should have a unique `sm_menu_opt.next_id` field value, since this field is automatically used as a fast path.

If you want two menu items to point to the same successor menu, one of the `next_id` fields should point to an alias, which in turn points to the successor menu.

Build aliases and fast paths by defining them in a stanza file and then processing the file with the **odmadd** command. Several menus, selectors, and dialogs can be defined in a single file. The **odmadd** command adds each alias definition to the specified object repository. The **/usr/lib/objrepos** directory is the default object repository for system information and can be used to store your compiled objects. At SMIT run time, the objects are automatically retrieved from a SMIT database.

Note: You should always back up the **/usr/lib/objrepos** directory before deleting or adding any objects or object classes. Unanticipated damage to objects or classes needed for system operations can cause system problems.

SMIT information command descriptors

The System Management Interface Tool (SMIT) can use several descriptors defined in its objects to get the information, such as current run time values, required to continue through the SMIT interface structure.

Each of these descriptors is assigned some form of command string to run and retrieve the needed data.

The descriptors that can be set to a command for discovery of required information are:

- The **cmd_to_discover** descriptor that is part of the **sm_cmd_hdr** object class used to define a dialog header.
- The **cmd_to_classify** descriptor that is part of the **sm_name_hdr** object class used to define a selector header.
- The **cmd_to_list** descriptor that is part of the **sm_cmd_opt** object class used to define a selector option list associated with a selector or a dialog command option list associated with a dialog entry field.

SMIT executes a command string specified by a **cmd_to_list**, **cmd_to_classify**, or **cmd_to_discover** descriptor by first creating a child process. The standard error (`strerr`) and standard output of the child process are redirected to SMIT via pipes. SMIT next executes a **setenv("ENV=")** subroutine in the child process to prevent commands specified in the **\$HOME/.env** file of the user from being run automatically when a new shell is invoked. Finally, SMIT calls the **execl** system subroutine to start a new **ksh** shell, using the command string as the **ksh -c** parameter value. If the exit value is not 0, SMIT notifies the user that the command failed.

SMIT makes the path names of the log files and the settings of the command line **verbose**, **trace**, and **debug** flags available in the shell environment of the commands it runs. These values are provided via the following environment variables:

- **_SMIT_LOG_FILE**
- **_SMIT_SCRIPT_FILE**
- **_SMIT_VERBOSE_FLAG**
- **_SMIT_TRACE_FLAG**
- **_SMIT_DEBUG_FLAG**

The presence or absence of the corresponding flag is indicated by a value of 0 or 1, respectively.

An easy way to view the current settings is to invoke the shell function after starting SMIT and then run the command string **env | grep _SMIT**.

All writes to the log files should be done as appends and should be immediately followed by flushes unless this occurs automatically.

The **cmd_to_discover** descriptor

When SMIT puts up a dialog, it gets the **sm_cmd_hdr** (dialog header) object and its associated dialog body (one or more **sm_cmd_opt** objects) from the object repository. However, the **sm_cmd_opt** objects can also be initialized with current run time values. If the `sm_cmd_hdr.cmd_to_discover` field is not empty (" "), SMIT runs the command specified in the field to obtain current run time values.

Any valid **ksh** command string can be used as a **cmd_to_discover** descriptor value. The command should generate the following output format as its standard output:

```
#name_1:name_2: ... :name_n\nvalue_1:value_2: ... :value_n
```

In the standard output of a command, the first character is always a # (pound sign). A \n (new line character) is always present to separate the name line from the value line. Multiple names and values are separated by : (colons). And any name or value can be an empty string (which in the output format appears as two colons with no space between them). SMIT maintains an internal current value set in this format that is used to pass name-value pairs from one screen to the next.

Note: If the value includes a : (colon), the : must be preceded by #! (pound sign, exclamation point). Otherwise, SMIT reads the : (colon) as a field separator.

When SMIT runs a command specified in a `cmd_to_discover` field, it captures the `stdout` of the command and loads these name-value pairs (`name_1` and `value_1` `name_2` and `value_2`, and so on) into the **disp_values** and **aix_values** descriptors of the **sm_cmd_opt** (dialog command option) objects by matching each name to a **sm_cmd_opt.disc_field_name** descriptor in each **sm_cmd_opt** object.

For a **sm_cmd_opt** (dialog command option) object that displays a value from a preceding selector, the **disc_field_name** descriptor for the dialog command option object must be set to `"_rawname"` or `"_cookedname"` (or whatever alternate name was used to override the default name) to indicate which value to use. In this case, the **disc_field_name** descriptor of the **sm_cmd_opt** (dialog command option) object should normally be a no-entry field. If a particular value should always be passed to the command, the **required** descriptor for the **sm_cmd_opt** (dialog command option) object must be set to `y` (yes), or one of the other alternatives.

A special case of option ring field initialization permits the current value for a **cmd_to_discover** descriptor (that is, any name-value pair from the current value set of a dialog) of a ring entry field to specify which pre-defined ring value to use as the default or initial value for the corresponding entry field. At dialog initialization time, when a dialog entry field matches a name in the current value set of the dialog (via **sm_cmd_opt.disc_field_name**), a check is made to determine if it is an option ring field (**sm_cmd_opt.op_type** = `"r"`) and if it has predefined ring values (**sm_cmd_opt.aix_values** != `""`). If so, this set of option ring values is compared with the current value for **disc_field_name** from the current value set. If a match is found, the matched option ring value becomes the default ring value (**sm_cmd_opt.value_index** is set to its index). The corresponding translated value (**sm_cmd_opt.disp_values**), if available, is displayed. If no match is found, the error is reported and the current value becomes the default and only value for the ring.

In many cases, discovery commands already exist. In the devices and storage areas, the general paradigms of add, remove, change, and show exist. For example, to add (**mk**), a dialog is needed to solicit characteristics. The dialog can have as its discovery command the show (**ls**) command with a parameter that requests default values. SMIT uses the standard output of the show (**ls**) command to fill in the suggested defaults. However, for objects with default values that are constants known at development time (that is, that are not based on the current state of a given machine), the defaults can be initialized in the dialog records themselves; in this case, no **cmd_to_discover** is needed. The dialog is then displayed. When all fields are filled in and the dialog is committed, the add (**mk**) command is executed.

As another example, a change (**ch**) dialog can have as its discovery command a show (**ls**) command to get current values for a given instance such as a particular device. SMIT uses the standard output of the show (**ls**) command to fill in the values before displaying the dialog. The show (**ls**) command used for discovery in this instance can be the same as the one used for discovery in the add (**mk**) example, except with a slightly different set of options.

The **cmd_to_*_postfix** descriptors

Associated with each occurrence of a **cmd_to_discover**, **cmd_to_classify**, or **cmd_to_list** descriptor is a second descriptor that defines the postfix for the command string defined by the **cmd_to_discover**, **cmd_to_classify**, or **cmd_to_list** descriptor. The postfix is a character string defining the flags and parameters that are appended to the command before it is executed.

The descriptors that can be used to define a postfix to be appended to a command are:

- The **cmd_to_discover_postfix** descriptor that defines the postfix for the **cmd_to_discover** descriptor in an **sm_cmd_hdr** object defining a dialog header.
- The **cmd_to_classify_postfix** descriptor that defines the postfix for the **cmd_to_classify** descriptor in an **sm_name_hdr** object defining a selector header.
- The **cmd_to_list_postfix** descriptor that defines the postfix for the **cmd_to_list** descriptor in an **sm_cmd_opt** object defining a selector entry field associated with a selector or a dialog entry field associated with a dialog.

The following is an example of how the postfix descriptors are used to specify parameter flags and values. The * (asterisk) in the example can be `list`, `classify`, or `discover`.

Assume that `cmd_to_*` equals "DEMO -a", that `cmd_to_*_postfix` equals "-l _rawname -n stuff -R _cookedname", and that the current value set is:

```
#name1: _rawname: _cookedname: :stuff\nvalue1:gigatronicundulator:parallel:xxx:47
```

Then the constructed command string would be:

```
DEMO -a -l 'gigatronicundulator' -n '47' -R 'parallel'
```

Surrounding ' ' (single-quotation marks) can be added around postfix descriptor values to permit handling of parameter values with embedded spaces.

SMIT command generation and execution

Each dialog in the System Management Interface Tool (SMIT) builds and executes a version of a standard command.

The command to be executed by the dialog is defined by the **cmd_to_exec** descriptor in the **sm_cmd_hdr** object that defines the dialog header.

Generating dialog defined tasks

In building the command defined in an **sm_cmd_hdr.cmd_to_exec** descriptor, SMIT uses a two-pass scan over the dialog set of **sm_cmd_opt** objects to collect prefix and parameter values. The parameter values collected include those that the user changed from their initially displayed values and those with the **sm_cmd_opt.required** descriptor set to "y".

The first pass gathers all of the values of the **sm_cmd_opt** objects (in order) for which the **prefix** descriptor is either an empty string ("") or starts with a - (a minus sign). These parameters are not position-sensitive and are added immediately following the command name, together with the contents of the **prefix** descriptor for the parameter.

The second pass gathers all of the values of the remaining **sm_cmd_opt** objects (in order) for which the **prefix** descriptor is -- (two dashes). These parameters are position-sensitive and are added after the flagged options collected in the first pass.

Note: SMIT executes the value of what you enter in the prefix field. If the value in the prefix field is a reserved shell character, for example, the * (asterisk), you must follow the character with a -- (dash dash single quotation mark). Then, when the system evaluates the character, it does not mistake it for a shell character.

Command parameter values in a dialog are filled in automatically when the **disc_field_name** descriptors of its **sm_cmd_opt** objects match names of values generated by preceding selectors or a preceding

discovery command. These parameter values are effectively default values and are normally not added to the command line. Initializing an **sm_cmd_opt.required** descriptor to "y" or "+" causes these values to be added to the command line even when they are not changed in the dialog. If the **sm_cmd_opt.required** descriptor value is "?", the corresponding values are used only if the associated entry field is non-empty. These parameter values are built into the command line as part of the regular two-pass process.

Leading and trailing white spaces (spaces and tabs) are removed from parameter values except when the **sm_cmd_opt.entry_type** descriptor is set to "r". If the resulting parameter value is an empty string, no further action is taken unless the **sm_cmd_opt.prefix** descriptor starts with an option flag. Surrounding single quotation marks are added to the parameter value if the **prefix** descriptor is not set to "--" (two dashes). Each parameter is placed immediately after the associated prefix, if any, with no intervening spaces. Also, if the **multi_select** descriptor is set to "m", tokens separated by white space in the entry field are treated as separate parameters.

Executing dialog defined tasks

SMIT runs the command string specified in a **sm_cmd_hdr.cmd_to_exec** descriptor by first creating a child process. The standard error and standard output of the child process are handled as specified by the contents of the **sm_cmd_hdr.exec_mode** descriptor. SMIT next runs a **setenv("ENV=")** subroutine in the child process to prevent commands specified in the **\$HOME/.env** file of the user from being run automatically when a new shell is invoked. Finally, SMIT calls the **execl** subroutine to start a **ksh** shell, using the command string as the **ksh -c** parameter value.

SMIT makes the path names of the log files and the settings of the command line verbose, trace, and debug flags available in the shell environment of the commands it runs. These values are provided with the following environment variables:

- **_SMIT_LOG_FILE**
- **_SMIT_SCRIPT_FILE**
- **_SMIT_VERBOSE_FLAG**
- **_SMIT_TRACE_FLAG**
- **_SMIT_DEBUG_FLAG**

The presence or absence of the corresponding flag is indicated by a value of 0 or 1, respectively.

Additionally, the **SMIT** environment variable provides information about which SMIT environment is active. The **SMIT** environment variable can have the following values:

Value	SMIT environment
a	SMIT in an ASCII interface
d	SMIT in the Distributed SMIT (DSMIT) interface
m	SMIT in a windows (also called Motif) interface

An easy way to view the current settings is to invoke the shell function after starting SMIT and then run the command string **env | grep SMIT**.

You can disable the function key F9=Shell by setting the environment variable **SMIT_SHELL=n**.

All writes to the log files should be done as appends and should immediately be followed by flushes where this does not occur automatically.

You can override SMIT default output redirection of the (child) task process by setting the **sm_cmd_hdr.exec_mode** field to "i". This setting gives output management control to the task, since the task process simply inherits the standard error and standard output file descriptors.

You can cause SMIT to shutdown and replace itself with the target task by setting the **sm_cmd_hdr.exec_mode** field to "e".

Adding tasks to the SMIT database

When developing new objects for the System Management Interface Tool (SMIT) database, it is recommended that you set up a separate test database for development.

Procedure

To create a test database, do the following:

1. Create a directory for testing use. For example, the following command creates a `/home/smit/test` directory:

```
mkdir /home/smit /home/smit/test
```

2. Make the test directory the current directory:

```
cd /home/smit/test
```

3. Define the test directory as the default object repository by setting the **ODMDIR** environment variable to `.` (the current directory):

```
export ODMDIR= .
```

4. Create a new SMIT database in the test directory:

```
cp /usr/lib/objrepos/sm_* .
```

To add tasks to the SMIT database:

1. Design the dialog for the command you want SMIT to build.
2. Design the hierarchy of menus and, optionally, of selectors needed to get a SMIT user to the dialog, and determine where and how this hierarchy should be linked into the existing SMIT database. The following strategy may save you time if you are developing SMIT database extensions for the first time:
 - a. Start SMIT (run the **smit** command), look for existing menu, selector, and dialog screens that perform tasks similar to the one you want to add, and find the menu screen(s) to which you will add the new task.
 - b. Exit from SMIT, then remove the existing SMIT log file. Instead of removing the log file, you can use the **-l** flag of the **smit** command to specify a different log file when starting SMIT in the following step. This enables you to isolate the trace output of your next SMIT session.
 - c. Start SMIT again with the **-t** command flags and again look at the screen to which you will add the new task. This logs the object IDs accessed for each screen for the next step.
 - d. Look at the SMIT log file to determine the ID for each object class used as part of the menu(s).
 - e. Use the object class IDs with the **odmget** command to retrieve the stanzas for these objects. The stanzas can be used as rough examples to guide your implementation and to learn from the experience of others.
 - f. Look in the SMIT log file for the command strings used when running through the screens to see if special tools are being utilized (such as **sed** or **awk** scripts, **ksh** shell functions, environment variable assignment, and so on). When entering command strings, keep in mind that they are processed twice: the first time by the **odmadd** command and the second time by the **ksh** shell. Be careful when using special escape meta-characters such as `\` or quotation characters (`'` and `"`). Note also that the output of the **odmget** command does not always match the input to the **odmadd** command, especially when these characters or multiline string values are used.
3. Code the dialog, menu, and selector objects by defining them in the ASCII object stanza file format required by the **odmadd** command.
4. Add the dialog, menu, and selector objects to the SMIT test database with the **odmadd** command, using the name of your ASCII object stanza file in place of `test_stanzas`:

```
odmadd test_stanzas
```

5. Test and debug your additions by running SMIT using the local test database:

```
smit -o
```

When you are finished testing, restore the `/etc/objrepos` directory as the default object repository by setting the **ODMDIR** environment variable to `/etc/objrepos`:

```
export ODMDIR=/etc/objrepos
```

Debugging SMIT database extensions

This section describes the process of debugging SMIT database extensions.

Prerequisite tasks or conditions

1. Add a task to the SMIT Database.
2. Test the task.

Procedure

1. Identify the problem using one of the following flags:
 - Run the **smit -v** command if the problem applies to the following SMIT descriptors:
 - **cmd_to_list**
 - **cmd_to_classify**
 - **cmd_to_discover**
 - Run the **smit -t** command if the problems applies to individual SMIT database records.
 - Run the **smit -l** command to generate an alternate log file. Use the alternate log file to isolate current session information.
2. Modify the SMIT database where the incorrect information resides.
3. Retest SMIT task.

Creating SMIT help information for a new task

System Management Interface Tool (SMIT) helps are an extension of the SMIT program.

They are a series of helps designed to give you online information about the components of SMIT used to construct dialogs and menus. SMIT helps reside in a database, just as the SMIT executable code resides in a database. SMIT has two ways to retrieve SMIT help information:

Each of these methods provides a different way to retrieve SMIT helps from the SMIT help database.

Man pages method

Prerequisite tasks or conditions: Create a SMIT task that requires help information.

1. Using any editor, create a file and enter help text inside the file. The file must adhere to the format specified by the **man** command. Put only one set of help information in a file.
2. Give the help text file a title as specified by the **man** command.
3. Place the help text file in the correct place in the **manual** subdirectory.
4. Test the newly created file to ensure it works using the **man** command.
5. Locate the file that contains the ASCII object stanza file format for the new SMIT task.
6. Locate the help descriptor fields in the object stanzas of the file.
7. Set the `help_msg_loc` help descriptor field equal to the title of the help text file. The title for the text file is also the parameter to pass to the **man** command. For example:

```
help_msg_loc = "xx", where "xx" = title string name
```


This example executes the **man** command with the `xx` title string name.

8. Leave the rest of the help descriptor fields empty.

Message catalog method

Prerequisite tasks or conditions: Create a new SMIT task that requires help information.

1. Use any editor to create a file and enter help messages inside the file. The **.msg** file must adhere to the format specified by the message facility.

Note: An existing **.msg** file can also be used.

2. Give each help message a set number (Set #) and a message number (MSG#). This allows the system to retrieve the proper help text.

3. Use the **gencat** command to convert the **.msg** file into a **.cat** file. Place the **.cat** file in the correct directory according to the **NLSPATH** environment variable.

4. Test the help messages using the **dspmsg** command.

5. Locate the file that contains the ASCII object stanza file format for the new SMIT task.

6. Locate the help descriptor fields in the object stanzas of the file.

7. For each object stanza, locate the `help_msg_id` help descriptor field. Enter the Set# and Msg# values for the message in the **.msg** file. These values must adhere to the Messages Facility format. For example, to retrieve message #14 for set #2, set:

```
help_msg_id - "2,14"
```

8. Set the `help_msg_loc` help descriptor field to the filename of the file containing the help text.

9. Leave the other help descriptor fields empty.

sm_menu_opt (SMIT menu) object class

Each item on a menu is specified by an **sm_menu_opt** object.

The displayed menu represents the set of objects that have the same value for **id** plus the **sm_menu_opt** object used for the title, which has a **next_id** value equal to the **id** value of the other objects.

Note: When coding an object in this object class, set unused empty strings to "" (double-quotation marks) and unused integer fields to 0.

The descriptors for **sm_menu_opt** objects are:

Descriptor	Definition
id	The ID or name of the object. The value of id is a string with a maximum length of 64 characters. IDs should be unique both to your application and unique within the particular SMIT database used. See the next_id and alias definitions for this object for related information.
id_seq_num	The position of this item in relation to other items on the menu. Non-title sm_menu_opt objects are sorted on this string field. The value of id_seq_num is a string with a maximum length of 16 characters.
next_id	The fast path name of the next menu, if the value for the next_type descriptor of this object is "m" (menu). The next_id of a menu should be unique both to your application and within the particular SMIT database used. All non-alias sm_menu_opt objects with id values matching the value of next_id form the set of selections for that menu. The value of next_id is a string with a maximum length of 64 characters.
text	The description of the task that is displayed as the menu item. The value of text is a string with a maximum length of 1024 characters. This string can be formatted with embedded \n (newline) characters.

Descriptor	Definition
text_msg_file	The file name (not the full path name) that is the Message Facility catalog for the string, text . The value of text_msg_file is a string with a maximum length of 1024 characters. Message catalogs required by an application program can be developed with the Message Facility. Set to "" if you are not using the Message Facility.
text_msg_set	The Message Facility set ID for the string, text . Set IDs can be used to indicate subsets of a single catalog. The value of text_msg_set is an integer. Set to 0 if you are not using the Message Facility.
text_msg_id	The Message Facility ID for the string, text . The value of text_msg_id is an integer. Set to 0 if you are not using the Message Facility.
next_type	The type of the next object if this item is selected. Valid values are: "m" Menu; the next object is a menu (sm_menu_opt). "d" Dialog; the next object is a dialog (sm_cmd_hdr). "n" Name; the next object is a selector (sm_name_hdr). "i" Info; this object is used to put blank or other separator lines in a menu, or to present a topic that does not lead to an executable task but about which help is provided via the help_msg_loc descriptor of this object.
alias	Defines whether or not the value of the id descriptor for this menu object is an alias for another existing fast path specified in the next_id field of this object. The value of the alias descriptor must be "n" for a menu object.
help_msg_id	Specifies a Message Facility message set number and message ID number with a comma as the separator or a numeric string equal to a SMIT identifier tag.
help_msg_loc	The file name sent as a parameter to the man command for retrieval of help text, or the file name of a file containing help text. The value of help_msg_loc is a string with a maximum length of 1024 characters.
help_msg_base	The fully qualified path name of a library that SMIT reads for the file names associated with the correct book.
help_msg_book	Contains the string with the value of the name file contained in the file library indicated by help_msg_base .

The **sm_menu_opt** Object Class Used for Aliases

A SMIT alias is specified by an **sm_menu_opt** object.

The descriptors for the **sm_menu_opt** object class and their settings to specify an alias are:

Descriptor	Definition
id	The ID or name of the new or alias fast path. The value of id is a string with a maximum length of 64 characters. IDs should be unique to your application and unique within the SMIT database in which they are used.
id_seq_num	Set to "" (empty string).
next_id	Specifies the id_seq_num of the menu object pointed to by the alias. The value of next_id is a string with a maximum length of 64 characters.
text	Set to "" (empty string).

Descriptor	Definition
text_msg_file	Set to "" (empty string).
text_msg_set	Set to 0.
text_msg_id	Set to 0.
next_type	The fast path screen type. The value of next_type is a string. Valid values are: "m" Menu; the next_id field specifies a menu screen fast path. "d" Dialog; the next_id field specifies a dialog screen fast path. "n" Name; the next_id field specifies a selector screen fast path.
alias	Defines this object as an alias fast path. The alias descriptor for an alias must be set to "y" (yes).
help_msg_id	Set to "" (empty string).
help_msg_loc	Set to "" (empty string).
help_msg_base	Set to "" (empty string).
help_msg_book	Set to "" (empty string).

sm_name_hdr (SMIT selector header) object class

A selector screen is specified by two objects: an **sm_name_hdr** object that specifies the screen title and other information, and an **sm_cmd_opt** object that specifies what type of data item is to be obtained.

Note: When coding an object in this object class, set unused empty strings to "" (double-quotation marks) and unused integer fields to 0.

In a SMIT Selector Header screen (**sm_name_hdr**) with type = "c", if you specify a value using a : (colon), (for example, tty:0), SMIT inserts a #! (pound sign, exclamation point) in front of the : to signify that the : is not a field separator. SMIT removes the #! after parsing the rest of the value, before passing it to the **cmd_to_classify** descriptor. To make any further additions to the **cmd_to_classify** descriptor, reinsert the #! in front of the :

The descriptors for the **sm_name_hdr** object class are:

Descriptor	Definition
id	The ID or name of the object. The id field can be externalized as a fast path ID unless has_name_select is set to "y" (yes). The value of id is a string with a maximum length of 64 characters. IDs should be unique to your application and unique within your system.
next_id	Specifies the header object for the subsequent screen; set to the value of the id field of the sm_cmd_hdr object or the sm_name_hdr object that follows this selector. The next_type field described below specifies which object class is indicated. The value of next_id is a string with a maximum length of 64 characters.
option_id	Specifies the body of this selector; set to the id field of the sm_cmd_opt object. The value of option_id is a string with a maximum length of 64 characters.

Descriptor

has_name_select

Definition

Specifies whether this screen must be preceded by a selector screen. Valid values are:

"" or "n"

No; this is the default case. The **id** of this object can be used as a fast path, even if preceded by a selector screen.

"y"

Yes; a selector must precede this object. This setting prevents the **id** of this object from being used as a fast path to the corresponding dialog screen.

name

The text displayed as the title of the selector screen. The value of **name** is a string with a maximum length of 1024 characters. The string can be formatted with embedded \n (newline) characters.

name_msg_file

The file name (not the full path name) that is the Message Facility catalog for the string, **name**. The value of **name_msg_file** is a string with a maximum length of 1024 characters. Message catalogs required by an application program can be developed with the Message Facility.

name_msg_set

The Message Facility set ID for the string, **name**. Set IDs can be used to indicate subsets of a single catalog. The value of **name_msg_set** is an integer.

name_msg_id

The Message Facility ID for the string, **name**. The value of **name_msg_id** is an integer.

type

The method to be used to process the selector. The value of **type** is a string with a maximum length of 1 character. Valid values are:

"" or "j"

Just next ID; the object following this object is always the object specified by the value of the **next_id** descriptor. The **next_id** descriptor is a fully-defined string initialized at development time.

"r"

Cat raw name; in this case, the **next_id** descriptor is defined partially at development time and partially at runtime by user input. The value of the **next_id** descriptor defined at development time is concatenated with the value selected by the user to create the **id** value to search for next (that of the dialog or selector to display).

"c"

Cat cooked name; the value selected by the user requires processing for more information. This value is passed to the command named in the **cmd_to_classify** descriptor, and then output from the command is concatenated with the value of the **next_id** descriptor to create the **id** descriptor to search for next (that of the dialog or selector to display).

Descriptor	Definition
ghost	<p>Specifies whether to display this selector screen or only the list pop-up panel produced by the command in the <code>cmd_to_list</code> field. The value of ghost is a string. Valid values are:</p> <p>"" or "n" No; display this selector screen.</p> <p>"y" Yes; display only the pop-up panel produced by the command string constructed using the <code>cmd_to_list</code> and <code>cmd_to_list_postfix</code> fields in the associated sm_cmd_opt object. If there is no cmd_to_list value, SMIT assumes this object is a super-ghost (nothing is displayed), runs the cmd_to_classify command, and proceeds.</p>
cmd_to_classify	<p>The command string to be used, if needed, to classify the value of the name field of the sm_cmd_opt object associated with this selector. The value of cmd_to_classify is a string with a maximum length of 1024 characters. The input to the cmd_to_classify taken from the <code>entry</code> field is called the "raw name" and the output of the cmd_to_classify is called the "cooked name". Previous to AIX Version 4.2.1, you could create only one value with cmd_to_classify. If that value included a colon, it was escaped automatically. In AIX 4.2.1 and later, you can create multiple values with cmd_to_classify, but the colons are no longer escaped. The colon is now being used as a delimiter by this command. If you use colons in your values, you must preserve them manually.</p>
cmd_to_classify_postfix	<p>The postfix to interpret and add to the command string in the <code>cmd_to_classify</code> field. The value of cmd_to_classify_postfix is a string with a maximum length of 1024 characters.</p>
raw_field_name	<p>The alternate name for the raw value. The value of raw_field_name is a string with a maximum length of 1024 characters. The default value is <code>"_rawname"</code>.</p>
cooked_field_name	<p>The alternate name for the cooked value. The value of cooked_field_name is a string with a maximum length of 1024 characters. The default value is <code>"cookedname"</code>.</p>
next_type	<p>The type of screen that follows this selector. Valid values are:</p> <p>"n" Name; a selector screen follows. See the description of next_id above for related information.</p> <p>"d" Dialog; a dialog screen follows. See the description of next_id above for related information.</p>
help_msg_id	<p>Specifies a Message Facility message set number and message ID number with a comma as the separator or a numeric string equal to a SMIT identifier tag.</p>
help_msg_loc	<p>The file name sent as a parameter to the man command for retrieval of help text, or the file name of a file containing help text. The value of help_msg_loc is a string with a maximum length of 1024 characters.</p>
help_msg_base	<p>The fully qualified path name of a library that SMIT reads for the file names associated with the correct book.</p>

Descriptor**help_msg_book****Definition**

Contains the string with the value of the name file contained in the file library indicated by **help_msg_base**.

sm_cmd_opt (SMIT dialog/selector command option) object class

Each object in a dialog, except the dialog header object, normally corresponds to a flag, option, or attribute of the command that the dialog performs.

One or more of these objects is created for each SMIT dialog; a ghost dialog can have no associated dialog command option objects. Each selector screen is composed of one selector header object and one selector command option object.

Note: When coding an object in this object class, set unused empty strings to "" (double-quotation marks) and unused integer fields to 0.

The dialog command option object and the selector command option object are both **sm_cmd_opt** objects. The descriptors for the **sm_cmd_opt** object class and their functions are:

Descriptor**Function****id**

The ID or name of the object. The **id** of the associated dialog or selector header object can be used as a fast path to this and other dialog objects in the dialog. The value of **id** is a string with a maximum length of 64 characters. All dialog objects that appear in one dialog must have the same ID. Also, IDs should be unique to your application and unique within the particular SMIT database used.

id_seq_num

The position of this item in relation to other items on the dialog; **sm_cmd_opt** objects in a dialog are sorted on this string field. The value of **id_seq_num** is a string with a maximum length of 16 characters. When this object is part of a dialog screen, the string "0" is not a valid value for this field. When this object is part of a selector screen, the **id_seq_num** descriptor must be set to 0.

disc_field_name

A string that should match one of the name fields in the output of the **cmd_to_discover** command in the associated dialog header. The value of **disc_field_name** is a string with a maximum length of 64 characters.

The value of the **disc_field_name** descriptor can be defined using the raw or cooked name from a preceding selector instead of the **cmd_to_discover** command in the associated header object. If the descriptor is defined with input from a preceding selector, it must be set to either "_rawname" or "_cookedname", or to the corresponding **sm_name_hdr.cooked_field_name** value or **sm_name_hdr.raw_field_name** value if this was used to redefine the default name.

name

The string that appears on the dialog or selector screen as the field name. It is the visual questioning or prompting part of the object, a natural language description of a flag, option or parameter of the command specified in the **cmd_to_exec** field of the associated dialog header object. The value of **name** is a string with a maximum length of 1024 characters.

name_msg_file

The file name (not the full path name) that is the Message Facility catalog for the string, **name**. The value of **name_msg_file** is a string with a maximum length of 1024 characters. Message catalogs required by an application program can be developed with the Message Facility. Set to "" (empty string) if not used.

Descriptor	Function
name_msg_set	The Message Facility set ID for the string, name . The value of name_msg_set is an integer. Set to 0 if not used.
name_msg_id	The Message Facility message ID for the string, name . The value of name_msg_id is an integer. Set to 0 if not used.
op_type	<p>The type of auxiliary operation supported for this field. The value of op_type is a string. Valid values are:</p> <p>"" or "n" - This is the default case. No auxiliary operations (list or ring selection) are supported for this field.</p> <p>"l" - List selection operation provided. A pop-up window displays a list of items produced by running the command in the cmd_to_list field of this object when the user selects the F4=List function of the SMIT interface.</p> <p>"r" - Ring selection operation provided. The string in the disp_values or aix_values field is interpreted as a comma-delimited set of valid entries. The user can tab or backtab through these values to make a selection. Also, the F4=List interface function can be used in this case, since SMIT will transform the ring into a list as needed.</p> <p>The values "N", "L", and "R" can be used as op_type values just as the lowercase values "n", "l", and "r". However, with the uppercase values, if the cmd_to_exec command is run and returns with an exit value of 0, then the corresponding entry field will be cleared to an empty string.</p>
entry_type	<p>The type of value required by the entry field. The value of entry_type is a string. Valid values are:</p> <p>"" or "n" - No entry; the current value cannot be modified via direct type-in. The field is informational only.</p> <p>"t" - Text entry; alphanumeric input can be entered.</p> <p>"#" - Numeric entry; only the numeric characters 0, 1, 2, 3, 4, 5, 6, 7, 8, or 9 can be entered. A - (minus sign) or + (plus sign) can be entered as the first character.</p> <p>"x" - Hex entry; hexadecimal input only can be entered.</p> <p>"f" - File entry; a file name should be entered.</p> <p>"r" - Raw text entry; alphanumeric input can be entered. Leading and trailing spaces are considered significant and are not stripped off the field.</p>
entry_size	Limits the number of characters the user can type in the entry field. The value of entry_size is an integer. A value of 0 defaults to the maximum allowed value size.

Descriptor	Function
required	<p>Defines if a command field must be sent to the cmd_to_exec command defined in the associated dialog header object. The value of required is a string. If the object is part of a selector screen, the required field should normally be set to "" (empty string). If the object is part of a dialog screen, valid values are:</p> <p>"" or "n" - No; the option is added to the command string in the cmd_to_exec command only if the user changes the initially-displayed value. This is the default case.</p> <p>"y" - Yes; the value of the prefix field and the value of the entry field are always sent to the cmd_to_exec command.</p> <p>"+" - The value of the prefix field and the value of the entry field are always sent to the cmd_to_exec command. The entry field must contain at least one non-blank character. SMIT will not allow the user to run the task until this condition is satisfied.</p> <p>"?" - Except when empty; the value of the prefix field and the value of the entry field are sent to the cmd_to_exec field unless the entry field is empty.</p>
prefix	<p>In the simplest case, defines the flag to send with the entry field value to the cmd_to_exec command defined in the associated dialog header object. The value of prefix is a string with a maximum length of 1024 characters.</p> <p>The use of this field depends on the setting of the required field, the contents of the prefix field, and the contents of the associated entry field.</p> <p>Note: If the prefix field is set to – (dash dash), the content of the associated entry field is appended to the end of the cmd_to_exec command. If the prefix field is set to –' (dash dash single quotation mark), the contents of the associated entry field is appended to the end of the cmd_to_exec command in single quotes.</p>
cmd_to_list_mode	<p>Defines how much of an item from a list should be used. The list is produced by the command specified in this object's cmd_to_list field. The value of cmd_to_list_mode is a string with a maximum length of 1 character. Valid values are:</p> <p>"" or "a" - Get all fields. This is the default case.</p> <p>"1" - Get the first field.</p> <p>"2" - Get the second field.</p> <p>"r" - Range; running the command string in the cmd_to_list field returns a range (such as 1..99) instead of a list. Ranges are for information only; they are displayed in a list pop-up, but do not change the associated entry field.</p>
cmd_to_list	<p>The command string used to get a list of valid values for the value field. The value of cmd_to_list is a string with a maximum length of 1024 characters. This command should output values that are separated by \n (new line) characters.</p>

Descriptor	Function
cmd_to_list_postfix	The postfix to interpret and add to the command string specified in the <code>cmd_to_list</code> field of the dialog object. The value of cmd_to_list_postfix is a string with a maximum length of 1024 characters. If the first line starts with <code>#</code> (pound sign) following a space, that entry will be made non-selectable. This is useful for column headings. Subsequent lines that start with a <code>#</code> , optionally preceded by spaces, are treated as a comment and as a continuation of the preceding entry.
multi_select	<p>Defines if the user can make multiple selections from a list of valid values produced by the command in the <code>cmd_to_list</code> field of the dialog object. The value of multi_select is a string. Valid values are:</p> <p>"" - No; a user can select only one value from a list. This is the default case.</p> <p>"," - Yes; a user can select multiple items from the list. When the command is built, a comma is inserted between each item.</p> <p>"y" - Yes; a user can select multiple values from the list. When the command is built, the option prefix is inserted once before the string of selected items.</p> <p>"m" - Yes; a user can select multiple items from the list. When the command is built, the option prefix is inserted before each selected item.</p>
value_index	For an option ring, the zero-origin index to the array of <code>disp_value</code> fields. The value_index number indicates the value that is displayed as the default in the entry field to the user. The value of entry_size is an integer.
disp_values	The array of valid values in an option ring to be presented to the user. The value of the <code>disp_values</code> fields is a string with a maximum length of 1024 characters. The field values are separated by <code>,</code> (commas) with no spaces preceding or following the commas.
values_msg_file	The file name (not the full path name) that is the Message Facility catalog for the values in the <code>disp_values</code> fields, if the values are initialized at development time. The value of the <code>values_msg_file</code> field is a string with a maximum length of 1024 characters. Message catalogs required by an application program can be developed with the Message Facility.
values_msg_set	The Message Facility set ID for the values in the <code>disp_values</code> fields. Set to 0 if not used.
values_msg_id	The Message Facility message ID for the values in the <code>disp_values</code> fields. Set to 0 if not used.
aix_values	If for an option ring, an array of values specified so that each element corresponds to the element in the disp_values array in the same position; use if the natural language values in disp_values are not the actual options to be used for the command. The value of the <code>aix_values</code> field is a string with a maximum length of 1024 characters.
help_msg_id	Specifies a Message Facility message set number and message ID number with a comma as the separator or a numeric string equal to a SMIT identifier tag.

Descriptor	Function
help_msg_loc	The file name sent as a parameter to the man command for retrieval of help text, or the file name of a file containing help text. The value of help_msg_loc is a string with a maximum length of 1024 characters.
help_msg_base	The fully qualified path name of a library that SMIT reads for the file names associated with the correct book.
help_msg_book	Contains the string with the value of the name file contained in the file library indicated by help_msg_base .

sm_cmd_hdr (SMIT dialog header) object class

A dialog header object is an **sm_cmd_hdr** object. A dialog header object is required for each dialog, and points to the dialog command option objects associated with the dialog.

Note: When coding an object in this object class, set unused empty strings to "" (double-quotation marks) and unused integer fields to 0.

The descriptors for the **sm_cmd_hdr** object class are:

Descriptor	Function
id	The ID or name of the object. The value of id is a string with a maximum length of 64 characters. The id field can be used as a fast path ID unless there is a selector associated with the dialog. IDs should be unique to your application and unique within your system.
option_id	The id of the sm_cmd_opt objects (the dialog fields) to which this header refers. The value of option_id is a string with a maximum length of 64 characters.
has_name_select	Specifies whether this screen must be preceded by a selector screen or a menu screen. Valid values are: "" or "n" No; this is the default case. "y" Yes; a selector precedes this object. This setting prevents the id of this object from being used as a fast path to the corresponding dialog screen.
name	The text displayed as the title of the dialog screen. The value of name is a string with a maximum length of 1024 characters. The text describes the task performed by the dialog. The string can be formatted with embedded \n (newline) characters.
name_msg_file	The file name (not the full path name) that is the Message Facility catalog for the string, name . The value of name_msg_file is a string with a maximum length of 1024 characters. Message catalogs required by an application program can be developed with the Message Facility.
name_msg_set	The Message Facility set ID for the string, name . Set IDs can be used to indicate subsets of a single catalog. The value of name_msg_set is an integer.
name_msg_id	The Message Facility ID for the string, name . Message IDs can be created by the message extractor tools owned by the Message Facility. The value of name_msg_id is an integer.

Descriptor

cmd_to_exec

Function

The initial part of the command string, which can be the command or the command and any fixed options that execute the task of the dialog. Other options are automatically appended through user interaction with the command option objects (**sm_cmd_opt**) associated with the dialog screen. The value of **cmd_to_exec** is a string with a maximum length of 1024 characters.

ask

Defines whether or not the user is prompted to reconsider the choice to execute the task. Valid values are:

"" or "n"

No; the user is not prompted for confirmation; the task is performed when the dialog is committed. This is the default setting for the **ask** descriptor.

"y"

Yes; the user is prompted to confirm that the task be performed; the task is performed only after user confirmation.

Prompting the user for execution confirmation is especially useful prior to performance of deletion tasks, where the deleted resource is either difficult or impossible to recover, or when there is no displayable dialog associated with the task (when the **ghost** field is set to "y").

Descriptor

exec_mode

Function

Defines the handling of standard input, standard output, and the **stderr** file during task execution. The value of **exec_mode** is a string. Valid values are:

"" or "p"

Pipe mode; the default setting for the **exec_mode** descriptor. The command executes with standard output and the **stderr** file redirected through pipes to SMIT. SMIT manages output from the command. The output is saved and is scrollable by the user after the task finishes running. While the task runs, output is scrolled as needed.

"n"

No scroll pipe mode; works like the "p" mode, except that the output is not scrolled while the task runs. The first screen of output will be shown as it is generated and then remains there while the task runs. The output is saved and is scrollable by the user after the task finishes running.

"i"

Inherit mode; the command executes with standard input, standard output, and the **stderr** file inherited by the child process in which the task runs. This mode gives input and output control to the executed command. This value is intended for commands that need to write to the **/dev/tty** file, perform cursor addressing, or use **libcur** or **libcurses** library operations.

"e"

Exit/exec mode; causes SMIT to run (do an **exec1** subroutine call on) the specified command string in the current process, which effectively terminates SMIT. This is intended for running commands that are incompatible with SMIT (which change display modes or font sizes, for instance). A warning is given that SMIT will exit before running the command.

"E"

Same as "e", but no warning is given before exiting SMIT.

"P" , "N" or "I"

Backup modes; work like the corresponding "p", "n", and "i" modes, except that if the **cmd_to_exec** command is run and returns with an exit value of 0, SMIT backs up to the nearest preceding menu (if any), or to the nearest preceding selector (if any), or to the command line.

ghost

Indicates if the normally displayed dialog should not be shown. The value of **ghost** is a string. Valid values are:

"" or "n"

No; the dialog associated with the task is displayed. This is the default setting.

"y"

Yes; the dialog associated with the task is not displayed because no further information is required from the user. The command specified in the **cmd_to_exec** descriptor is executed as soon as the user selects the task.

Descriptor	Function
cmd_to_discover	The command string used to discover the default or current values of the object being manipulated. The value of cmd_to_discover is a string with a maximum length of 1024 characters. The command is executed before the dialog is displayed, and its output is retrieved. Output of the command must be in colon format.
cmd_to_discover_postfix	The postfix to interpret and add to the command string in the <code>cmd_to_discover</code> field. The value of cmd_to_discover_postfix is a string with a maximum length of 1024 characters.
help_msg_id	Specifies a Message Facility message set number and message ID number with a comma as the separator or a numeric string equal to a SMIT identifier tag.
help_msg_loc	The file name sent as a parameter to the man command for retrieval of help text, or the file name of a file containing help text. The value of help_msg_loc is a string with a maximum length of 1024 characters.
help_msg_base	The fully qualified path name of a library that SMIT reads for the file names associated with the correct book.
help_msg_book	Contains the string with the value of the name file contained in the file library indicated by help_msg_base .

SMIT example program

The following example program is designed to help you write your own stanzas.

If you add these stanzas to the SMIT directory that comes with the operating system, they will be accessible through SMIT by selecting the **Applications** item in the SMIT main menu. All of the demos are functional except for Demo 3, which does not install any languages.

```

#-----
# Intro:
# Unless you are creating a new SMIT database, first you need
# to decide where to insert the menu for your application.
# Your new menu will point to other menus, name headers, and
# dialogs. For this example, we are inserting a pointer to the
# demo menu under the "Applications" menu option. The next_id for
# the Applications menu item is "apps", so we begin by creating a
# menu_opt with "apps" as its id.
#-----
sm_menu_opt:
    id           = "apps"
    id_seq_num   = "010"
    next_id      = "demo"
    text         = "SMIT Demos"
    next_type    = "m"

sm_menu_opt:
    id           = "demo"
    id_seq_num   = "010"
    next_id      = "demo_queue"
    text         = "Demo 1: Add a Print Queue"
    next_type    = "n"

sm_menu_opt:
    id           = "demo"
    id_seq_num   = "020"
    next_id      = "demo_mle_inst_lang_hdr"
    text         = "Demo 2: Add Language for Application Already Installed"
    next_type    = "n"

#---
# Since demo_mle_inst_lang_hdr is a descriptive, but not very

```

```

# memorable name, an alias with a simpler name can be made to
# point to the same place.
#----
sm_menu_opt:
    id          = "demo_lang"
    next_id     = "demo_mle_inst_lang_hdr"
    next_type   = "n"
    alias      = "y"

sm_menu_opt:
    id_seq_num  = "030"
    id          = "demo"
    next_id     = "demo_lspv"
    text        = "Demo 3: List Contents of a Physical Volume"
    text_msg_file = "smit.cat"
    next_type   = "n"

sm_menu_opt:
    id_seq_num  = "040"
    id          = "demo"
    next_id     = "demo_date"
    text        = "Demo 4: Change / Show Date, Time"
    text_msg_file = "smit.cat"
    next_type   = "n"

#-----
# Demo 1
# -----
# Goal: Add a Print Queue. If the printers.rte package is not
#       installed, install it automatically. If the user is
#       running MSMIT (SMIT in a windows interface), launch a
#       graphical program for this task. Otherwise, branch to
#       the SMIT print queue task.
#
# Topics:      1. cooked output & cmd_to_classify
#              2. SMIT environment variable (msmit vs. ascii)
#              3. ghost name_hdr
#              4. super-ghost name_hdr
#              5. creating an "OK / cancel" option
#              6. dspmsg for translations
#              7. exit/exec mode
#              8. id_seq_num for a name_hdr option
#-----
#----
# Topics: 1,4
# Note that the next_id is the same as the id. Remember that the
# output of the cmd_to_classify is appended to the next_id,
# since the type is "c", for cooked. So, the next_id will be
# either demo_queue1 or demo_queue2. None of the output of the
# name_hdr is displayed, and there is no cmd_to_list in the
# demo_queue_dummy_opt, making this name_hdr a super-ghost.
#----
sm_name_hdr:
    id          = "demo_queue"
    next_id     = "demo_queue"
    option_id   = "demo_queue_dummy_opt"
    name        = "Add a Print Queue"
    name_msg_file = "smit.cat"
    name_msg_set = 52
    name_msg_id  = 41
    type        = "c"
    ghost       = "y"
    cmd_to_classify = "\
x()
{
    # Check to see if the printer file is installed.
    lslpp -l printers.rte 2>/dev/null 1>/dev/null
    if [[ $? != 0 ]]
    then
    echo 2
    else
    echo 1
    fi
}
x"
    next_type   = "n"

#----
# Topics: 2,4
# Having determined the printer software is installed, we want
# to know if the gui program should be run or if we should
# branch to the ascii SMIT screen for this task. To do this, we

```

```

# check the value of the environment variable SMIT, which is "m"
# for windows (Motif) or "a" for ascii. Here again we tack the
# output of the cmd_to_classify onto the next_id.
#----
sm_name_hdr:
  id           = "demo_queue1"
  next_id      = "mkpq"
  option_id    = "demo_queue_dummy_opt"
  has_name_select = ""
  ghost       = "y"
  next_type   = "n"
  type       = "c"
  cmd_to_classify = "\
gui_check()
{
  if [ $SMIT = \"m\" ]; then
    echo gui
  fi
}
gui_check"

sm_name_hdr:
  id       = "mkpqgui"
  next_id  = "invoke_gui"
  next_type = "d"
  option_id = "demo_queue_dummy_opt"
  ghost    = "y"

#----
# Topics: 7
# Note: the exec_mode of this command is "e", which
# exits SMIT before running the cmd_to_exec.
#----
sm_cmd_hdr:
  id           = "invoke_gui"
  cmd_to_exec  = "/usr/bin/X11/xpntm"
  exec_mode    = "e"
  ghost       = "y"

sm_cmd_opt:
  id           = "demo_queue_dummy_opt"
  id_seq_num  = 0

#----
# Topics: 3,5
# The printer software is not installed. Install the software
# and loop back to demo_queue1 to check the SMIT environment
# variable. This is a ghost name_hdr. The cmd_to_list of the
# sm_cmd_opt is displayed immediately as a pop-up option
# instead of waiting for the user to input a response. In this
# ghost, the cmd_opt is a simple OK/cancel box that prompts the
# user to press return.
#----
sm_name_hdr:
  id           = "demo_queue2"
  next_id      = "demo_queue1"
  option_id    = "demo_queue_opt"
  name         = "Add a Print Queue"
  name_msg_file = "smit.cat"
  name_msg_set = 52
  name_msg_id  = 41
  ghost       = "y"
  cmd_to_classify = "\
install_printers ()
{
  # Install the printer package.
  /usr/lib/assist/install_pkg \"printers.rte\" 2>&1 >/dev/null
  if [[ $? != 0 ]]
  then
    echo "Error installing printers.rte"
    exit 1
  else
    exit 0
  fi
}
install_printers "
  next_type = "n"

#----
# Topics: 5,6,8
# Here a cmd_opt is used as an OK/cancel box. Note also that the

```

```

# command dspmsg is used to display the text for the option. This
# allows for translation of the messages.
# Note: the id_seq_num for the option is 0. Only one option is
# allowed per name_hdr, and its id_seq_num must be 0.
#----
sm_cmd_opt:
    id                = "demo_queue_opt"
    id_seq_num        = "0"
    disc_field_name   = ""
    name              = "Add a Print Queue"
    name_msg_file     = "smit.cat"
    name_msg_set      = 52
    name_msg_id       = 41
    op_type           = "1"
    cmd_to_list       = "x()\
{
if [ $SMIT = \"a\" ] \n\
then \n\
dspmsg -s 52 smit.cat 56 \
'Press Enter to automatically install the printer software.\n\
Press F3 to cancel.\n\
'\n\
else \n\
dspmsg -s 52 smit.cat 57 'Click on this item to automatically install
the printer software.\n' \n\
fi\n\
}\n\
x"
    entry_type        = "t"
    multi_select      = "n"

#-----
#
# Demo 2
# -----
# Goal: Add a Language for an Application Already Installed. It
# is often clearer to the user to get some information
# before displaying the dialog screen. Name Headers
# (sm_name_hdr) can be used for this purpose. In this
# example, two name headers are used to determine the
# language to install and the installation device. The
# dialog has entries for the rest of the information needed
# to perform the task.
#
# Topics:
# 1. Saving output from successive name_hdrs with
#    cooked_field_name
# 2. Using getopt inside cmd_to_exec to process cmd_opt
#    info
# 3. Ring list vs. cmd_to_list for displaying values
#    cmd_opts
#-----

#----
# Topic: 1
# This is the first name_hdr. It is called by the menu_opt for
# this function. We want to save the user's input for later use
# in the dialog. The parameter passed into the cmd_to_classify
# comes from the user's selection/entry. Cmd_to_classify cleans
# up the output and stores it in the variable specified by
# cooked_field_name. This overrides the default value for the
# cmd_to_classify output, which is _cookedname. The default must
# be overridden because we also need to save the output of the
# next name_hdr.
#----
sm_name_hdr:
    id                = "demo_mle_inst_lang_hdr"
    next_id           = "demo_mle_inst_lang"
    option_id         = "demo_mle_inst_lang_select"
    name              = "Add Language for Application Already Installed"
    name_msg_file     = "smit.cat"
    name_msg_set      = 53
    name_msg_id       = 35
    type              = "j"
    ghost             = "n"
    cmd_to_classify   = "\
foo() {
    echo $1 | sed -n \"s/[^]*\\[[^]*\\].*/\\1/p\"
}
foo"
    cooked_field_name = "add_lang_language"
    next_type         = "n"

```



```

    help_msg_id           = "2850325"

sm_cmd_opt:
  id                     = "demo_mle_inst_lang_select"
  id_seq_num             = "0"
  disc_field_name       = "add_lang_language"
  name                   = "LANGUAGE translation to install"
  name_msg_file         = "smit.cat"
  name_msg_set          = 53
  name_msg_id           = 20
  op_type               = "1"
  entry_type            = "n"
  entry_size            = 0
  required              = ""
  prefix                = "-l "
  cmd_to_list_mode      = "a"
  cmd_to_list           = "/usr/lib/nls/lsmle -l"
  help_msg_id           = "2850328"

#----
# Topic:1
# This is the second name_hdr. Here the user's input is passed
# directly through the cmd_to_classify and stored in the
# variable add_lang_input.
#----
sm_name_hdr:
  id                     = "demo_mle_inst_lang"
  next_id               = "demo_dialog_add_lang"
  option_id             = "demo_add_input_select"
  has_name_select       = "y"
  name                  = "Add Language for Application Already Installed"
  name_msg_file         = "smit.cat"
  name_msg_set          = 53
  name_msg_id           = 35
  type                  = "j"
  ghost                 = "n"
  cmd_to_classify       = "\
    foo() {
      echo $1
    }
    foo"
  cooked_field_name     = "add_lang_input"
  next_type             = "d"
  help_msg_id           = "2850328"

sm_cmd_opt:
  id                     = "demo_add_input_select"
  id_seq_num            = "0"
  disc_field_name       = "add_lang_input"
  name                  = "INPUT device/directory for software"
  name_msg_file         = "smit.cat"
  name_msg_set          = 53
  name_msg_id           = 11
  op_type               = "1"
  entry_type            = "t"
  entry_size            = 0
  required              = "y"
  prefix                = "-d "
  cmd_to_list_mode      = "1"
  cmd_to_list           = "/usr/lib/instl/sm_inst list_devices"
  help_msg_id           = "2850313"

#----
# Topic: 2
# Each of the cmd_opts formats its information for processing
# by the getopt command (a dash and a single character, followed
# by an optional parameter). The colon following the letter in
# the getopt command means that a parameter is expected after
# the dash option. This is a nice way to process the cmd_opt
# information if there are several options, especially if one of
# the options could be left out, causing the sequence of $1, $2,
# etc. to get out of order.
#----
sm_cmd_hdr:
  id                     = "demo_dialog_add_lang"
  option_id             = "demo_mle_add_app_lang"
  has_name_select       = ""
  name                  = "Add Language for Application Already Installed"
  name_msg_file         = "smit.cat"
  name_msg_set          = 53
  name_msg_id           = 35
  cmd_to_exec           = "\

```

```

foo()
{
while getopts d:l:S:X Option \ "$@"
do
    case $Option in
        d) device=$OPTARG;;
        l) language=$OPTARG;;
        S) software=$OPTARG;;
        X) extend_fs="-X";;
    esac
done

if [[ ` /usr/lib/assist/check_cd -d $device ` = '1' ]]
then
    /usr/lib/assist/mount_cd $device
    CD_MOUNTED=true
fi

if [[ $software = \ "ALL\ " ]]
then
    echo "Installing all software for $language..."
else
    echo "Installing $software for $language..."
fi
exit $RC
}
foo"

ask                = "y"
ghost              = "n"
help_msg_id        = "2850325"

sm_cmd_opt:
id                 = "demo_mle_add_app_lang"
id_seq_num         = "0"
disc_field_name    = "add_lang_language"
name               = "LANGUAGE translation to install"
name_msg_file      = "smit.cat"
name_msg_set       = 53
name_msg_id        = 20
entry_type         = "n"
entry_size         = 0
required           = "y"
prefix             = "-l "
cmd_to_list_mode   = "a"
help_msg_id        = "2850328"

#---
# Topic: 2
# The prefix field precedes the value selected by the user, and
# both the prefix and the user-selected value are passed into
# the cmd_to_exec for getopts processing.
#---
sm_cmd_opt:
id                 = "demo_mle_add_app_lang"
id_seq_num         = "020"
disc_field_name    = "add_lang_input"
name               = "INPUT device/directory for software"
name_msg_file      = "smit.cat"
name_msg_set       = 53
name_msg_id        = 11
entry_type         = "n"
entry_size         = 0
required           = "y"
prefix             = "-d "
cmd_to_list_mode   = "1"
cmd_to_list        = "/usr/lib/instl/sm_inst list_devices"
help_msg_id        = "2850313"

sm_cmd_opt:
id                 = "demo_mle_add_app_lang"
id_seq_num         = "030"
name               = "Installed APPLICATION"
name_msg_file      = "smit.cat"
name_msg_set       = 53
name_msg_id        = 43
op_type            = "l"
entry_type         = "n"
entry_size         = 0
required           = "y"
prefix             = "-S "
cmd_to_list_mode   = ""
cmd_to_list        = "\

```

```

list_messages (
{
    language=$1
    device=$2
    lslpp -Lqc | cut -f2,3 -d':'
}
list_messages"
cmd_to_list_postfix = "add_lang_language add_lang_input"
multi_select       = ", "
value_index        = 0
disp_values        = "ALL"
help_msg_id        = "2850329"

#----
# Topic: 3
# Here, instead of a cmd_to_list, there is a comma-delimited set
# of Ring values in the disp_values field. This list is displayed
# one item at a time as the user presses tab in the cmd_opt entry
# field. However, instead of passing a yes or no to the cmd_hdr,
# it is more useful to use the aix_values field to pass either
# a -X or nothing. The list in the aix_values field must match
# one-to-one with the list in the disp_values field.
#----
sm_cmd_opt:
    id_seq_num = "40"
    id = "demo_mle_add_app_lang"
    disc_field_name = ""
    name = "EXTEND file systems if space needed?"
    name_msg_file = "smit.cat"
    name_msg_set = 53
    name_msg_id = 12
    op_type = "r"
    entry_type = "n"
    entry_size = 0
    required = "y"
    multi_select = "n"
    value_index = 0
    disp_values = "yes,no"
        values_msg_file = "sm_inst.cat"
    values_msg_set = 1
    values_msg_id = 51
    aix_values = "-X,"
    help_msg_id = "0503005"

#-----
#
# Demo 3
# -----
# Goal: Show Characteristics of a Logical Volume. The name of the
# logical volume is entered by the user and passed to the
# cmd_hdr as _rawname.
#
# Topics:          1. _rawname
#                  2. Ringlist & aix_values
#-----

#----
# Topic: 1
# No rawname is needed because we have only one name_hdr and
# we can use the default variable name _rawname.
#----
sm_name_hdr:
    id = "demo_lspv"
    next_id = "demo_lspvd"
    option_id = "demo_cmdlvmpvns"
    has_name_select = ""
    name = "List Contents of a Physical Volume"
    name_msg_file = "smit.cat"
    name_msg_set = 15
    name_msg_id = 100
    type = "j"
    ghost = ""
    cmd_to_classify = ""
    raw_field_name = ""
    cooked_field_name = ""
    next_type = "d"
    help_msg_id = "0516100"

sm_cmd_opt:
    id_seq_num = "0"
    id = "demo_cmdlvmpvns"
    disc_field_name = "PVName"

```

```

name = "PHYSICAL VOLUME name"
name_msg_file = "smit.cat"
name_msg_set = 15
name_msg_id = 101
op_type = "l"
entry_type = "t"
entry_size = 0
required = "+"
cmd_to_list_mode = "1"
cmd_to_list = "lsvg -o|lsvg -i -p|grep -v '[:P]|\ \
    cut -f1 -d' '"
cmd_to_list_postfix = ""
multi_select = "n"
help_msg_id = "0516021"

#---
# Topic: 1
# The cmd_to_discover_postfix passes in the name of the physical
# volume, which is the raw data selected by the user in the
# name_hdr - _rawname.
#---
sm_cmd_hdr:
id = "demo_lspvd"
option_id = "demo_cmdlvmlspv"
has_name_select = "y"
name = "List Contents of a Physical Volume"
name_msg_file = "smit.cat"
name_msg_set = 15
name_msg_id = 100
cmd_to_exec = "lspv"
ask = "n"
cmd_to_discover_postfix = "_rawname"
help_msg_id = "0516100"

sm_cmd_opt:
id_seq_num = "01"
id = "demo_cmdlvmlspv"
disc_field_name = "_rawname"
name = "PHYSICAL VOLUME name"
name_msg_file = "smit.cat"
name_msg_set = 15
name_msg_id = 101
op_type = "l"
entry_type = "t"
entry_size = 0
required = "+"
cmd_to_list_mode = "1"
cmd_to_list = "lsvg -o|lsvg -i -p|grep -v '[:P]|\ \
    cut -f1 -d' '"
help_msg_id = "0516021"

#---
# Topic: 2
# Here a ringlist of 3 values matches with the aix_values we
# want to pass to the sm_cmd_hdr's cmd_to_exec.
#---
sm_cmd_opt:
id_seq_num = "02"
id = "demo_cmdlvmlspv"
disc_field_name = "Option"
name = "List OPTION"
name_msg_file = "smit.cat"
name_msg_set = 15
name_msg_id = 92
op_type = "r"
entry_type = "n"
entry_size = 0
required = "n"
value_index = 0
disp_values = "status,logical volumes,physical \
    partitions"
values_msg_file = "smit.cat"
values_msg_set = 15
values_msg_id = 103
aix_values = " , -l, -p"
help_msg_id = "0516102"

#-----
#
# Demo 4
# -----
# Goal: Change / Show Date & Time

```

```

#
# Topics:      1. Using a ghost name header to get variable
#              values for the next dialog screen.
#              2. Using a cmd_to_discover to fill more than one
#              cmd_opt with initial values.
#              3. Re-ordering parameters in a cmd_to_exec.
#-----

#----
# Topic: 1
# This ghost name_hdr gets two values and stores them in the
# variables daylight_y_n and time_zone for use in the cmd_opts
# for the next dialog. The output of cmd_to_classify is colon-
# delimited, as is the list of field names in cooked_field_name.
#----
sm_name_hdr:
  id = "demo_date"
  next_id = "demo_date_dial"
  option_id = "date_sel_opt"
  name_msg_set = 0
  name_msg_id = 0
  ghost = "y"
  cmd_to_classify = " \
if [ $(echo $TZ | awk '{ \
if (length($1) <=6) {printf(\"2\")} \
else {printf(\"1\")} }') = 1 ] \n\
then\n\
  echo $(dspmsg smit.cat -s 30 18 'yes')\" :$TZ\" \n\
else\n\
  echo $(dspmsg smit.cat -s 30 19 'no')\" :$TZ\" \n\
fi #"
  cooked_field_name = "daylight_y_n:time_zone"

sm_cmd_opt:
  id_seq_num = "0"
  id = "date_sel_opt"

#----
# Topic: 2,3
# Here the cmd_to_discover gets six values, one for each of the
# editable sm_cmd_opts for this screen. The cmd_to_discover
# output is two lines, the first with a # followed by a list of
# variable names, and the second line the list of values. Both
# lists are colon-delimited. We also see here the cmd_to_exec
# taking the parameters from the cmd_opts and reordering them
# when calling the command.
#----
sm_cmd_hdr:
  id = "demo_date_dial"
  option_id = "demo_chtz_opts"
  has_name_select = "y"
  name = "Change / Show Date & Time"
  name_msg_file = "smit.cat"
  name_msg_set = 30
  name_msg_id = 21
  cmd_to_exec = "date_proc () \
# MM dd hh mm ss yy\n\
# dialogue param order # 3 4 5 6 7 2\n\
{\n\
date \"$3$4$5$6.$7$2\" \n\
}\n\
date_proc "
  exec_mode = "P"
  cmd_to_discover = "disc_proc() \n\
{\n\
TZ=\"$1\" \n\
echo '#cur_month:cur_day:cur_hour:cur_min:cur_sec:cur_year' \n\
date +%m:%d:%H:%M:%S:%y \n\
}\n\
disc_proc"
  cmd_to_discover_postfix = ""
  help_msg_id = "055101"

#----
# The first two cmd_opts get their initial values
# (disc_field_name) from the name_hdr.
#----
sm_cmd_opt:
  id_seq_num = "04"
  id = "demo_chtz_opts"
  disc_field_name = "time_zone"
  name = "Time zone"

```

```

name_msg_file = "smit.cat"
name_msg_set = 30
name_msg_id = 16
required = "y"

sm_cmd_opt:
  id_seq_num = "08"
  id = "demo_chtz_opts"
  disc_field_name = "daylight_y_n"
  name = "Does this time zone go on Daylight Saving Time?\n"
  name_msg_file = "smit.cat"
  name_msg_set = 30
  name_msg_id = 17
  entry_size = 0

#---
# The last six cmd_opts get their values from the
# cmd_to_discover.
#---
sm_cmd_opt:
  id_seq_num = "10"
  id = "demo_chtz_opts"
  disc_field_name = "cur_year"
  name = "YEAR (00-99)"
  name_msg_file = "smit.cat"
  name_msg_set = 30
  name_msg_id = 10
  entry_type = "#"
  entry_size = 2
  required = "+"
  help_msg_id = "055102"

sm_cmd_opt:
  id_seq_num = "20"
  id = "demo_chtz_opts"
  disc_field_name = "cur_month"
  name = "MONTH (01-12)"
  name_msg_file = "smit.cat"
  name_msg_set = 30
  name_msg_id = 11
  entry_type = "#"
  entry_size = 2
  required = "+"
  help_msg_id = "055132"

sm_cmd_opt:
  id_seq_num = "30"
  id = "demo_chtz_opts"
  disc_field_name = "cur_day"
  name = "DAY (01-31)\n"
  name_msg_file = "smit.cat"
  name_msg_set = 30
  name_msg_id = 12
  entry_type = "#"
  entry_size = 2
  required = "+"
  help_msg_id = "055133"

sm_cmd_opt:
  id_seq_num = "40"
  id = "demo_chtz_opts"
  disc_field_name = "cur_hour"
  name = "HOUR (00-23)"
  name_msg_file = "smit.cat"
  name_msg_set = 30
  name_msg_id = 13
  entry_type = "#"
  entry_size = 2
  required = "+"
  help_msg_id = "055134"

sm_cmd_opt:
  id_seq_num = "50"
  id = "demo_chtz_opts"
  disc_field_name = "cur_min"
  name = "MINUTES (00-59)"
  name_msg_file = "smit.cat"
  name_msg_set = 30
  name_msg_id = 14
  entry_type = "#"
  entry_size = 2
  required = "+"

```

```
help_msg_id = "055135"

sm_cmd_opt:
  id_seq_num = "60"
  id = "demo_chtz_opts"
  disc_field_name = "cur_sec"
  name = "SECONDS (00-59)"
  name_msg_file = "smit.cat"
  name_msg_set = 30
  name_msg_id = 15
  entry_type = "#"
  entry_size = 2
  required = "+"
  help_msg_id = "055136"
```

System resource controller

This article provides information about the System Resource Controller (SRC), which facilitates the management and control of complex subsystems.

The SRC is a subsystem controller. Subsystem programmers who own one or more daemon processes can use SRC services to define a consistent system management interface for their applications. The SRC provides a single set of commands to start, stop, trace, refresh, and query the status of a subsystem.

In addition, the SRC provides an error notification facility. You can use this facility to incorporate subsystem-specific recovery methods. The type of recovery information included is limited only by the requirement that the notify method is a string in a file and is executable.

Refer to the following information to learn more about SRC programming requirements:

Subsystem interaction with the SRC

The SRC defines a subsystem as a program or set of related programs designed as a unit to perform related functions. See "[System Resource Controller](#)" in *Operating system and device management* for a more detailed description of the characteristics of a subsystem.

A subserver is a process that belongs to and is controlled by a subsystem.

The SRC operates on objects in the SRC object class. Subsystems are defined to the SRC as subsystem objects; subservers, as subserver-type objects. The structures associated with each type of object are predefined in the [usr/include/sys/srcobj.h](#) file.

The SRC can issue SRC commands against objects at the subsystem, subserver, and subsystem-group levels. A subsystem group is a group of any user-specified subsystems. Grouping subsystems allows multiple subsystems to be controlled by invoking a single command. Groups of subsystems may also share a common notification method.

The SRC communicates with subsystems by sending signals and exchanging request and reply packets. In addition to signals, the SRC recognizes the sockets and IPC message-queue communication types. A number of subroutines are available as an SRC API to assist in programming communication between subsystems and the SRC. The SRC API also supports programming communication between client programs and the SRC.

The SRC and the `init` command

The SRC is operationally independent of the `init` command. However, the SRC is intended to extend the process-spawning functionality provided by this command. In addition to providing a single point of control to start, stop, trace, refresh, and query the status of subsystems, the SRC can control the operations of individual subsystems, support remote system control, and log subsystem failures.

Operationally, the only time the `init` command and the SRC interact occurs when the `srcmstr` (SRC master) daemon is embedded within the `inittab` file. By default, the `srcmstr` daemon is in the `inittab` file. In this case, the `init` command starts the `srcmstr` daemon at system startup, as with all other processes. You must have root user authority or be in the system group to invoke the `srcmstr` daemon.

Compiling programs to interact with the `srcmstr` daemon

To enable programs to interact with the **srcmstr** daemon, the **/usr/include/spc.h** file should be included and the program should be compiled with the **libsrc.a** library. This support is not needed if the subsystem uses signals to communicate with the SRC.

SRC operations

To make use of SRC functionality, a subsystem must interact with the **srcmstr** daemon in two ways:

- A subsystem object must be created for the subsystem in the SRC subsystem object class.
- If a subsystem uses signals, it does not need to use SRC subroutines. However, if it uses message queues or sockets, it must respond to stop requests using the SRC subroutines.

All SRC subsystems must support the **stopsrc** command. The SRC uses this command to stop subsystems and their subserver with the **SIGNORM** (stop normal), **SIGFORCE** (stop force), or **SIGCANCEL** (cancel systems) signals.

Subsystem support is optional for the **startsrc**, **lssrc -l**, **traceson**, **tracesoff**, and **refresh** commands, long status and subserver status reporting, and the SRC notification mechanism.

SRC capabilities

The SRC provides the following support for the subsystem programmer:

- A common command interface to support starting, stopping, and sending requests to a subsystem
- A central point of control for subsystems and groups of subsystems
- A common format for requests to the subsystem
- A definition of subserver so that each subserver can be managed as it is uniquely defined to the subsystem
- The ability to define subsystem-specific error notification methods
- The ability to define subsystem-specific responses to requests for status, trace support, and configuration refresh
- A single point of control for servicing subsystem requests in a network computing environment

Related information

[mknotify](#)

[mkserver](#)

[mksys](#)

[refresh](#)

[rmnotify](#)

[srcmstr](#)

[tracesoff](#)

[traceson](#)

[spc.h](#)

[srcobj.h](#)

SRC objects

The System Resource Controller (SRC) defines and manages three object classes:

Together, these object classes represent the domain in which the SRC performs its functions. A predefined set of object-class descriptors comprise the possible set of subsystem configurations supported by the SRC.

Note: Only the SRC Subsystem object class is required. Use of the Subserver Type and Notify object classes is subsystem-dependent.

Subsystem object class

The subsystem object class contains the descriptors for all SRC subsystems. A subsystem must be configured in this class before it can be recognized by the SRC.

The descriptors for the Subsystem object class are defined in the **SRCsubsys** structure of the **/usr/include/sys/srcobj.h** file. The Subsystem Object Descriptors and Default Values table provides a short-form illustration of the subsystem descriptors as well as the **mkssys** and **chssys** command flags associated with each descriptor.

<i>Table 96. Subsystem Object Descriptors and Default Values</i>		
Descriptors	Default values	Flags
Subsystem name		-s
Path to subsystem command		-p
Command arguments		-a
Execution priority	20	-E
Multiple instance	NO	-Q -q
User ID		-u
Synonym name (key)		-t
Start action	ONCE	-O -R
stdin	/dev/console	-i
stdout	/dev/console	-o
stderr	/dev/console	-e
Communication type	Sockets	-K -I -S
Subsystem message type		-m
Communication IPC queue key		-l
Group name		-G
SIGNORM signal		-n
SIGFORCE signal		-f
Display	Yes	-D -d
Wait time	20 seconds	-w
Auditid		

The subsystem object descriptors are defined as follows:

Object descriptors	Definition
Subsystem name	Specifies the name of the subsystem object. The name cannot exceed 30 bytes, including the null terminator (29 characters for single-byte character sets, or 14 characters for multibyte character sets). This descriptor must be POSIX-compliant. This field is required.
Subsystem command path	Specifies the full path name for the program executed by the subsystem start command. The path name cannot exceed 200 bytes, including the null terminator (199 characters for single-byte character sets, or 99 characters for multibyte character sets). The path name must be POSIX-compliant. This field is required.

Object descriptors	Definition
Command arguments	Specifies any arguments that must be passed to the command that starts the subsystem. The arguments cannot exceed 200 bytes, including the null terminator (199 characters for single-byte character sets, or 99 characters for multibyte character sets). The arguments are parsed by the srcmstr daemon according to the same rules used by shells. For example, quoted strings are passed as a single argument, and blanks outside quoted strings delimit arguments.
Execution priority	Specifies the process priority of the subsystem to be run. Subsystems started by the srcmstr daemon run with this priority. The default value is 20.
Multiple instance	Specifies the number of instances of a subsystem that can run at one time. A value of NO (the -Q flag) specifies that only one instance of the subsystem can run at one time. Attempts to start this subsystem if it is already running will fail, as will attempts to start a subsystem on the same IPC message queue key. A value of YES (the -q flag) specifies that multiple subsystems may use the same IPC message queue and that there can be multiple instances of the same subsystem. The default value is NO.
User ID	Specifies the user ID (numeric) under which the subsystem is run. A value of 0 indicates the root user. This field is required.
Synonym name	Specifies a character string to be used as an alternate name for the subsystem. The character string cannot exceed 30 bytes, including the null terminator (29 characters for single-byte character sets, or 14 characters for multibyte character sets). This field is optional.
Start action	Specifies whether the srcmstr daemon should restart the subsystem after an abnormal end. A value of RESPAWN (the -R flag) specifies the srcmstr daemon should restart the subsystem. A value of ONCE (the -O flag) specifies the srcmstr daemon should not attempt to restart the failed system. There is a respawn limit of two restarts within a specified wait time. If the failed subsystem cannot be successfully restarted, the notification method option is consulted. The default value is ONCE.
Standard Input File/Device	Specifies the file or device from which the subsystem receives its input. The default is /dev/console . This field cannot exceed 200 bytes, including the null terminator (199 characters for single-byte character sets, or 99 characters for multibyte character sets). This field is ignored if the communication type is sockets.
Standard Output File/Device	Specifies the file or device to which the subsystem sends its output. This field cannot exceed 200 bytes, including the null terminator (199 characters for single-byte character sets, or 99 characters for multibyte character sets). The default is /dev/console .
Standard Error File/Device	Specifies the file or device to which the subsystem writes its error messages. This field cannot exceed 200 bytes, including the null terminator (199 characters for single-byte character sets, or 99 characters for multibyte character sets). Failures are handled as part of the notify method. The default is /dev/console .
	Note: Catastrophic errors are sent to the error log.
Communication type	Specifies the communication method between the srcmstr daemon and the subsystem. Three types can be defined: IPC (-I), sockets (-K), or signals (-S). The default is sockets.

Object descriptors	Definition
Communication IPC queue key	Specifies a decimal value that corresponds to the IPC message queue key that the srcmstr daemon uses to communicate to the subsystem. This field is required for subsystems that communicate using IPC message queues. Use the ftok subroutine with a fully qualified path name and an ID parameter to ensure that this key is unique. The srcmstr daemon creates the message queue prior to starting the subsystem.
Group name	Designates the subsystem as a member of a group. This field cannot exceed 30 bytes, including the null terminator (29 characters for single-byte character sets, or 14 characters for multibyte character sets). This field is optional.
Subsystem message type	Specifies the mtype of the message that is placed on the subsystem's message queue. The subsystem uses this value to retrieve messages by using the msgrcv or msgxrcv subroutine. This field is required if you are using message queues.
SIGNORM signal value	Specifies the value to be sent to the subsystem when a stop normal request is sent. This field is required of subsystems using the signals communication type.
SIGFORCE signal value	Specifies the value to be sent to the subsystem when a stop force request is sent. This field is required of subsystems using the signals communication type.
Display value	Indicates whether the status of an inoperative subsystem can be displayed on lssrc -a or lssrc -g output. The -d flag indicates display; the -D flag indicates do not display. The default is -d (display).
Wait time	Specifies the time in seconds that a subsystem has to complete a restart or stop request before alternate action is taken. The default is 20 seconds.
Auditid	Specifies the subsystem audit ID. Created automatically by the srcmstr daemon when a subsystem is defined, this field is used by the security system, if configured. This field cannot be set or changed by a program.

Subserver type object class

An object must be configured in this class if a subsystem has subservers and the subsystem expects to receive subserver-related commands from the **srcmstr** daemon.

This object class contains three descriptors, which are defined in the **SRCsubsvr** structure of the **srcobj.h** file:

Descriptor	Definition
Subserver ID (key)	Specifies the name of the subserver type object identifier. The set of subserver type names defines the allowable values for the -t flag of the subserver commands. The name length cannot exceed 30 bytes, including the terminating null (29 characters for single-byte character sets, or 14 characters for multibyte character sets).
Owning subsystem name	Specifies the name of the subsystem that owns the subserver object. This field is defined as a link to the SRC subsystem object class.

Descriptor	Definition
Code point	Specifies a decimal number that identifies the subserver. The code point is passed to the subsystem controlling the subserver in the <code>object</code> field of the subreq structure of the SRC request structure. If a subserver object name is also provided in the command, the srcmstr daemon forwards the code point to the subsystem in the <code>objname</code> field of the subreq structure. See the " SRC Request Structure Example " in the spc.h file documentation for examples of these elements.

The commands that reference subservers identify each subserver as a named type of subserver and can also append a name to each instance of a subserver type. The SRC daemon uses the subserver type to determine the controlling subsystem for the subserver, but does not examine the subserver name.

Notify object class

This class provides a mechanism for the **srcmstr** daemon to invoke subsystem-provided routines when the failure of a subsystem is detected. When the SRC daemon receives a **SIGCHLD** signal indicating the termination of a subsystem process, it checks the status of the subsystem (maintained by the **srcmstr** daemon) to determine if the termination was caused by a **stopsrc** command. If no **stopsrc** command was issued, the termination is interpreted as an abnormal termination. If the restart action in the definition does not specify `respawn`, or if `respawn` attempts fail, the **srcmstr** daemon attempts to read an object associated with the subsystem name from the Notify object class. If such an object is found, the method associated with the subsystem is run.

If no subsystem object is found in the Notify object class, the **srcmstr** daemon determines whether the subsystem belongs to a group. If so, the **srcmstr** daemon attempts to read an object of that group name from the Notify object class. If such an object is found, the method associated with it is invoked. In this way, groups of subsystems can share a common method.

Note: The subsystem notify method takes precedence over the group notify method. Therefore, a subsystem can belong to a group that is started together, but still have a specific recovery or cleanup routine defined.

Notify objects are defined by two descriptors:

Descriptor	Definition
Subsystem name or group name	Specifies the name of the subsystem or group for which a notify method is defined.
Notify method	Specifies the full path name to the routine that is executed when the srcmstr daemon detects abnormal termination of the subsystem or group.

Such notification is useful when specific recovery or clean-up work needs to be performed before a subsystem can be restarted. It is also a tool for information gathering to determine why a subsystem abnormally stopped.

Notify objects are created with the **mknotify** command. To modify a notify method, the existing notify object must be removed using the **rmnotify** command, and then a new notify object created.

Command	Definition
mknotify	Adds a notify method to the SRC configuration database
rmnotify	Removes a notify method from the SRC configuration database

The **srcmstr** daemon logs subsystem recovery activity. The subsystem is responsible for reporting subsystem failures.

SRC communication types

The System Resource Controller (SRC) supports three communication types: signals, sockets, and interprocess communication (IPC) message queues.

The communication type chosen determines to what degree the subsystem takes advantage of SRC functions.

Note: All subsystems, regardless of the communication type specified in the subsystem environment object, must be capable of supporting limited signals communication. A signal-catcher routine must be defined to handle **SIGTERM** (stop cancel) signals. The **SIGTERM** signal indicates a subsystem should clean up all resources and terminate.

Refer to the following sections to learn more about SRC communication types:

The Communications Between the srcmstr Daemon and Subsystems table summarizes communication type actions associated with SRC functions.

Function	Using IPC or sockets	Using signals
start		
subsystem	SRC forks and execs to create subsystem process.	SRC forks and execs to create subsystem process.
subserver	Uses IPC message queue or socket to send request to subsystem.	Not supported
stop normal		
subsystem	Uses IPC message queue or socket to send request to subsystem.	Sends SIGNORM to the subsystem.
subserver	Uses IPC message queue or socket to send request to subsystem.	Not supported.
stop forced		
subsystem	Uses IPC message queue or socket to send request to subsystem.	Sends SIGFORCE to the subsystem.
subserver	Uses IPC message queue or socket to send request to subsystem.	Not supported.
stop cancel		
subsystem	Sends SIGTERM followed by SIGKILL to the process group of the subsystem.	Sends SIGTERM followed by SIGKILL to the process group of the subsystem.
status short		
subsystem	Implemented by SRC (no subsystem request).	Implemented by SRC (no subsystem request).
subserver	Uses IPC message queue or socket to send request to subsystem.	Not supported.
status long		

Function	Using IPC or sockets	Using signals
subsystem	Uses IPC message queue or socket to send request to subsystem.	Not supported.
subserver	Uses IPC message queue or socket to send request to subsystem.	Not supported.
traceon/traceoff		
subsystem	Uses IPC message queue or socket to send request to subsystem.	Not supported.
subserver	Uses IPC message queue or socket to send request to subsystem.	Not supported.
refresh		
subsystem	Uses IPC message queue or socket to send request to subsystem.	Not supported.
subserver	Uses IPC message queue or socket to send request to subsystem.	Not supported.
notify		
subsystem	Implemented by subsystem-provided method.	Implemented by subsystem-provided method.

Signals communication

The most basic type of communication between a subsystem and the **srcmstr** daemon is accomplished with signals. Because signals constitute a one-way communication scheme, the only SRC command that signals subsystems recognize is a stop request. Subsystems using signals do not recognize long status, refresh, or trace requests. Nor do they recognize subservers.

Signals subsystems must implement a signal-catcher routine, such as the **sigaction**, **sigvec**, or **signal** subroutine, to handle **SIGNORM** and **SIGFORCE** requests.

Signals subsystems are specified in the SRC subsystem object class by issuing a **mkssys -Snf** command string or by using the **defssys** and **addssys** subroutines.

Item	Descriptor
addssys	Adds a subsystem definition to the SRC configuration database
defssys	Initializes a new subsystem definition with default values
mkssys	Adds a subsystem definition to the SRC configuration database

Sockets communication

Increasingly, the communication option of choice for subsystem programmers is sockets. Sockets are also the default communication type for the **srcmstr** daemon. See the [Sockets Overview Communications Programming Concepts](#) for more information.

The **srcmstr** daemon uses sockets to receive work requests from a command process. When this communication type is chosen, the **srcmstr** daemon creates the subsystem socket in which the subsystem will receive **srcmstr** daemon requests. UNIX sockets (**AF_UNIX**) are created for local

subsystems. Internet sockets (**AF_INET**) are created for remote subsystems. The following steps describe the command processing sequence:

1. The command process accepts a command from the input device, constructs a work-request message, and sends the work-request UDP datagram to the **srcmstr** daemon on the well-known SRC port. The **AF_INET** is identified in the **/etc/services** file.
2. The **srcmstr** daemon listens on the well-known SRC port for work requests. Upon receiving a work request, it tells the system to fill the **socket** subroutine's **sockaddr** structure to obtain the originating system's address and appends the address and port number to the work request.
3. The **srcmstr** daemon uses the **srcrrqs** and **srcsrpy** subroutines. It processes only those requests that it can process and then sends the information back to the command process. Other requests are forwarded to the appropriate subsystem on the port that the subsystem has specified for its work requests.
4. The subsystem listens on the port previously obtained by the **srcmstr** daemon for the subsystem. (Each subsystem inherits a port when the **srcmstr** daemon starts a subsystem.) The subsystem processes the work request and sends a reply back to the command process.
5. The command process listens for the response on the specified port.

The file access permissions and addresses of the sockets used by the **srcmstr** daemon are maintained in the **/dev/SRC** and **/dev/.SRC-unix** temporary directories. Though displayable using the **ls** command, the information contained in these directories is for internal SRC use only.

Message queues and sockets offer equal subsystem functionality.

Item	Descriptor
srcrrqs	Saves the destination address of your subsystem's response to a received packet. (Also see threadsafe version srcrrqs_r)
srcsrpy	Sends your subsystem response packet to a request that your subsystem received.

IPC message queue communication

IPC message queue functionality is similar to sockets functionality. Both communication types support a full-function SRC environment.

When the communication type is IPC message queue, the **srcmstr** daemon uses sockets to receive work requests from a command process, then uses an IPC message queue in which the subsystem receives SRC messages. The message queue is created when the subsystem is started, and is used thereafter. Message queue subsystems use the following command-processing sequence to communicate with the **srcmstr** daemon:

1. The **srcmstr** daemon gets the message queue ID from the SRC subsystem object and sends the message to the subsystem.
2. The subsystem waits for the message queue and issues a **msgrcv** subroutine to receive the command from the message queue in the form of the **subreq** structure required of subsystem requests.
3. The subsystem calls the **srcrrqs** subroutine to get a tag ID that is used in responding to the message.
4. The subsystem interprets and processes the received command. Depending upon the command, the subsystem creates either a **svreply** or **statcode** data structure to return a reply to the command process. Refer to the **/usr/include/spc.h** file for more information on these structures.
5. The subsystem calls the **srcsrpy** subroutine to send back a reply buffer to the command process.

Programming subsystem communication with the SRC

System Resource Controller (SRC) commands are executable programs that take options from the command line.

After the command syntax has been verified, the commands call SRC run-time subroutines to construct a User Datagram Protocol (UDP) datagram and send it to the **srcmstr** daemon.

The following sections provide more information about SRC subroutines and how they can be used by subsystems to communicate with the SRC main process:

Programming subsystems to receive SRC requests

The programming tasks associated with receiving SRC requests vary with the communication type specified for the subsystem. The **srcmstr** daemon uses sockets to receive work requests from a command process and constructs the necessary socket or message queue to forward work requests. Each subsystem needs to verify the creation of its socket or message queue. Read the following sections for information on communication type-specific guidelines on programming your subsystem to receive SRC request packets.

Note: All subsystems, regardless of communication type, must define a signal-catcher routine to handle the **SIGTERM** request.

Receiving SRC signals

Subsystems that use signals as their communication type must define a signal-catcher routine to catch the **SIGNORM** and **SIGFORCE** signals. The signal-catching method used is subsystem-dependent. Following are two examples of the types of subroutines that can be used for this purpose.

Subroutine	Description
<u>sigaction</u> , <u>sigvec</u> , or <u>signal</u> subroutine	Specifies the action to take upon the delivery of a signal.
<u>sigset</u> , <u>sighold</u> , <u>sigelse</u> , or <u>sigignore</u> subroutine	Enhances the signal facility and provides signal management for application processes.

Receiving SRC request packets using sockets

Use the following guidelines when programming sockets subsystems to receive SRC request packets:

- Include the SRC subsystem structure in your subsystem code by specifying the /usr/include/spc.h file. This file contains the structures the subsystem uses to respond to SRC commands. In addition, the **spc.h** file includes the **srcerrno.h** file, which does not need to be included separately. The **srcerrno.h** file contains error-code definitions for daemon support.
- When a sockets subsystem is started, the socket on which the subsystem receives SRC request packets is set as file descriptor 0. The subsystem should verify this by calling the getsockname subroutine, which returns the address of the subsystem's socket. If file descriptor 0 is not a socket, the subsystem should log an error and then exit. See "Reading Internet Datagrams Example Program" in *Communications Programming Concepts* for information on how the getsockname subroutine can be used to return the address of a subsystem socket.
- If a subsystem polls more than one socket, use the **select** subroutine to determine which socket has something to read. See "Checking for Pending Connections Example Program" in *Communications Programming Concepts* for more information on how the **select** subroutine can be used for this purpose.
- Use the recvfrom subroutine to get the request packet from the socket.

Note: The return address for the subsystem response packet is in the received SRC request packet. This address should not be confused with the address that the **recvfrom** subroutine returns as one of its parameters.

After the **recvfrom** subroutine completes and the packet has been received, use the srcrrqs subroutine to return a pointer to a static **srchdr** structure. This pointer contains the return address for the subsystem's reply. This structure is overwritten each time the **srcrrqs** subroutine is called, so its contents should be stored elsewhere if they will be needed after the next call to the **srcrrqs** subroutine.

Receiving SRC request packets using message queues

Use the following guidelines when programming message queue subsystems to receive SRC request packets:

- Include the SRC subsystem structure in your subsystem code by specifying the **/usr/include/spc.h** file. This file contains the structures the subsystem uses to respond to SRC commands. In addition, the **spc.h** file includes the **srcerrno.h** include file, which does not need to be included separately. The **srcerrno.h** file contains error-code definitions for daemon support.
- Specify **-DSRCBYQUEUE** as a compile option. This places a message type (mtype) field as the first field in the **srcreq** structure. This structure should be used any time an SRC packet is received.
- When the subsystem has been started, use the **msgget** subroutine to verify that a message queue was created at system startup. The subsystem should log an error and exit if a message queue was not created.
- If a subsystem polls more than one message queue, use the **select** subroutine to determine which message queue has something to read. See "Checking for Pending Connections Example Program" in *Communications Programming Concepts* for information on how the **select** subroutine can be used for this purpose.
- Use the **msgrcv** or **msgxrcv** subroutine to get the packet from the message queue. The return address for the subsystem response packet is in the received packet.
- When the **msgrcv** or **msgxrcv** subroutine completes and the packet has been received, call the **srcrrqs** subroutine to finish the reception process. The **srcrrqs** subroutine returns a pointer to a static **srchdr** structure that is overwritten each time the **srcrrqs** subroutine is called. This pointer contains the return address for the subsystem's reply.

Programming subsystems to process SRC request packets

Subsystems must be capable of processing stop requests. Optionally, subsystems may support start, status, trace, and refresh requests.

Processing request packets involves a two-step process:

Reading SRC request packets

SRC request packets are received by subsystems in the form of a **srcreq** structure as defined in the **/usr/include/spc.h** file. The subsystem request resides in the **subreq** structure of the **srcreq** structure:

```
struct subreq
  short object;          /*object to act on*/
  short action;         /*action START, STOP, STATUS, TRACE,\
                        REFRESH*/
  short parm1;          /*reserved for variables*/
  short parm2;          /*reserved for variables*/
  char objname;         /*object name*/
```

The **object** field of the **subreq** structure indicates the object to which the request applies. When the request applies to a subsystem, the **object** field is set to the **SUBSYSTEM** constant. Otherwise, the **object** field is set to the subserver code point or the **objname** field is set to the subserver PID as a character string. It is the subsystem's responsibility to determine the object to which the request applies.

The **action** field specifies the action requested of the subsystem. Subsystems should understand the **START**, **STOP**, and **STATUS** action codes. The **TRACE** and **REFRESH** action codes are optional.

The **parm1** and **parm2** fields are used differently by each of the actions.

Action	parm1	parm2
STOP	NORMAL or FORCE	
STATUS	LONGSTAT or SHORTSTAT	
TRACE	LONGTRACE or SHORT-TRACE	TRACEON or TRACEOFF

The **START** subserver and **REFRESH** actions do not use the **parm1** and **parm2** fields.

Programming subsystem response to SRC requests

The appropriate subsystem actions for the majority of SRC requests are programmed when the subsystem object is defined to the SRC. The structures that subsystems use to respond to SRC requests are defined in the `/usr/include/spc.h` file. Subsystems may use the following SRC run-time subroutines to meet command processing requirements:

Subroutine Description

- `srcrrqs`** Allows a subsystem to store the header from a request.
- `srcsrpy`** Allows a subsystem to send a reply to a request.

Status-request processing requires a combination of tasks and subroutines.

When subsystems receive requests they cannot process or that are invalid, they must send an error packet with an error code of **SRC_SUBICMD** in response to the unknown, or invalid, request. SRC reserves action codes 0-255 for SRC internal use. If your subsystem receives a request containing an action code that is not valid, your subsystem must return an error code of **SRC_SUBICMD**. Valid action codes supported by SRC are defined in the `spc.h` file. You can also define subsystem-specific action codes. An action code is not valid if it is not defined by the SRC or your subsystem.

Note: Action codes 0-255 are reserved for SRC use.

Processing SRC status requests

Subsystems may be requested to provide three types of status reports: long subsystem status, short subserver status, and long subserver status.

Note: Short subsystem status reporting is performed by the `srcmstr` daemon. Statcode and reply-status value constants for this type of report are defined in the `/usr/include/spc.h` file. The Status Value Constants table lists required and suggested reply-status value codes.

Reply status value codes

Value	Meaning	Subsystem	Subserver
SRCWARN	Received a request to stop. (Will be stopped within 20 seconds.)	X	X
SRCACT	Started and active.	X	X
SRCINAC	Not active.		
SRCINOP	Inoperative.	X	X
SRCLOSD	Closed.		
SRCCLSPN	In the process of being closed.		
SRCNOSTAT	Idle.		
SRCOBIN	Open, but not active.		
SRCOPND	Open.		
SRCOPPN	In the process of being opened.		
SRCSTAR	Starting.		X
SRCSTPG	Stopping.	X	X
SRCTST	TEST active.		
SRCTSTPN	TEST pending.		

The SRC **lssrc** command displays the received information on standard output. The information returned by subsystems in response to a long status request is left to the discretion of the subsystem. Subsystems that own subservers are responsible for tracking and reporting the state changes of subservers, if desired. Use the **srcstathdr** subroutine to retrieve a standard status header to pass back at the beginning of your status data.

The following steps are recommended in processing status requests:

1. To return status from a subsystem (short or long), allocate an array of **statcode** structures plus a **srchdr** structure. The **srchdr** structure must start the buffer that you are sending in response to the status request. The **statcode** structure is defined in the **/usr/include/spc.h** file.

```
struct statcode
{
    short objtype;
    short status;
    char objtext [65];
    char objname [30];
};
```

2. Fill in the **objtype** field with the **SUBSYSTEM** constant to indicate that the status is for a subsystem, or with a subserver code point to indicate that the status is for a subserver.
3. Fill in the **status** field with one of the SRC status constants defined in the **spc.h** file.
4. Fill in the **objtext** field with the NLS text that you wish displayed as status. This field must be a NULL terminated string.
5. Fill in the **objname** field with the name of the subsystem or subserver for which the **objtext** field applies. This field must be a NULL terminated string.

Note: The subsystem and requester can agree to send other subsystem-defined information back to the requester.

Programming subsystems to send reply packets

The packet that a subsystem returns to the SRC should be in the form of the **srcrep** structure as defined in the **/usr/include/spc.h** file. The **svrreply** structure that is part of the **srcrep** structure will contain the subsystem reply:

```
struct svrreply
{
    short rtncode;           /*return code from the subsystem*/
    short objtype;          /*SUBSYSTEM or SUBSERVER*/
    char objtext[65];        /*object description*/
    char objname[20];        /*object name*/
    char rtnmsg[256];        /*returned message*/
};
```

Use the **srcsrpy** subroutine to return a packet to the requester.

Creating a reply

To program a subsystem reply, use the following procedure:

1. Fill in the **rtncode** field with the SRC error code that applies. Use **SRC_SUBMSG** as the **rtncode** field to return a subsystem-specific NLS message.
2. Fill in the **objtype** field with the **SUBSYSTEM** constant to indicate that the reply is for a subsystem, or with the subserver code point to indicate that the reply is for a subserver.
3. Fill in the **objname** field with the subsystem name, subserver type, or subserver object that applies to the reply.
4. Fill in the **rtnmsg** field with the subsystem-specific NLS message.
5. Key the appropriate entry in the **srcsrpy** *Continued* parameter. See "srcsrpy Continuation Packets" for more information.

Note: The last packet from the subsystem must always have **END** specified in the *Continued* parameter to the **srcsrpy** subroutine.

srcsrpy continuation packets

Subsystem responses to SRC requests are made in the form of continuation packets. Two types of continuation packets may be specified: Informative message, and reply packets.

The informative message is not passed back to the client. Instead, it is printed to the client's standard output. The message must consist of NLS text, with message tokens filled in by the sending subsystem. To send this type of continuation packet, specify CONTINUED in the **srcsrpy** subroutine *Continued* parameter.

Note: The STOP subsystem action does not allow any kind of continuation. However, all other action requests received by the subsystem from the SRC may be sent an informative message.

The reply packet is passed back to the client for further processing. Therefore, the packet must be agreed upon by the subsystem and the requester. One example of this type of continuation is a status request. When responding to subsystem status requests, specify STATCONTINUED in the **srcsrpy** *Continued* parameter. When status reporting has completed, or all subsystem-defined reply packets have been sent, specify END in the **srcsrpy** *Continued* parameter. The packet is then passed to the client to indicate the end of the reply.

Programming subsystems to return SRC error packets

Subsystems are required to return error packets for both SRC errors and non-SRC errors.

When returning an SRC error, the reply packet that the subsystem returns should be in the form of the **svrreply** structure of the **srcrep** structure, with the objname field filled in with the subsystem name, subserver type, or subserver object in error. If the NLS message associated with the SRC error number does not include any tokens, the error packet is returned in short form. This means the error packet contains the SRC error number only. However, if tokens are associated with the error number, standard NLS message text from the message catalog should be returned.

When returning a non-SRC error, the reply packet should be a **svrreply** structure with the rtncode field set to the SRC_SUBMSG constant and the rtnmsg field set to a subsystem-specific NLS message. The rtnmsg field is printed to the client's standard output.

Responding to trace requests

Support for the **traceson** and **tracesoff** commands is subsystem-dependent. If you choose to support these commands, trace actions can be specified for subsystems and subservers.

Subsystem trace requests will arrive in the following form: A subsystem trace request will have the **subreq** action field set to the TRACE constant and the **subreq** object field set to the SUBSYSTEM constant. The trace action uses parm1 to indicate LONGTRACE or SHORTTRACE trace, and parm2 to indicate TRACEON or TRACEOFF.

When the subsystem receives a trace subsystem packet with parm1 set to SHORTTRACE and parm2 set to TRACEON, the subsystem should turn short tracing on. Conversely, when the subsystem receives a trace subsystem packet with parm1 set to LONGTRACE and parm2 set to TRACEON, the subsystem should turn long tracing on. When the subsystem receives a trace subsystem packet with parm2 set to TRACEOFF, the subsystem should turn subsystem tracing off.

Subserver trace requests will arrive in the following form: the subserver trace request will have the **subreq** action field set to the TRACE constant and the **subreq** object field set to the subserver code point of the subserver to send status on. The trace action uses parm1 to indicate LONGTRACE or SHORTTRACE, and parm2 to indicate TRACEON or TRACEOFF.

When the subsystem receives a trace subserver packet with parm1 set to SHORTTRACE and parm2 set to TRACEON, the subsystem should turn subserver short tracing on. Conversely, when the subsystem receives a trace subserver packet with parm1 set to LONGTRACE and parm2 set to TRACEON, the subsystem should turn subserver long tracing on. When the subsystem receives a trace subserver packet with parm2 set to TRACEOFF, the subsystem should turn subserver tracing off.

Responding to refresh requests

Support for subsystem refresh requests is subsystem-dependent. Subsystem programmers that choose to support the **refresh** command should program their subsystems to interact with the SRC in the following manner:

- A subsystem refresh request will have the **subreq** structure `action` field set to the REFRESH constant and the **subreq** structure `object` field set to the SUBSYSTEM constant. The refresh subsystem action does not use `parm1` or `parm2`.
- When the subsystem receives the refresh request, the subsystem should reconfigure itself.

Defining your subsystem to the SRC

Subsystems are defined to the SRC object class as subsystem objects. Subservers are defined in the SRC configuration database as subserver type objects.

The structures associated with each type of object are predefined in the [sys/srcobj.h](#) file.

A subsystem object is created with the **mkssys** command or the **addssys** subroutine. A subserver type object is created with the **mkserver** command. You are not required to specify all possible options and parameters using the configuration commands and subroutines. The SRC offers pre-set defaults. You must specify only the required fields and any fields in which you want some value other than the default.

Descriptors can be added or modified at the command line by writing a shell script. They can also be added or modified using the C interface. Commands and subroutines are available for configuring and modifying the SRC objects.

Note: The choice of programming interfaces is provided for convenience only.

At the command line use the following commands:

Command	Description
<u>mkssys</u>	Adds a subsystem definition to the SRC configuration database.
<u>mkserver</u>	Adds a subserver definition to the SRC configuration database.
<u>chssys</u>	Changes a subsystem definition in the SRC configuration database.
<u>chserver</u>	Changes a subserver definition in the SRC configuration database.
<u>rmssys</u>	Removes a subsystem definition from the SRC configuration database.
<u>rmserver</u>	Removes a subserver definition from the SRC configuration database.

When using the C interface, use the following subroutines:

Subroutine	Description
<u>addssys</u>	Adds a subsystem definition to the SRC configuration database
<u>chssys</u>	Changes a subsystem definition in the SRC configuration database
<u>defssys</u>	Initializes a new subsystem definition with default values
<u>delssys</u>	Deletes an existing subsystem definition from the SRC configuration database

Note: The object code running with the **chssys** subroutine must be running with the group system.

<u>getssys</u>	Gets a subsystem definition from the SRC configuration database
<u>getsubsvr</u>	Gets a subserver definition from the SRC configuration database

The **mkssys** and **mkserver** commands call the **defssys** subroutine internally to determine subsystem and subserver default values prior to adding or modifying any values entered at the command line.

The **getssys** and **getsubsvr** subroutines are used when the SRC master program or a subsystem program needs to retrieve data from the SRC configuration files.

List of additional SRC subroutines

Use the following subroutines to program communication with the SRC and the subsystems controlled by the SRC:

Subroutine	Description
<u>src_err_msg</u>	Returns message text for SRC errors encountered by SRC library routines. (Also see threadsafe version <u>src_err_msg_r</u>)
<u>srcsbuf</u>	Requests status from the subsystem in printable format. (Also see threadsafe version <u>srcsbuf_r</u>)
<u>srcsrqt</u>	Sends a message or request to the subsystem. (Also see threadsafe version <u>srcsrqt_r</u>)
<u>srcstat</u>	Requests short subsystem status. (Also see threadsafe version <u>srcstat_r</u>)
<u>srcstathdr</u>	Gets the title text for SRC status.
<u>srcstattxt</u>	Gets the text representation for an SRC status code. (Also see threadsafe version <u>srcstattxt_r</u>)
<u>srcstop</u>	Requests termination of the subsystem.
<u>srcstrt</u>	Requests the start of a subsystem.

Trace facility

The trace facility helps you isolate system problems by monitoring selected system events or selected processes. Events that can be monitored include: entry and exit to selected subroutines, kernel routines, kernel extension routines, and interrupt handlers.

Trace can also be restricted to tracing a set of running processes or threads, or it can be used to initiate and trace a program.

When the trace facility is active, information is recorded in a system trace log file. The trace facility includes commands for activating and controlling traces and generating trace reports. Applications and kernel extensions can use several subroutines to record additional events.

For more information on the trace facility, refer to the following:

The trace facility overview

The trace facility is in the **bos.sysmgt.trace** file set. To see if this file set is installed, type the following on the command line:

```
ls1pp -l | grep bos.sysmgt.trace
```

If a line is produced which includes **bos.sysmgt.trace** then the file set is installed, otherwise you must install it.

The system trace facility records trace events which can be formatted later by the trace report command. Trace events are compiled into kernel or application code, but are only traced if tracing is active.

Tracing is activated with the **trace** command or the **trcstart** subroutine. Tracing is stopped with either the **trcstop** command or the **trcstop** subroutine. While active, tracing can be suspended or resumed with the **trcoff** and **trcon** commands, or the **trcoff** and **trcon** subroutines.

Once the trace has been stopped with **trcstop**, a trace report can then be generated with the **trcprt** command. This command uses a template file, **/etc/trcfmt**, to know how to format the entries. The templates are installed with the **trcupdate** command. For a discussion of the templates, see the [trcupdate](#) command.

Controlling the trace

The **trace** command starts the tracing of system events and controls the trace buffer and log file sizes.

There are three methods of gathering trace data.

1. The default method is to use 2 buffers to continuously gather trace data, writing one buffer while data is being put into the other buffer. The log file wraps when it becomes full.
2. The circular method gathers trace data continuously, but only writes the data to the log file when the trace is stopped. This is particularly useful for debugging a problem where you know when the problem is happening and you just want to capture the data at that time. You can start the trace at any time, and then stop it right after the problem occurs and you'll have captured the events around the problem. This method is enabled with the **-l** trace daemon flag.
3. The third option only uses one trace buffer, and quits tracing when that buffer fills, and writes the buffer to the log file. The trace is not stopped at this point, rather tracing is turned off as if a **trcoff** command had been issued. At this point you will usually want to stop the trace with the **trcstop** command. This option is most often used to gather performance data where we don't want trace to do i/o or buffer swapping until the data has been gathered. Use the **-f** flag to enable this option.

You will usually want to run the trace command asynchronously, in other words, you want to enter the trace command and then continue with other work. To run the trace asynchronously, use the **-a** flag or the **-x** flag. If you use the **-a** flag, you must then stop the trace with the **trcstop** command. If you use the **-x** flag, trace automatically stops when the program finishes.

It is usually desirable to limit the information that is traced. Use the **-j events** or **-k events** flags to specify a set of events to include (**-j**) or exclude (**-k**).

Note: When you limit the trace to specific processes or threads, you also limit the amount of information traced.

To display the program names associated with trace hooks, certain hooks must be enabled. These are specified using the **tidhk** trace event group. For example, if you want to trace the **mbuf** hook, 254, and show program names also, you need to run **trace** as follows:

```
trace -aJ tidhk -j 254
```

Tracing occurs. To stop tracing, type the following on a command line:

```
trcstop  
trcrpt -O exec=on
```

The **-O exec=on** trcrpt option shows the program names, see the **trcrpt** command for more information.

It is often desirable to specify the buffer size and the maximum log file size. The trace buffers require real memory to be available so that no paging is necessary to record trace hooks. The log file will fill to the maximum size specified, and then wrap around, discarding the oldest trace data. The **-T size** and **-L size** flags specify the size of the memory buffers and the maximum size of the trace data in the log file in bytes.

Note: Because the trace facility pins the data collection buffers, making this amount of memory unavailable to the rest of the system, the trace facility can impact performance in a memory-constrained environment. If the application being monitored is not memory-constrained, or if the percentage of memory consumed by the trace routine is small compared to what is available in the system, the impact of trace "stolen" memory should be small. If you do not specify a value, trace uses the default sizes.

Tracing can also be controlled from an application. See the **trcstart**, and **trcstop** articles.

Recording trace event data

There are two types of trace data.

generic data

consists of a data word, a buffer of opaque data and the opaque data's length. This is useful for tracing items such as path names. See the Generic Trace Channels article in the **Trace Facility Overview**. It can be found in **Trace Facility**.

Note: Tracing of specific processes or threads is only supported for channel 0. It is not supported for generic trace channels.

Non-generic data

This is what is normally traced by the AIX operating system. Each entry of this type consists of a hook word and up to 5 words of trace data. For a 64-bit application these are 8-byte words. The C programmer should use the macros TRCHKL0 through TRCHKL5, and TRCHKL0T through TRCHKL5T defined in the `/usr/include/sys/trcmacros.h` file, to record non-generic data. If these macros can not be used, see the article on the [utrchook](#) subroutine.

Generating a trace report

See the **trcrpt** command article for a full description of **trcrpt**. This command is used to generate a readable trace report from the log file generated by the **trace** command. By default the command formats data from the default log file, `/var/adm/ras/trcfile`. The **trcrpt** output is written to standard output.

To generate a trace report from the default log file, and write it to `/tmp/rptout`, enter

```
trcrpt >/tmp/rptout
```

To generate a trace report from the log file `/tmp/tlog` to `/tmp/rptout`, which includes program names and system call names, use

```
trcrpt -O exec=on,svc=on /tmp/tlog >/tmp/rptout
```

Extracting trace data from a dump

If trace was active when the system takes a dump, the trace can usually be retrieved with the **trcdead** command. To avoid overwriting the default trace log file on the current system, use the **-o output-file** option.

For example:

```
trcdead -o /tmp/tlog /var/adm/ras/vmcore.0
```

creates a trace log file `/tmp/tlog` which may then be formatted with the following:

```
trcrpt /tmp/tlog
```

Trace facility commands

The following commands are part of the trace facility:

Command	Function
<u>trace</u>	Starts the tracing of system events. With this command, you can control the size and manage the trace log file as well as the internal trace buffers that collect trace event data.
<u>trcdead</u>	Extracts trace information from a system dump. If the system halts while the trace facilities are active, the contents of the internal trace buffers are captured. This command extracts the trace event data from the dump and writes it to the trace log file.

Command

trcnm

Function

Generates a kernel name list used by the **trcrpt** command. A kernel name list is composed of a symbol table and a loader symbol table of an object file. The **trcrpt** command uses the kernel name list file to interpret addresses when formatting a report from a trace log file.

Note: It is recommended that you use the **-n trace** option instead of **trcnm**. This puts name list information into the trace log file instead of a separate file, and includes symbols from kernel extensions.

trcrpt

Formats reports of trace event data contained in the trace log file. You can specify the events to be included (or omitted) in the report, as well as determine the presentation of the output with this command. The **trcrpt** command uses the trace formatting templates stored in the **/etc/trcfmt** file to determine how to interpret the data recorded for each event.

trcstop

Stops the tracing of system events.

trcupdate

Updates the trace formatting templates stored in the **/etc/trcfmt** file. When you add applications or kernel extensions that record trace events, templates for these events must be added to the **/etc/trcfmt** file. The **trcrpt** command will use the trace formatting templates to determine how to interpret the data recorded for each event. Software products that record events usually run the **trcupdate** command as part of the installation process.

Trace facility calls and subroutines

The following calls and subroutines are part of the trace facility:

Subroutine

trcgen, trcgent

Description

Records trace events of more than five words of data. The **trcgen** subroutine can be used to record an event as part of the system event trace (trace channel 0) or to record an event on a generic trace channel (channels 1 through 7). Specify the channel number in a subroutine parameter when you record the trace event. The **trcgent** subroutine appends a time stamp to the event data. When using AIX 5L Version 5.3 with the 5300-05 Technology Level and above, the time stamp is always appended to the event data regardless of the subroutine used. Use **trcgenk** and **trcgenkt** in the kernel. C programmers should always use the **TRCGEN** and **TRCGENK** macros.

Subroutine

[utrhook](#), [utrhook64](#)

Description

Records trace events of up to five words of data. These subroutines can be used to record an event as part of the system event trace (trace channel 0). Kernel programmers can use **trhook** and **trhook64**. C programmers should always use the **TRCHKLO - TRCHKL5** and **TRCHKLOT - TRCHKL5T** macros.

If you are not using these macros, you need to build your own trace hook word. The format is documented with the [/etc/trcfmt](#) file. Note that the 32-bit and 64-bit traces have different hook word formats.

[trcoff](#)

Suspends the collection of trace data on either the system event trace channel (channel 0) or a generic trace channel (1 through 7). The trace channel remains active and trace data collection can be resumed by using the **trcon** subroutine.

[trcon](#)

Starts the collection of trace data on a trace channel. The channel can be either the system event trace channel (0) or a generic channel (1 through 7). The trace channel, however, must have been previously activated by using the **trace** command or the **trcstart** subroutine. You can suspend trace data collection by using the **trcoff** subroutine.

[trcstart](#)

Provides a library interface to the trace command. It returns the channel number of the trace it starts. If a generic channel is requested, the channel number is one of the following numbers: 1,2,3,4,5,6, 7. Otherwise the channel number is 0.

[trcstop](#)

Frees and deactivates a generic trace channel.

Trace facility files

File

[/etc/trcfmt](#)

Description

Contains the trace formatting templates used by the **trcrpt** command to determine how to interpret the data recorded for each event.

[/var/adm/ras/trcfile](#)

Contains the default trace log file. The **trace** command allows you to specify a different trace log file.

[/usr/include/sys/trchkid.h](#)

Contains trace hook identifier definitions.

[/usr/include/sys/trcmacros.h](#)

Contains commonly used macros for recording trace events.

Trace event data

See the [/etc/trcfmt](#) file for the format of the trace event data.

Trace hook identifiers

A trace hook identifier is a three- or four-digit hexadecimal number that identifies an event being traced. Prior to AIX 7.1 and on 32-bit applications running on AIX 7.1 and above, only three-digit hook identifiers can be used. When using a tracing macro such as `TRCHKL1`, you specify the trace hook as:

```
hhh00000
```

where `hhh` is the hook id.

On 64-bit applications and kernel routines running on AIX 7.1 and above, three- and four-digit identifiers can be used. When using a tracing macro such as `TRCHKL1`, you specify the trace hook as:

```
hhhh0000
```

where `hhhh` is the hook id.

Note: If a four-digit identifier is used and the identifier is less than `0x1000`, the least-significant digit must be 0 (of the form `0x0hh0`).

A three-digit identifier has an implicit 0 in its least-significant digit so that a 32-bit hook identifier is equivalent to a 64-bit hook of the form `hhh0`.

Most trace hook identifiers are defined in the `/usr/include/sys/trchkid.h` file. The values `0x0100` through `0x0FF0` are available for use by 64-bit user applications. The values `0x010` through `0x0FF` are available for use by 32-bit user applications. All other values are reserved for system use. The currently defined trace hook identifiers can be listed using the `trcrpt -j` command.

Trace facility generic trace channels

The trace facility supports up to eight active trace sessions at a time. Each trace session uses a channel of the multiplexed trace special file, `/dev/systrace`. Channel 0 is used by the trace facility to record system events. The tracing of system events is started and stopped by the `trace` and `trcstop` commands. If you trace specific processes or threads, or if a program is traced, only channel 0 is used. Channels 1 through 7 are referred to as generic trace channels and can only be used by subsystems for other types of tracing such as data link tracing.

To implement tracing using the generic trace channels of the trace facility, a subsystem calls the `trcstart` subroutine to activate a trace channel and to determine the channel number. The subsystem modules can then record trace events using the `TRCGEN` or `TRCGENT` macros, or if necessary, `trcgen`, `trcgent`, `trcgenk`, or `trcgenkt` subroutine. The channel number returned by the `trcstart` subroutine is one of the parameters that must be passed to these subroutines. The subsystem can suspend and resume trace data collection using the `trcoff` and `trcon` subroutines and can deactivate a trace channel using the `trcstop` subroutine. The subsystem must provide the user interface to activate and deactivate subsystem tracing.

The trace hook IDs, most of which are stored in the `/usr/include/sys/trchkid.h` file, and the trace formatting templates, which are stored in the `/etc/trcfmt` file, are shared by all the trace channels.

Related information

[trace](#)

[trcdead](#)

[trcnm](#)

[trcrpt](#)

[trcstop](#)

[trcupdate](#)

Start the trace facility

Use the following procedures to configure and start a system trace:

Configuring the trace command

The **trace** command starts the tracing of system events and controls the size of and manages the **trace** log file, as well as the internal trace buffers that collect trace event data. For information on the syntax of this command, see [trace Daemon](#).

Recording trace event data

The data recorded for each traced event consist of a word containing the trace hook identifier and the hook type followed by a variable number of words of trace data optionally followed by a time stamp. The word containing the trace hook identifier and the hook type is called the hook word. The remaining two bytes of the hook word are called hook data and are available for recording event data.

Trace hook identifiers

A trace hook identifier is a three- or four-digit hexadecimal number that identifies an event being traced. Prior to AIX 6.1 and on 32-bit applications running on AIX 6.1 and above, only three-digit hook identifiers can be used. When using a tracing macro such as TRCHKL1, the trace hook is specified as:

```
hhh0000
```

where hhh is the hook identifier.

On 64-bit applications and kernel routines running on AIX 6.1 and above, three- and four-digit hook identifiers can be used. When using a tracing macro such as TRCHKL1, the trace hook is specified as:

```
hhhh0000
```

where hhhh is the hook identifier.

Note: If a four-digit identifier is used and the identifier is less than 0x1000, the least-significant digit must be 0 (of the form 0x0hh0).

A three-digit identifier has an implicit 0 in its least-significant digit such that a 32-bit hook identifier of hhh is equivalent to a 64-bit hook identifier of hhh0.

Most trace hook identifiers are defined in the `/usr/include/sys/trchkid.h` file. The values 0x0100 through 0x0FF0 are available for use by 64-bit applications. The values 0x010 through 0x0FF are available for use by 32-bit user applications. All other values are reserved for system use. The currently defined trace hook identifiers can be listed with the **trcrpt -j** command.

Using generic trace channels

The trace facility supports up to eight active trace sessions at a time. Each trace session uses a channel of the multiplexed trace special file, `/dev/systrace`. Channel 0 is used by the trace facility to record system events. The tracing of system events is started and stopped by the **trace** and **trcstop** commands. Channels 1 through 7 are referred to as generic trace channels and may be used by subsystems for other types of tracing such as data link tracing.

To implement tracing using the generic trace channels of the trace facility, a subsystem calls the **trcstart** subroutine to activate a trace channel and to determine the channel number. The subsystem modules can then record trace events TRCGEN or TRCGENT macros, or the `trcgen`, `trcgent`, `trcgenk`, or `trcgenkt` subroutines. The channel number returned by the **trcstart** subroutine is one of the parameters that must be passed to these subroutines. The subsystem can suspend and resume trace data collection using the **trcoff** and **trcon** subroutines and can deactivate a trace channel using the **trcstop** subroutine. The trace events for each channel are written to a separate **trace** log file, which, by default, is named `/var/adm/ras/trcfile.n`, where *n* is the channel number. The subsystem must provide the user interface to activate and deactivate subsystem tracing.

Starting a trace

Use one of the following procedures to start the trace facility.

- Start the trace facility by using the **trace** command.

Start the trace asynchronously. For example:

```
trace -a
mycmd
trcstop
```

When using the trace facility asynchronously, use the **trace** daemon to trace the selected system events (such as the **mycmd** command); then, use the **trcstop** command to stop the trace.

OR

Start the trace interactively. For example:

```
trace
->!mycmd
->quit
```

When using the trace facility interactively, get into the interactive mode as denoted by the **->** prompt, and use the **trace** subcommands (such as **!**) to trace the selected system events. Use the **quit** subcommand to stop the trace.

- Use **smit trace**, and choose the **Start Trace** option.

```
smit trace
```

Stopping a trace

Use one of the following procedures to stop the trace you started earlier.

- When using **trace** asynchronously at the command line, use the **trcstop** command:

```
trace -a
mycmd
trcstop
```

When using the trace facility asynchronously, use the **trace** daemon to trace the selected system events (such as the **mycmd** command); then, use the **trcstop** command to stop the trace.

- When using **trace** interactively at the command line, use the **quit** subcommand:

```
trace
->!mycmd
->quit
```

The interactive mode is denoted by the **->** prompt. Use the **trace** subcommands (such as **!**) to trace the selected system events. Use the **quit** subcommand to stop the trace.

- Use **smit trace** and choose the **Stop Trace** option:

```
smit trace
```

Generating a trace report

Use either of the following procedures to generate a report of events that have been traced.

- Use the **trcrpt** command:

```
trcrpt>/tmp/NewFile
```

The previous example formats the **trace** log file and sends the report to **/tmp/newfile**. The **trcrpt** command reads the **trace** log file, formats the trace entries, and writes a report.

- Use the **smit trcrpt** command:

```
smit trcrpt
```

Tracing user application

This topic describes tracing user applications.

The tracing operation depends on three logically different processes: the traced process, the controller process, and the analyzer process. A process can simultaneously be the traced process, the controller process and the analyzer process. When the traced process is running and a trace point is reached, a trace event is recorded into the trace streams created for that process if the trace event type identifier associated with this process is not filtered out.

The controller process controls the recording of trace events into trace streams. The controller process performs the following operations for an active trace stream:

- Initializes the attributes of a trace stream.
- Creates the trace stream for a specified traced process using the initialized attributes.
- Starts and stops tracing for a trace stream.
- Filters the type of trace events to be recorded.
- Ends a trace stream.

The analyzer process retrieves the traced events either at runtime, when the trace stream is active and is recording trace events; or after opening a trace log that is previously recorded and closed.

The **posix_trace_create**, **posix_trace_create_withlog**, and the **posix_trace_open** subroutines create a trace stream identifier. The **posix_trace_create** and **posix_trace_create_withlog** subroutines are used only by a controller process. The **posix_trace_open** subroutine is used only by an analyzer process.

A traced process contains a mapping of the trace event names to the trace event type identifiers that are defined for the process. An active trace event records the system predefined trace event types, such as the `POSIX_TRACE_START`, and the trace event types that are defined for the traced processes but not filtered out by the trace stream. To define the mapping, call the **posix_trace_eventid_open** subroutine from the instrumented application, or call the **posix_trace_trid_eventid_open** subroutine from the controller process. For a prerecorded trace stream, the list of trace event types is obtained from the prerecorded trace log.

The tracing subroutines can be used in debugging the possibly pre-instrumented code and postmortem fault analysis. Debugging pre-instrumented code might require pre-filtering capabilities to avoid overwhelming the trace stream and permits focusing on expected information. Postmortem fault analysis requires comprehensive trace capabilities to be able to record all types of information.

The events to be traced belong to the following two classes:

- User trace events that the instrumented application generated.
- System trace events that the operating system generated in correspondence with the trace control operation.

In a file associated with an active trace stream, the `st_ctime` and `st_mtime` fields are marked for update every time any of the tracing operations modifies that file.

In a file associated with a trace stream, the `st_atime` field is marked for update every time any of the tracing operations causes data to be read from that file.

If an application performs any operation on a file descriptor associated with an active or prerecorded trace stream, the results are not defined until calling the **posix_trace_shutdown** or **posix_trace_close** subroutines for that trace stream.

Tracing data structures

This section describes tracing data structure.

The `<trace.h>` header file defines the **posix_trace_status_info** structure and **posix_trace_event_info** structures.

posix_trace_status_info structure

To facilitate the control of a trace stream, call the **posix_trace_get_status** subroutine to dynamically obtain the information about the current state of an active trace stream.

The **posix_trace_status_info** structure defined in the <trace.h> file contains the following members:

Member type	Member name	Description
int	posix_stream_status	The operating mode of the trace stream.
int	posix_stream_full_status	The full status of the trace stream.
int	posix_stream_overrun_status	Indicates whether trace events were lost in the trace stream.
int	posix_stream_flush_status	Indicates whether a flush is in progress.
int	posix_stream_flush_error	Indicates whether any error occurred during the last flush operation.
int	posix_log_overrun_status	Indicates whether trace events were lost in the trace log.
int	posix_log_full_status	The full status of the trace log.

The **posix_stream_status** member indicates the operating mode of the trace stream. It can have one of the following values defined by the manifest constants in the <trace.h> header:

POSIX_TRACE_RUNNING

The tracing is in progress. The trace stream is accepting trace events.

POSIX_TRACE_SUSPENDED

The trace stream is not accepting trace events. The tracing operation has not started or has stopped, either following a **posix_trace_stop** subroutine call or because the trace resources are exhausted.

The **posix_stream_full_status** member indicates the full status of the trace stream. It can have one of the following values defined by manifest constants in the <trace.h> header:

POSIX_TRACE_FULL

The space in the trace stream for trace events is exhausted.

POSIX_TRACE_NOT_FULL

The space in the trace stream is not full.

The combination of the **posix_stream_status** and **posix_stream_full_status** members indicate the actual status of the stream. The status can be interpreted as follows:

POSIX_TRACE_RUNNING and POSIX_TRACE_NOT_FULL

The tracing is in progress, and the space is available for recording more trace events.

POSIX_TRACE_RUNNING and POSIX_TRACE_FULL

The tracing is in progress, and the trace stream is full of trace events. If the stream full policy is set to **POSIX_TRACE_LOOP**, this status might occur. The trace stream contains trace events recorded during a moving time window of previous trace events, and some previous trace events might be overwritten and lost.

POSIX_TRACE_SUSPENDED and POSIX_TRACE_NOT_FULL

The tracing has not started, the **posix_trace_stop** subroutine stopped the tracing, or the **posix_trace_clear** subroutine cleared the tracing.

POSIX_TRACE_SUSPENDED and POSIX_TRACE_FULL

The tracing stops. It stops because the stream full policy attribute is set to **POSIX_TRACE_UNTIL_FULL** and the trace resources is exhausted, or because calling the **posix_trace_stop** subroutine when the trace resource is exhausted.

The **posix_stream_overrun_status** member indicates whether trace events are lost in the trace stream. This member can have one of the following values that are defined by manifest constants in the <trace.h> header:

POSIX_TRACE_OVERRUN

At least one trace event is lost and is not recorded in the trace stream.

POSIX_TRACE_NO_OVERRUN

No trace events are lost.

When the corresponding trace stream is created, the **posix_stream_outrun_status** member is set to **POSIX_TRACE_NO_OVERRUN**. When an overrun occurs, the status is set to **POSIX_TRACE_OVERRUN**.

An overrun occurs in the following situations:

- The policy is **POSIX_TRACE_LOOP** and a recorded trace event is overwritten.
- The policy is **POSIX_TRACE_UNTIL_FULL** and the trace stream is full when a trace event is generated.
- The policy is **POSIX_TRACE_FLUSH** and more than one trace event is lost while flushing the trace stream to the trace log.

The **posix_stream_outrun_status** member is reset to zero after its value is read.

The **posix_stream_flush_status** member indicates whether a flush operation is being performed and can have one of the following values that the manifest constants in the <trace.h> header file define:

POSIX_TRACE_FLUSHING

The trace stream is being flushed to the trace log.

POSIX_TRACE_NOT_FLUSHING

The **posix_stream_flush_status** member is set to **POSIX_TRACE_NOT_FLUSHING** if no flush operation is in progress.

The **posix_stream_flush_status** member is set to **POSIX_TRACE_FLUSHING** in the following situations:

- A flush operation is in progress because calling the **posix_trace_flush** subroutine.
- A flush operation is in progress because the trace stream becomes full with the stream-full-policy attribute set to **POSIX_TRACE_FLUSH**.

The **posix_stream_flush_error** member is set to zero if no error occurs in the flushing. If an error occurs in a previous flushing operation, the **posix_stream_flush_error** member is set to the value of the first error that occurs. If more than one error occurs in the flushing operation, the first error value is used and the others are discarded. The **posix_stream_flush_error** member is reset to zero after its value is read.

The **posix_log_outrun_status** member indicates whether trace events are lost in the trace log. This member can have one of the following values that the manifest constants in the <trace.h> header defines:

POSIX_TRACE_OVERRUN

At least one trace event is lost.

POSIX_TRACE_NO_OVERRUN

No trace event is lost.

Whenever an overrun occurs, the **posix_log_outrun_status** member is set to **POSIX_TRACE_OVERRUN**. When the corresponding trace stream is created, the **posix_log_outrun_status** member is set to **POSIX_TRACE_NO_OVERRUN**.

The **posix_log_outrun_status** member is reset to zero after its value is read.

If an active trace stream is created by the **posix_trace_create** subroutine and it has no log, the **posix_log_outrun_status** member is set to **POSIX_TRACE_NO_OVERRUN**.

The **posix_log_full_status** member indicates the full status of the trace log, and it can have one of the following values defined by manifest constants in the <trace.h> header:

POSIX_TRACE_FULL

The space in the trace log is full.

POSIX_TRACE_NOT_FULL

The space in the trace log is not full.

The **posix_log_full_status** member is meaningful only when the log-full-policy attribute is either **POSIX_TRACE_UNTIL_FULL** or **POSIX_TRACE_LOOP**.

If an active trace stream is created by the **posix_trace_create** subroutine and it has no log, the **posix_log_full_status** member is set to `POSIX_TRACE_NOT_FULL`.

posix_trace_event_info Structure

The trace event structure **posix_trace_event_info** contains the information for a recorded trace event. The **posix_trace_getnext_event**, **posix_trace_timedgetnext_event**, and **posix_trace_trygetnext_event** subroutines return this subroutine.

The **posix_trace_event_info** structure that is defined in the `<trace.h>` header file contains the following members:

Member type	Member name	Description
trace_event_id_t	posix_event_id	Trace event type identification.
pid_t	posix_pid	Process ID of the process that generates the trace event.
void*	posix_prog_address	Address where the trace point is invoked.
int	posix_truncation_status	Truncation status of the data associated with this trace event.
struct timespec	posix_timestamp	Time at which the trace event is generated.
pthread_t	posix_thread_id	ID of the thread that generates the trace event.

The **posix_event_id** member represents the identification of the trace event type. You cannot define the member value directly. To return the identification, run one of the following subroutines:

- `posix_trace_trid_eventid_open`
- `posix_trace_eventtypelist_getnext_id`
- `posix_trace_eventid_open`

To obtain the name of the trace event type, run the **posix_trace_eventid_get_name** subroutine.

The **posix_pid** is the process identifier of the traced process that generates the trace event. If the **posix_event_id** member is one of the system trace events and the trace event is not associated with any process, the **posix_pid** member is set to zero.

For a user trace event, the **posix_prog_address** member is the mapped address of a process and the point at which the associated call to the **posix_trace_event** subroutine is made.

The **posix_truncation_status** member defines the truncation status of the data associated with a trace event. The **posix_truncation_status** member can have one of the following values that the manifest constants in the `<trace.h>` header define:

POSIX_TRACE_NOT_TRUNCATED

All of the traced data is available.

POSIX_TRACE_TRUNCATED_RECORD

Data is truncated when the trace event is generated.

POSIX_TRACE_TRUNCATED_READ

Data is truncated when the trace event is read from a trace stream or a trace log. This truncation status overrides the `POSIX_TRACE_TRUNCATED_RECORD` status.

The **posix_timestamp** member defines the time when the trace event is generated. The clock used is `CLOCK_REALTIME`. To retrieve the resolution of this clock, call the **posix_trace_attr_getclockres** subroutine.

The **posix_thread_id** member is the identifier of the thread that generates the trace event. If the **posix_event_id** member is one of the system trace events and that trace event is not associated with any thread, the **posix_thread_id** member is set to zero.

Trace stream attributes

The following trace stream attributes compose the **posix_trace_attr_t** trace stream attributes object:

- The *genversion* attribute identifies the origin and version of the trace system.
- The *tracename* attribute is a character string defined by the trace controller. It identifies the trace stream.
- The *creation-time* attribute represents the time when the trace stream is created.
- The *clock-resolution* attribute defines the clock resolution of the clock that is used to generate timestamps.
- The *stream-min-size* attribute defines the minimum size in byte of the trace stream. The size is strictly reserved for the trace events.
- The *stream-full-policy* attribute defines the policy to follow when the trace stream is full. Its value is `POSIX_TRACE_LOOP`, `POSIX_TRACE_UNTIL_FULL`, or `POSIX_TRACE_FLUSH`.
- The *max-data-size* attribute defines the maximum size in byte of the record of a trace event.
- The *inheritance* attribute specifies whether a newly created trace stream inherits tracing in its parent's process trace stream. The value of this attribute is either `POSIX_TRACE_INHERITED` or `POSIX_TRACE_CLOSE_FOR_CHILD`.
- The *log-max-size* attribute defines the maximum size in byte of the trace log associated with an active trace stream. Other stream data is not included in this size.
- The *log-full-policy* attribute defines the policy of a trace log associated with an active trace stream to be `POSIX_TRACE_LOOP`, `POSIX_TRACE_UNTIL_FULL`, or `POSIX_TRACE_APPEND`.

Trace event type definitions

The `<trace.h>` header file defines system trace event types and user trace event types. This topic describes the trace event types and the definitions.

System trace event type definitions

To interpret a trace stream or a trace log, a trace analyzer process requires information on trace events and information on system trace events that report the invocation of trace operations.

The following system trace event types track the invocation of the trace operations:

POSIX_TRACE_START

Tracks the invocation of a trace start operation.

POSIX_TRACE_STOP

Tracks the invocation of a trace stop operation.

POSIX_TRACE_FILTER

Tracks the invocation of a change operation on a trace event type.

The following system trace event types report the trace events that are recorded by the Posix Trace Library subroutines:

POSIX_TRACE_OVERFLOW

A trace overflow begins.

POSIX_TRACE_RESUME

A trace overflow ends.

POSIX_TRACE_FLUSH_START

A flush operation begins.

POSIX_TRACE_FLUSH_STOP

A flush operation ends.

POSIX_TRACE_ERROR

A trace error occurs.

The POSIX_TRACE_START and POSIX_TRACE_STOP trace events specify the elapsed time that the trace stream is running.

The POSIX_TRACE_STOP trace event with a value of zero indicates a call of the **posix_trace_stop** subroutine.

The POSIX_TRACE_STOP trace event with a value of nonzero indicates an automatic stop of the trace stream. For more information, see the **posix_trace_attr_getstreamfullpolicy** subroutine.

The POSIX_TRACE_FILTER trace event indicates that the value of the filter of the trace event types changes when the trace stream is running.

The POSIX_TRACE_ERROR indicates that an internal error of the trace system occurs.

The POSIX_TRACE_OVERFLOW trace event is reported with a timestamp equal to the timestamp of the first trace event that is overwritten. This trace event indicates that some generated trace events are lost.

The POSIX_TRACE_RESUME trace event indicates that the trace system is recording trace events after an overflow.

A constant with the name of a trace event and a **trace_event_id_t** constant defines a trace event type. Trace event data is associated with some of these trace events.

The System Trace Events table describes the system trace events that are predefined.

Event name	Constant	Associated data	Data type
posix_trace_error	POSIX_TRACE_ERROR	Error	int
posix_trace_start	POSIX_TRACE_START	event_filter	trace_event_set_t
posix_trace_stop	POSIX_TRACE_STOP	Auto	int
posix_trace_filter	POSIX_TRACE_FILTER	old_event_filter, new_event_filter	trace_event_set_t
posix_trace_overflow	POSIX_TRACE_OVERFLOW	none	none
posix_trace_resume	POSIX_TRACE_RESUME	none	none
posix_trace_flush_start	POSIX_TRACE_FLUSH_START	none	none
posix_trace_flush_stop	POSIX_TRACE_FLUSH_STOP	none	none

User trace event type definitions

The <trace.h> header file defines the user trace event POSIX_TRACE_UNNAMED_USEREVENT. If the limit of TRACE_USER_EVENT_MAX is reached, the POSIX_TRACE_UNNAMED_USEREVENT user event is returned when an application tries to register more events than it is allowed. No data is associated with this user trace event.

The POSIX_TRACE_UNNAMED_USEREVENT constant is predefined for the **posix_trace_unnamed_userevent** trace event name.

Trace subroutines

This section describes the trace subroutines and the different tracing roles they support.

The trace interface is built and structured to improve portability by using trace data of opaque type.

To set up and customize resources for a trace controller process to run a trace stream, use the subroutines in the following table.

Subroutine	Purpose
posix_trace_attr_init	Initializes attributes
posix_trace_attr_destroy	Destroys attributes
posix_trace_attr_getgenversion	Gets the version of a trace stream
posix_trace_attr_getname	Gets the trace name
posix_trace_attr_setname	Sets the trace name
posix_trace_attr_getinherited	Gets the inheritance policy of a trace stream
posix_trace_attr_setinherited	Sets the inheritance policy of a trace stream
posix_trace_attr_setstreamfullpolicy	Sets the stream full policy
posix_trace_attr_getstreamfullpolicy	Gets the stream full policy
posix_trace_attr_setlogfullpolicy	Sets the log full policy of a trace stream
posix_trace_attr_getlogfullpolicy	Gets the log full policy of a trace stream
posix_trace_attr_setlogsize	Sets the size of the log of a trace stream
posix_trace_attr_setmaxdatasize	Sets the maximum size of the user trace event data
posix_trace_attr_setstreamsize	Sets the size of a trace stream
posix_trace_attr_getmaxusereventsize	Gets the maximum size of an user event for a given length
posix_trace_create	Creates an active trace stream
posix_trace_create_withlog	Creates an active trace stream and associates it with a trace log
posix_trace_flush	Copies the content of a trace stream into the associated trace log of the trace stream
posix_trace_shutdown	Shuts down a trace stream
posix_trace_clear	Clears a trace stream and a trace log
posix_trace_trid_eventid_open	Associates a trace event type identifier to a user trace event name
posix_trace_eventid_equal	Compares two trace event type identifiers
posix_trace_eventid_get_name	Retrieves the trace event name from a trace event type identifier
posix_trace_eventtypelist_getnext_id, posix_trace_eventtypelist_rewind	Iterate over a mapping of trace event types
posix_trace_eventset_add	Adds a trace event type in a trace event type set
posix_trace_eventset_empty	Empties a trace event type set
posix_trace_eventset_del	Deletes a trace event type from a trace event type set
posix_trace_eventset_fill	Fills in a trace event type set
posix_trace_eventset_ismember	Tests if the trace event type is included in the trace event type set

Subroutine	Purpose
posix_trace_get_filter	Gets the filter of an initialized trace stream
posix_trace_set_filter	Sets the filter of an initialized trace stream
posix_trace_start	Starts a trace stream
posix_trace_stop	Stops a trace stream
posix_trace_get_attr	Reads trace stream information
posix_trace_get_status	Gets trace attributes or trace status

To set trace point for a traced process, use the **posix_trace_event** and **posix_trace_eventid_open** subroutines. These subroutines define identifiers of trace event types and insert trace points.

The following table shows the subroutines that retrieve information from a trace stream and trace log for a trace analyzer process.

Subroutine	Purpose
posix_trace_attr_getname	Gets the trace name
posix_trace_attr_getgenversion	Reads identification information
posix_trace_attr_getcreatetime	Gets the creation time of a trace stream
posix_trace_attr_getinherited	Gets the inheritance policy of a trace stream
posix_trace_attr_getstreamfullpolicy	Gets the stream full policy
posix_trace_attr_getlogfullpolicy	Gets the log full policy of a trace stream
posix_trace_attr_getmaxusereventsize	Gets the maximum size of an user event for a given length
posix_trace_attr_getmaxsystemeventsize	Gets the maximum size of a system trace event
posix_trace_attr_getlogsize	Gets the size of the log of a trace stream
posix_trace_attr_getmaxdatasize	Gets the maximum size of user trace event data
posix_trace_attr_getstreamsize	Gets the size of a trace stream
posix_trace_trid_eventid_open	Associates a trace event type identifier to a user trace event name
posix_trace_eventid_equal	Compares two trace event type identifiers
posix_trace_eventid_get_name	Gets the trace event name from a trace event type identifier
posix_trace_eventtypelist_getnext_id, posix_trace_eventtypelist_rewind	Iterate over a mapping of trace event types
posix_trace_open	Opens a trace log
posix_trace_close	Closes a trace log
posix_trace_rewind	Re-initializes a trace log for reading
posix_trace_get_attr	Reads the information of a trace stream
posix_trace_get_status	Reads the status of a trace stream
posix_trace_getnext_event	Reads a trace event

Subroutine	Purpose
<code>posix_trace_timedgetnext_event</code> , <code>posix_trace_trygetnext_event</code>	Gets a trace event

tty subsystem

AIX is a multiuser operating system that allows user access from local or remotely attached devices. The communication layer that supports this function is the tty subsystem.

The communication between terminal devices and the programs that read and write to them is controlled by the tty interface. Examples of tty devices are:

- Modems
- ASCII terminals
- System consoles
- Serial printers
- Xterms or aixterms under X-Windows

This overview provides information on the following topics:

TTY subsystem objectives

The tty subsystem is responsible for:

- Controlling the physical flow of data on asynchronous lines (including the transmission speed, character size, and line availability)
- Interpreting the data by recognizing special characters and adapting to national languages
- Controlling jobs and terminal access by using the concept of controlling terminal

A controlling terminal manages the input and output operations of a group of processes. The **tty** special file supports the controlling terminal interface. In practice, user programs seldom open terminal files, such as **dev/tty5**. These files are opened by a **getty** or **rlogind** command and become the user's standard input and output devices.

See **tty Special File** in *Files Reference* for more information about the controlling terminal.

tty subsystem modules

To perform its tasks, the tty subsystem is composed of modules, or disciplines. A module is a set of processing rules that govern the interface for communication between the computer and an asynchronous device. Modules can be added and removed dynamically for each tty.

The tty subsystem supports three main types of modules:

tty drivers

tty drivers, or hardware disciplines, directly control the hardware (tty devices) or pseudo-hardware (pty devices). They perform the actual input and output to the adapter by providing services to the modules above it. The services are flow control and special semantics when a port is being opened.

The following tty drivers are provided:

Driver	Description
cxma	128-port asynchronous PCI controller.
cxpa	8-port asynchronous PCI controller.
lft	Low-function terminal. The tty name is /dev/lftY , where $Y \geq 0$.
pty	<u>pseudo-terminal device driver</u> .

Driver	Description
sa	2-port asynchronous EIA-232 PCI adapter.
sf	Universal asynchronous receiver/transceivers (UARTs) on system planar.

The “TTY drivers” on page 788 section provides more information.

Line disciplines

Line disciplines provide editing, job control, and special character interpretation. They perform all transformations that occur on the inbound and outbound data streams. Line disciplines also perform most of the error handling and status monitoring for the tty driver.

The following line disciplines are provided:

Subroutine	Description
ldterm	Terminal devices
spttr	Serial printer (splp command)
slip	Serial Line Internet Protocol (slattach command)

Converter modules

Converter modules, or *mapping disciplines*, translate, or map, input and output characters.

The following converter modules are provided:

Converter	Description
nls	National language support for terminal mapping; this converter translates incoming and outgoing characters on the data stream, based on the input and output maps defined for the port (see the setmaps command)
lc_sjis and uc_sjis	Upper and lower converter used to translate multibyte characters between the Shifted Japanese Industrial Standard (SJIS) and the Advanced Japanese EUC Code (AJEC) handled by the ldterm line discipline.

“Converter modules” on page 787 provides more information on converters.

TTY subsystem structure

The tty subsystem is based on STREAMS. This STREAMS-based structure provides modularity and flexibility, and enables the following features:

- Easy customizing; users can customize their terminal subsystem environment by adding and removing modules of their choice.
- Reusable modules; for example, the same line discipline module can be used on many tty devices with different configurations.
- Easy addition of new features to the terminal subsystem.
- Providing an homogeneous tty interface on heterogeneous devices.

The structure of a tty stream is made up of the following modules:

- The stream head, processing the user's requests. The stream head is the same for all tty devices, regardless of what line discipline or tty driver is in use.
- An optional upper converter (**uc_sjis** for example), a converter module pushed above the line discipline to convert upstream and downstream data.
- The line discipline.

- An optional lower converter (**lc_sjis** for example), a converter module pushed below the line discipline to convert upstream and downstream data.
- An optional character mapping module (**nls**), a converter module pushed above the tty driver to support input and output terminal mapping.
- The stream end: a tty driver.

Unless required, the internationalization modules are not present in the tty stream.

For a serial printer, the internationalization modules are usually not present on the stream; therefore, the structure is simpler.

Common services

The **/usr/include/sys/ioctl.h** and **/usr/include/termios.h** files describe the interface to the common services provided by the tty subsystem. The **ioctl.h** file, which is used by all of the modules, includes the **winsize** structure, as well as several ioctl commands. The **termios.h** file includes the POSIX compliant subroutines and data types.

The provided services are grouped and discussed here according to their specific functions.

Hardware control services

The following subroutines are provided for hardware control:

Subroutine	Description
<u>cfgetispeed</u>	Gets input baud rate
<u>cfgetospeed</u>	Gets output baud rate
<u>cfsetispeed</u>	Sets input baud rate
<u>cfsetospeed</u>	Sets output baud rate
<u>tcsendbreak</u>	Sends a break on an asynchronous serial data line

Flow control services

The following subroutines are provided for flow control:

Subroutine	Description
<u>tcdrain</u>	Waits for output to complete
<u>tcflow</u>	Performs flow control functions
<u>tcflush</u>	Discards data from the specified queue

Terminal information and control

The following subroutines are provided for terminal information and control:

Subroutine	Description
<u>isatty</u>	Determines if the device is a terminal
<u>setcsmap</u>	Reads a code set map file and assigns it to the standard input device
<u>tcgetattr</u>	Gets terminal state
<u>tcsetattr</u>	Sets terminal state
<u>ttylock</u>, <u>ttywait</u>, <u>ttyunlock</u>, or <u>ttylocked</u>	Controls tty locking functions

Subroutine	Description
<u>ttyname</u>	Gets the name of a terminal

Window and terminal size services

The kernel stores the **winsize** structure to provide a consistent interface for the current terminal or window size. The **winsize** structure contains the following fields:

Field	Description
<code>ws_row</code>	Indicates the number of rows (in characters) on the window or terminal
<code>ws_col</code>	Indicates the number of columns (in characters) on the window or terminal
<code>ws_xpixel</code>	Indicates the horizontal size (in pixels) of the window or terminal
<code>ws_ypixel</code>	Indicates the vertical size (in pixels) of the window or terminal

By convention, a value of 0 in all of the **winsize** structure fields indicates that the structure has not yet been set up.

Subroutine	Description
<u>termdef</u>	Queries terminal characteristics.
TIOCGWINSZ	Gets the window size. The argument to this ioctl operation is a pointer to a winsize structure, into which the current terminal or window size is placed.
TIOCSWINSZ	Sets the window size. The argument to this ioctl operation is a pointer to a winsize structure, which is used to set the current terminal or window size information. If the new information differs from the previous, a SIGWINCH signal is sent to the terminal process group.

Process group management services

The following subroutines are provided for process group management:

Subroutine	Description
<u>tcgetpgrp</u>	Gets foreground process group ID
<u>tcsetpgrp</u>	Sets foreground process group ID

Buffer size operations

The following ioctl operations are used for setting the size of the terminal input and output buffers. The argument to these operations is a pointer to an integer specifying the size of the buffer.

Operations	Description
TXSETIHOG	Sets the hog limit for the number of input characters that can be received and stored in the internal tty buffers before the process reads them. The default hog limit is 8192 characters. Once the hog limit plus one character is reached, an error is logged in the error log and the input buffer is flushed. The hog number should not be too large, since the buffer is allocated from the system-pinned memory.

Operations	Description
TXSETOHOG	Sets the hog limit for the number of output characters buffered to echo input. The default hog limit is 8192 characters. Once the hog output limit is reached, input characters are no longer echoed. The hog number should not be too large, since the buffer is allocated from the system-pinned memory.

Synchronization

The tty subsystem takes advantage of the synchronization provided by STREAMS. The tty stream modules are configured with the queue pair level synchronization. This synchronization allows the parallelization of the processing for two different streams.

Related information

[rlogind](#)
[setmaps](#)
[stty](#)
[xdm](#)
[eucioctl.h](#)
[lft](#)
[pty](#)
[setmaps file](#)
[termios.h](#)
[tty file](#)

Line discipline module (ldterm)

The **ldterm** line discipline is the common line discipline for terminals.

This line discipline is POSIX compliant and also ensures compatibility with the BSD interface. The latter line discipline is supported only for compatibility with older applications. For portability reasons, it is strongly recommended that you use the POSIX interface in new applications.

This section describes the features provided by the **ldterm** line discipline. For more information about controlling **ldterm**, see "[termios.h File](#)" in *Files Reference*

Terminal parameters

The parameters that control certain terminal I/O characteristics are specified in the **termios** structure as defined in the [termios.h](#) file. The **termios** structure includes (but is not limited to) the following members:

tcflag_t c_iflag

Input modes

tcflag_t c_oflag

Output modes

tcflag_t c_cflag

Control modes

tcflag_t c_lflag

Local modes

cc_t c_cc[NCCS]

Control characters.

The **tcflag_t** and **cc_t** unsigned integer types are defined in the [termios.h](#) file. The **NCCS** symbol is also defined in the [termios.h](#) file.

Process group session management (job control)

A controlling terminal distinguishes one process group in the session, with which it is associated, to be the foreground process group. All other process groups in the session are designated as background process groups. The foreground process group plays a special role in handling signals.

Command interpreter processes that support job control, such as the Korn shell (the **ksh** command) and the C shell (the **cs**h command), can allocate the terminal to different *jobs*, or process groups, by placing related processes in a single process group and associating this process group with the terminal. A terminal's foreground process group can be set or examined by a process, assuming the permission requirements are met. The terminal driver assists in job allocation by restricting access to the terminal by processes that are not in the foreground process group.

Terminal access contro

If a process that is not in the foreground process group of its controlling terminal attempts to read from the controlling terminal, the process group of that process is sent a **SIGTTIN** signal. However, if the reading process is ignoring or blocking the **SIGTTIN** signal, or if the process group of the reading process is orphaned, the read request returns a value of -1, sets the **errno** global variable to **EIO**, and does not send a signal.

If a process that is not in the foreground process group of its controlling terminal attempts to write to the controlling terminal, the process group of that process is sent a **SIGTTOU** signal. However, the management of the **SIGTTOU** signal depends on the **TOSTOP** flag which is defined in the **c_lflag** field of the **termios** structure. If the **TOSTOP** flag is not set, or if the **TOSTOP** flag is set and the process is ignoring or blocking the **SIGTTOU** signal, the process is allowed to write to the terminal, and the **SIGTTOU** signal is not sent. If the **TOSTOP** flag is set, the process group of the writing process is orphaned, and the writing process is not ignoring or blocking the **SIGTTOU** signal, then the write request returns a value of -1, sets the **errno** global variable to **EIO**, and does not send a signal.

Certain functions that set terminal parameters (**tcsetattr**, **tcsendbreak**, **tcflow**, and **tcflush**) are treated in the same manner as write requests, except that the **TOSTOP** flag is ignored. The effect is identical to that of terminal write requests when the **TOSTOP** flag is set.

Reading data and input processing

Two general kinds of input processing are available, depending on whether the terminal device file is in canonical or noncanonical mode. Additionally, input characters are processed according to the **c_iflag** and **c_lflag** fields. Such processing can include *echoing*, or the transmitting of input characters immediately back to the terminal that sent them. Echoing is useful for terminals that can operate in full-duplex mode.

A read request can be handled in two ways, depending on whether the **O_NONBLOCK** flag is set by an **open** or **fcntl** subroutine. If the **O_NONBLOCK** flag is not set, the read request is blocked until data is available or until a signal is received. If the **O_NONBLOCK** flag is set, the read request is completed, without blocking, in one of three ways:

- If there is enough data available to satisfy the entire request, the read request completes successfully and returns the number of bytes read.
- If there is not enough data available to satisfy the entire request, the read request completes successfully, having read as much data as possible, and returns the number of bytes it was able to read.
- If there is no data available, the read request returns a value of -1 and sets the **errno** global variable to **EAGAIN**.

The availability of data depends on whether the input processing mode is canonical or noncanonical. The canonical or noncanonical modes can be set with the **stty** command.

Canonical mode input processing

In canonical mode input processing (**ICANON** flag set in **c_lflag** field of **termios** structure), terminal input is processed in units of lines. A line is delimited by a new-line (ASCII LF) character, an end-of-file (EOF) character, or an end-of-line (EOL) character. This means that a program attempting to read is blocked until an entire line has been typed or a signal has been received. Also, regardless of how many characters are specified in the read request, no more than one line is returned. It is not, however,

necessary to read an entire line at once. Any number of characters can be specified in a read request without losing information. During input, erase and kill processing is done.

ERASE character

(Backspace, by default) erases the last character typed

WERASE character

(Ctrl-W key sequence, by default) erases the last word typed in the current line, but not any preceding spaces or tabs

(A *word* is defined as a sequence of nonblank characters; tabs are regarded as blanks.) Neither the ERASE nor the WERASE character erases beyond the beginning of the line.

KILL character

(Ctrl-U sequence, by default) deletes the entire input line and, optionally, outputs a new-line character

All of these characters operate on a keystroke basis, independent of any backspacing or tabbing that might have been done.

REPRINT character

(Ctrl-R sequence, by default) prints a new line followed by the characters from the previous line that have not been read

Reprinting also occurs automatically if characters that would normally be erased from the screen are fouled by program output. The characters are reprinted as if they were being echoed. Consequently, if the **ECHO** flag is not set in the `c_lflag` field of the **termios** structure, the characters are not printed. The ERASE and KILL characters can be entered literally by preceding them with the escape character `\` (backslash), in which case, the escape character is not read. The ERASE, WERASE, and KILL characters can be changed.

Noncanonical mode input processing

In noncanonical mode input processing (**-ICANON** flag set in `c_lflag` field of **termios** structure), input bytes are not assembled into lines, and erase and kill processing does not occur.

MIN

Represents the minimum number of bytes that should be received when the read request is successful

TIME

A timer of 0.1-second granularity that is used to time-out burst and short-term data transmissions

The values of the MIN and TIME members of the `c_cc` array are used to determine how to process the bytes received. The MIN and TIME values can be set with the **stty** command. MIN and MAX have values from 0 to 255. The four possible combinations for MIN and TIME and their interactions are described in the subsequent paragraphs.

Case A: MINO, TIMEO

In this case, TIME serves as an interbyte timer, which is activated after the first byte is received and reset each time a byte is received. If MIN bytes are received before the interbyte timer expires, the read request is satisfied. If the timer expires before MIN bytes are received, the characters received to that point are returned to the user. If TIME expires, at least one byte is returned. (The timer would not have been enabled unless a byte was received.) The read operation blocks until the MIN and TIME mechanisms are activated by the receipt of the first byte or until a signal is received.

Case B: MINO, TIME = 0

In this case, only MIN is significant; the timer is not significant (the value of TIME is 0). A pending read request is not satisfied (blocks) until MIN bytes are received or until a signal is received. A program that uses this case to read record-based terminal I/O can block indefinitely in the read operation.

Case C: MIN = 0, TIMEO

In this case, because the value of MIN is 0, TIME no longer represents an interbyte timer. It now serves as a read timer that is activated as soon as the read request is processed. A read request is satisfied as soon as a byte is received or when the read timer expires. Note that if the timer expires, no bytes are returned.

If the timer does not expire, the read request can be satisfied only if a byte is received. In this case, the read operation does not block indefinitely, waiting for a byte. If, after the read request is initiated, no byte is received within the period specified by `TIME` multiplied by 0.1 seconds, the read request returns a value of 0, having read no data.

Case D: MIN = 0, TIME = 0

In this case, the minimum of either the number of bytes requested or the number of bytes currently available is returned without waiting for more bytes to be input. If no characters are available, the read request returns a value of 0, having read no data.

Cases A and B exist to handle burst-mode activity, such as file transfer programs, where a program needs to process at least the number of characters specified by the `MIN` variable at one time. In Case A, the interbyte timer is activated as a safety measure. In Case B, the timer is turned off.

Cases C and D exist to handle single-character, limited transfers. These cases are readily adaptable to screen-based applications that need to know if a character is present in the input queue before refreshing the screen. In Case C, the timer is activated. In Case D, the timer is turned off. Case D can lead to performance issues if overused; but it is better to use it than doing a read request with setting the `O_NONBLOCK` flag.

Writing data and output processing

When one or more characters are written, they are transmitted to the terminal as soon as previously written characters are displayed. (Input characters are echoed by putting them into the output queue as they arrive.) If a process produces characters more rapidly than they can be displayed, the process is suspended when its output queue exceeds a certain limit. When the queue has drained down to a certain threshold, the program is resumed.

Modem management

If the `CLOCAL` flag is set in the `c_cflag` field of the `termios` structure, a connection does not depend on the state of the modem status lines. If the `CLOCAL` flag is clear, the modem status lines are monitored. Under normal circumstances, an `open` function waits for the modem connection to complete. However, if the `O_NONBLOCK` or `CLOCAL` flag is set, the `open` function returns immediately without waiting for the connection.

If the `CLOCAL` flag is not set in the `c_cflag` field of the `termios` structure and a modem disconnect is detected by the terminal interface for a controlling terminal, the `SIGHUP` signal is sent to the controlling process associated with the terminal. Unless other arrangements have been made, this signal causes the process to terminate. If the `SIGHUP` signal is ignored or caught, any subsequent read request returns an end-of-file indication until the terminal is closed. Any subsequent write request to the terminal returns a value of -1 and sets the `errno` global variable to `EIO` until the device is closed.

Closing a terminal device file

The last process to close a terminal device file causes any output to be sent to the device and any input to be discarded. Then, if the `HUPCL` flag is set in the `c_cflag` field of the `termios` structure and the communications port supports a disconnect function, the terminal device performs a disconnect.

Converter modules

Converter modules are optional modules; they are pushed onto a tty stream only if required.

They are usually provided for internationalization purposes and perform various character mapping.

The following converter modules are shipped:

- The `nls` module
- The `uc_sjis` and `lc_sjis` modules

NLS module

The `nls` module is a lower converter module that can be pushed onto a tty stream below the line discipline. The `nls` module ensures terminal mapping; it executes the mapping of input and output

characters for nonstandard terminals (that is, for terminals that do not support the basic codeset ISO 8859 of the system).

The mapping rules are specified in two map files located in the `/usr/lib/nls/termmap` directory. The `.in` files contain the mapping rules for the keyboard inputs. The `.out` files contain the mapping rules for the display outputs. The files format is specified in the `setmaps` file format "setmaps File Format" in *Files Reference*.

SJIS Modules

The `uc_sjis` and `lc_sjis` modules are converter modules that can be pushed onto a `tty` stream. They ensure code set handling; they execute the conversion of multibyte characters between the shifted Japanese industrial standard (SJIS) format and the advanced Japanese EUC code (AJEC) format, supported by the line disciplines. These modules are needed when the user process and the hardware terminal uses the IBM-943 code set.

AJEC is a Japanese implementation of the extended UNIX code (EUC) encoding method, which allows combination of ASCII, phonetic Kana, and ideographic Kanji characters. AJEC is a superset of UNIX Japanese industrial standard (UJIS), a common Japanese implementation of EUC.

Japanese-encoded data consist of characters from up to four code sets:

Code set	Contained characters
ASCII	Roman letters, digits, punctuation and control characters
JIS X0201	Phonetic Kana
JIS X0208	Ideographic Kanji
JIS X0212	Supplemental Kanji.

AJEC makes use of all four code sets. SJIS makes use only of ASCII, JIS X0201, and JIS X0208 code sets. Therefore, the `uc_sjis` and `lc_sjis` modules convert:

- All SJIS characters into AJEC characters
- AJEC characters from ASCII, JIS X0201, and JIS X0208 code sets into SJIS characters
- AJEC characters from JIS X0212 code set into the SJIS undefined character

The `uc_sjis` and `lc_sjis` modules are always used together. The `uc_sjis` upper converter is pushed onto the `tty` stream above the line discipline; the `lc_sjis` lower converter is pushed onto the stream below the line discipline. The `uc_sjis` and `lc_sjis` modules are automatically pushed onto the `tty` stream by the `setmaps` command and the `setcsmap` subroutine. They are also controlled by the EUC ioctl operations described in the `eucioctl.h` file in the *Files Reference*.

TTY drivers

A `tty` driver is a STREAMS driver managing the actual connection to the hardware terminal.

Depending on the connection, three kinds of `tty` drivers are provided: asynchronous line drivers, the `pty` driver, and the LFT driver.

Asynchronous line drivers

The asynchronous line drivers are provided to support devices (usually ASCII terminals) directly connected to the system through asynchronous lines, including modems.

The asynchronous line drivers provide the interface to the line control hardware:

- The `cxma` driver supports the 128-port PCI adapter card.
- The `cxpa` driver supports the 8-port PCI adapter card.
- The `sf` driver supports the native ports on the system planar.
- The `sa` driver supports the 2-port PCI adapter card.

The asynchronous line drivers are responsible for setting parameters, such as the baud rate, the character size, and the parity checking. The user can control these parameters through the `c_cflag` field of the **termios** structure.

The asynchronous line drivers also provide the following features:

- The hardware and software flow control, or pacing discipline, specifies how the connection is managed to prevent a buffer overflow. The user can control this feature through the `c_iflag` field of the **termios** structure (software flow control) and the `x_hflag` field of the **termiox** structure (hardware flow control).
- The open discipline specifies how to establish a connection. This feature is controlled at configuration time through the `x_sflag` field of the **termiox** structure.

Pseudo-terminal driver

The pseudo-terminal (pty) driver is provided to support terminals that need special processing, such as X terminals or remote systems connected through a network.

A pty driver just transmits the input and output data from the application to a server process through a second stream. The server process, running in the user space, is usually a daemon, such as the **rlogind** daemon or the **xdm** daemon. It manages the actual communication with the terminal.

Other optional modules may be pushed on either user or server stream.

Loader domains

In some programming environments, it is desirable to have shared libraries loaded at the same virtual address in each process.

Due to the dynamic nature of shared libraries maintained by the AIX system loader, this condition cannot be guaranteed. Loader domains provide a means of loading shared libraries at the same virtual address in a set of processes.

The system loader loads shared libraries into multiple global shared library regions. One region is called the shared library text region, which contains the executable instructions for loaded shared libraries. The shared library text region is mapped to the same virtual address in every process. The other region is the shared library data region. This region contains the data for shared libraries. Because shared library data is read/write, each process has its own private region that is a copy of the global shared library region. This private region is mapped to the same virtual address in every process.

Since the global shared library regions are mapped at the same virtual address in every process, shared libraries are loaded at the same virtual address in most cases. The case where this is not true is when there is more than one version of a shared library loaded in the system. This happens whenever a shared library that is in use is modified, or any shared libraries it depends on are modified. When this happens, the loader must create a new version of the modified shared library and all other shared libraries that depend on the modified shared library. Note that all shared libraries ultimately depend on the *Kernel Name Space*. The *Kernel Name Space* contains all the system calls defined by the kernel and can be modified any time a kernel extension is dynamically loaded or unloaded. When the system loader creates a new version of a shared library, the new version must be located at a different location in the global shared library segments. Therefore, processes that use the new version have the shared libraries loaded at a different virtual address than processes that use the previous versions of the shared libraries.

A loader domain is a subset of all the shared libraries that are loaded in the system. The set of all shared libraries loaded in the system is called the *global loader domain*. This global loader domain can be subdivided into smaller user-defined loader domains. A user-defined loader domain contains one version of any particular shared library. Processes can specify a loader domain. If a process specifies a loader domain, the process uses the shared libraries contained in the loader domain. If more than one process specifies the same loader domain, they use the same set of shared libraries. Since a loader domain contains one version of any particular shared library, all processes that specify the same loader domain use the same version of shared libraries and have their shared libraries loaded at the same virtual address.

Using loader domains

If a process uses a loader domain, it must be specified at exec time. The loader domain specified is in effect and used for the entire duration of the process. When a process that specifies a loader domain calls the **exec** system call, the system loader takes the following actions:

Finds/creates loader domain

The access permissions associated with the loader domain are checked to determine if this process can use the loader domain. If the process does not have sufficient privilege to access (read or write) the loader domain, no domain is used by the process. If the process does have sufficient privilege, the list of loader domains maintained by the system loader is searched for the loader domain specified by the process. If the loader domain specified is not found, it is created if the process has sufficient privilege. If the process does not have sufficient privilege to create the loader domain, then the **exec** call fails, and an error is returned.

Uses loader domain to limit search

If the process needs any shared libraries that are already listed in the loader domain, the version of the library specified in the domain is used. The version of the shared library in the loader domain is used regardless of other versions of the shared library that may exist in the global loader domain.

Adds shared libraries to loader domain

If the process needs a library that is not in the loader domain, the loader loads the library into the process image by following the normal loader convention of loading the most recent version. If the process has sufficient privilege, this version of the library is also added to the loader domain. If the process does not have sufficient privilege to add an entry, the **exec** call fails, and an error is returned.

Shared libraries can also be explicitly loaded with the **load()** system call. When a shared library is explicitly loaded, the data for these modules is normally put at the current break value of the process for a 32-bit process. For a 64-bit process, the data for the modules is put in the region's privately loaded modules. If a process uses a loader domain, the system loader puts the data in the shared library data region. The virtual address of this explicitly loaded module is the same for all processes that load the module. If the process has sufficient privilege, the shared library is added to the loader domain. If the process does not have sufficient privilege to add an entry, the **load** call fails, and an error is returned.

A loader domain can be associated with any regular file. It is important to note that a loader domain is associated with the file, NOT the path name of the file. The mode (access permissions) of the file determines the operations that can be performed on the loader domain. Access permissions on the file associated with the loader domain and the operations allowed on the loader domain are as follows:

- If the process is able to read the file, the process can specify the loader domain to limit the set of shared libraries it uses.
- If the process is able to write to the file, the process is able to add shared libraries to the loader domain and create the loader domain associated with the file.

If a process attempts to create or add entries to a loader domain without sufficient privilege, the operation in progress (**exec** or **load**) fails, and an error is returned.

Loader domains are specified as part of the **LIBPATH** information. **LIBPATH** information is a colon (:) separated list of directory path names used to locate shared libraries. **LIBPATH** information can come from either the **LIBPATH** environment variable or the **LIBPATH** string specified in the loader section of the executable file. If the first path name in the **LIBPATH** information is a regular file, a loader domain associated with the file is specified. For example:

- If `/etc/loader_domain/00domain_1` is a regular file, then setting the **LIBPATH** environment variable to the string

```
/etc/loader_domain/00domain_1:/lib:/usr/lib
```

causes processes to create and use the loader domain associated with the `/etc/loader_domain/00domain_1` file.

- If `/etc/loader_domain/00domain_1` is a regular file, then the `ldom` program is built with the following command:


```
cc -o ldom ldom.c -L/etc/loader_domain/00domain_1
```

The path name `/etc/loader_domain/00domain_1` is inserted as the first entry in the **LIBPATH** information of the loader section for the `ldom` file. When `ldom` is executed, it creates and uses the loader domain associated with the `/etc/loader_domain/00domain_1` file.

Creating/Deleting loader domains

A loader domain is created the first time a process with sufficient privilege attempts to use the domain. Access to a loader domain is controlled by access to the regular file associated with the domain. Application writers are responsible for managing the regular files associated with loader domains used by their applications. Loader domains are associated with regular files NOT the path names of the files. The following examples illustrate this point:

- The `ap1` application has specified loader domain `domain01` in its **LIBPATH** information. The `ap1` application is then executed. The current working directory is `/home/user1`, and it contains a regular file `domain1` that is writable by `ap1`. A new loader domain associated with the file `/home/user1/domain01` is created. `ap1` is executed again. This time `/home/user2` is the current working directory, and it also contains a regular file `domain01` that is writable by `ap1`. A new loader domain associated with the file `/home/user1/domain02` is created.
- Application `ap1` has specified loader domain `/etc/1_domain/domain01` in its **LIBPATH** information. `ap1` is then executed. `/etc/1_domain/domain01` is a regular file that is writable by `ap1`. A new loader domain associated with the file `/etc/1_domain/domain01` is created.

`/home/user1/my_domain` is a symbolic link to file `/etc/1_domain/domain01`.

Application `ap2` has specified loader domain `/home/user1/my_domain` in its **LIBPATH** information. `ap2` is then executed. The system loader notices that `/home/user1/my_domain` refers to the same file as `/etc/1_domain/domain01`. A loader domain is already associated with file `/etc/1_domain/domain01`; therefore, this loader domain is used by application `ap2`.

- Application `ap1` has specified loader domain `/etc/1_domain/domain01` in its **LIBPATH** information. `ap1` is then executed. `/etc/1_domain/domain01` is a regular file that is writable by `ap1`. A new loader domain associated with the file `/etc/1_domain/domain01` is created.

File `/etc/1_domain/domain01` is deleted and recreated as a regular file.

Application `ap1` is executed again. There is no longer any way to access the regular file that is associated with the original loader domain `/etc/1_domain/domain01`. Therefore, a new loader domain associated with the file `/etc/1_domain/domain01` is created.

Loader domains are dynamic structures. During the life of a loader domain, shared libraries are added and deleted. A shared library is added to a loader domain when a process that specified the loader domain needs a shared library that does not already exist in the domain. Of course, this assumes the process has sufficient privilege to add the shared library to the loader domain.

A separate use count is kept for each shared library that is a member of a loader domain. This use count keeps track of how many processes with loader domains are using the shared library. When this use count drops to zero, the shared library is deleted from the loader domain.

Data management application programming interface

AIX provides a data management application programming interface (DMAPI), which is an implementation of the "System Management: Data Storage Management (XDSM) API" X/Open standard that is published by The Open Group.

The DMAPI allows software vendors to develop data management applications using a set of functions and semantics not found in POSIX-compliant systems. It does not provide direct functionality to the end user. Complete documentation of the DMAPI is found in the Publications section of The Open Group's Web site.

The DMAPI provided by AIX is a general purpose implementation. The level of support for optional interfaces and functionality is determined by the underlying file system implementation and is documented in a separate section for the specific file system.

The intent of the DMAPI is to support a single product on any single file system. The DMAPI does not preclude different products from different vendors operating on the same file system, but it is not recommended. Different products on different file systems are fully supported by the DMAPI with regard to event delivery, subject to the following restrictions:

- Multiple sessions cannot register disposition for the same event on the same object.
- Event messages are targeted to and enqueued on sessions; there is no explicit targeting of an event to a specific process.
- If no session has registered a disposition for a particular event other than the mount event, the DMAPI will not generate an event and allow the process to proceed as if there is no event enabled. If no session has registered a disposition for the mount event, which is always enabled, the DMAPI will fail the event and not allow the file system to be mounted.

The DMAPI is implemented in an abstract layer within AIX, allowing any underlying file system to define its individual level of support and implementation options. The journaled file system (JFS) file system does not provide any support for the DMAPI. The enhanced journaled file system (JFS2) behaviors for implementation options, limits, and other specifics described by the X/Open standard are outlined in [DMAPI Considerations for the Enhanced Journaled File System](#).

The `dm_init_service` function returns 0 when the AIX DMAPI is correctly initialized and -1 if the initialization fails. Use of any other DMAPI function after the initialization fails will also fail.

The AIX DMAPI does not provide support for the following optional DMAPI functions:

- `dm_downgrade_right`
- `dm_upgrade_right`
- `dm_obj_ref_*` family
- `dm_pending`

Other optional interfaces might not be supported by the underlying file system implementation and are indicated in the DMAPI documentation for the specific file system.

When a data management (DM) application specifies that it wants to block until a right becomes available, the DM application is blocked uninterruptibly.

AIX allows multiple, non-overlapping persistent managed regions. Only regular files are allowed to have managed regions. Whether or not managed regions are reordered or coalesced is determined by the underlying file system implementation.

When no session has registered to receive a particular event for which an object is enabled and activity occurs that would otherwise trigger the event, AIX does not generate the event and allows the process to proceed as if there is no event enabled.

Executing the `dm_set_eventlist` function causes a persistent event list to be stored with the object. If an event list was previously set for the entire file system and a subsequent event list for an object in that file system includes an event that was set for the file system, events will continue to be generated based on the event list for the file system until such time as that event is disabled, in which case the event list for the object will come into play.

When a process generating an event is blocked waiting on a response from a DM application, the sleep is interruptible.

AIX adopts a reasonably reliable model of asynchronous message delivery. The number of undelivered asynchronous messages is limited by the amount of available memory (real or virtual) configured on the system. If the number of messages exceeds the amount of available memory, undelivered asynchronous messages will be lost. Asynchronous delivery of namespace event messages is determined by the underlying file system implementation.

For AIX, `DM_SESSION_INFO_LEN` is 256, and `DM_ATTR_NAME_SIZE` is 8.

For DMAPI interfaces that return data to a user buffer and fill in a user variable with the resulting size of the buffer, both the contents of the buffer and the user size variable are undefined when the interface fails with an error other than E2BIG. For any such error, the contents of the user buffer must be ignored. When the interface fails and errno is E2BIG, the content of the user size variable will be set to indicate the required size, in which case the application can retry the interface with a resized buffer.

DMAPI considerations for the enhanced journaled file system

Note: Internal snapshots cannot be used with DMAPI-managed file systems.

In addition to the functionality provided by the general AIX implementation of the DMAPI, the JFS2 implementation provides the following functionality and restrictions.

The `dm_get_config` function returns the following values for the JFS2 implementation options and limits:

DM_CONFIG_BULKALL

Supported

DM_CONFIG_LEGACY

Supported

DM_CONFIG_PERS_ATTRIBUTES

Supported

DM_CONFIG_PERS_EVENTS

Supported

DM_CONFIG_PERS_INHERIT_ATTRIBS

Supported

DM_CONFIG_PERS_MANAGED_REGIONS

Supported

DM_CONFIG_PUNCH_HOLE

Supported

DM_CONFIG_WILL_RETRY

Supported

DM_CONFIG_CREATE_BY_HANDLE

Not supported

DM_CONFIG_LOCK_UPGRADE

Not supported

DM_CONFIG_OBJ_REF

Not supported

DM_CONFIG_PENDING

Not supported

DM_CONFIG_DTIME_OVERLOAD

TRUE

DM_CONFIG_MAX_ATTR_ON_DESTROY

128

DM_CONFIG_MAX_ATTRIBUTE_SIZE

4072

DM_CONFIG_MAX_HANDLE_SIZE

32

DM_CONFIG_MAX_MANAGED_REGIONS

167

DM_CONFIG_MAX_MESSAGE_DATA

65536

DM_CONFIG_TOTAL_ATTRIBUTE_SPACE

4072

In the JFS2 implementation, all DM attribute values share the same allocation. Consequently, the size of any one attribute's value cannot exceed `DM_CONFIG_MAX_ATTRIBUTE_SIZE` and is further restricted by the sum of the value sizes of all DM attributes associated with an object, which is also limited to `DM_CONFIG_MAX_ATTRIBUTE_SIZE`.

In addition to the optional interfaces not supported by AIX, the JFS2 implementation does not support the optional DMAPI cancel and debut events, nor the additional optional `dm_getall_dmattr`, `dm_create_by_handle`, and `dm_symlink_by_handle` functions.

Due to the current implementation of JFS2's extended attribute support, the `dm_set_region` function causes the file's `ctime` to be modified. JFS2 does not attempt to reorder nor coalesce managed regions.

JFS2 generates asynchronous namespace event messages for all corresponding operations, whether they succeed or fail.

JFS2 provides interfaces that allow pre-allocation and direct control of metadata within a file system. Use of these interfaces with either the `MM_ALLOC` or `MM_RECORD` modes generates a DMAPI write event for the specified offset and length.

If a value is not specified for the mask to the `dm_get_bulkall`, `dm_get_bulkattr`, `dm_get_dirattr`, and `dm_get_fileattr` functions (that is, it is set to zero), JFS2 will return all fields in the `dm_stat` structure. If the mask is set to a specific value, only the fields requested by the mask are returned; the values for the fields not specified by the mask are undefined.

JFS2 does not use the `respbufp` parameter of the `dm_respond_event` function. If specified, the content of the buffer is undefined when the functions returns.

Because JFS2 overloads `dm_ctime` and `dm_dtime` (that is, `DM_CONFIG_DTIME_OVERLOAD` is true), the `setdtime` parameter of the `dm_set_dmattr` function is ignored.

At the time a file is memory mapped (that is, when the `mmap(2)` call is executed), any non-resident portions of a file must be made resident by the DM application. To notify the application of the mapping, JFS2 will generate a read or write event corresponding to the mode and region being mapped.

Activating the DMAPI on a JFS2 file system

To activate the DMAPI on a JFS2 file system, type the following:

```
chfs -a managed=yes mountpoint
```

If the file system is currently mounted when the `chfs` command is issued, there must be a DMAPI-enabled application listening for and responding to mount events when the *managed* parameter is set; the success of the `chfs` command will depend on how the application responds to the mount event.

To deactivate the DMAPI on a JFS2 file system, type the following:

```
chfs -a managed=no mountpoint
```

If the file system is currently mounted when the `chfs` command is issued, there must be a DMAPI-enabled application listening for and responding to pre-unmount events when the *managed* parameter is set. The success of the `chfs` command depends on how the application responds to the pre-unmount event.

Using the DMAPI on JFS2 encrypted file systems

When you perform invisible I/O operations on encrypted files in a JFS2 encrypted file system, the same offset and length alignment restrictions apply as if you performed RAW mode I/O on the file. In particular, the offset and length of the I/O must be block-aligned according to the block size of the file system. The size of the encrypted data is always a multiple of file system blocks, even when the decrypted file size is not; when the file size is not block aligned, the file contains encrypted data beyond the file size.

Note: The `stat` subroutine and DMAPI interfaces, such as the `dm_get_fileattr` function, report the clear text (decrypted) file size, while the `statx` subroutine reports the block-aligned encrypted data size when you pass `STX_EFSRAW` as the *command* parameter.

The **dm_read_invis** function and the **dm_write_invis** function must meet the following requirements to make the operation successful:

dm_read_invis

For encrypted files, both the *off* and the *len* parameters must be file system block-size aligned, or the operation fails with the **EINVAL** error code.

dm_write_invis

For encrypted files, both the *off* and the *len* parameters must be file system block-size aligned and the operation must not attempt to extend the file, or the operation fails with the **EINVAL** error code.

Using the DMAPI on AIX workload partitions

You must add the **PV_FS_DMAPI** privilege to the set of privileges and assign them to the processes running in the Workload Partition (WPAR) to run the DMAPI applications within a WPAR. You can add and assign the set of privileges to a WPAR when it is being created, or you can modify the set of privileges later.

Examples

```
mkwpar -S privs+=PV_FS_DMAPI -n wparname
```

```
chwpar -S privs+=PV_FS_DMAPI wparname
```

By default, only the root processes obtain the privilege to run in a WPAR. In a root-disabled system or in an Trusted AIX installation, where root is disabled by default, nonroot processes obtain this privilege by using the **privcmds** table in a global or WPAR system. For more information, see [RBAC privileges](#).

AIX transactional memory programming

Transactional memory (TM) is a shared-memory synchronization construction that allows process-threads to perform storage operations that appear to be atomic to other process-threads or applications.

Overview

TM is a construct that allows execution of lock-based critical sections of the code without acquiring a lock. The IBM POWER 8 processor is the first processor that implements TM programming.

Use the TM facility in some of the following scenarios:

- Optimistic Execution of Lock-Based Applications – TM supports speculative execution of critical sections of code without acquiring a lock. This method provides the benefits of fine-grained locking to applications by using the current locks that are not tuned for performance.
- Transactional Programming in High-Level Languages – The transactional programming model is a growing industry-wide standard that offers productivity gains relative to lock-based shared memory programs.
- Checkpoint/Rollback Usage – TM is used as a checkpoint to restore architectural state. This method enables speculative compiler optimizations during runtime code optimization or generation and simulation of checkpoints.

To use the TM facility, a process-thread marks the beginning and end of the sequence of storage accesses or transaction with the `tbegin.` and `tend.` instructions. The `tbegin.` instruction initiates the transactional execution, during which the loads and stores appear to occur atomically. The `tend.` instruction ends the transactional execution.

If a transaction is prematurely stopped, the storage updates that were made after executing the `tbegin.` instruction are rolled back. Correspondingly, the contents of a subset of the registers are also rolled back to the state before the `tbegin.` instruction was executed. When a transaction is prematurely stopped, a software failure handler is started. The failure can be of the transient type or the persistent type. The failure handler can retry the transaction or choose to employ a different locking construct or logic path that depends on the nature of the failure.

The AIX operating system supports the usage of TM including handling of TM state management across context switches and interrupts.

Checkpoint state

When a transaction is initiated a set of registers are saved that represent the checkpoint state of the processor. In case of a transaction failure, a set of registers are restored to the point before the start of the transaction. The checkpoint state of the processor is also called as the pre-transactional state. The checkpoint state includes the problem-state writable registers except for the CR0, FXCC, EBBHR, EBBRR, BESCR registers, the performance monitor registers, and the TM SPRs.

Note: The checkpoint state cannot be directly accessed through the supervisor state or the problem state.

The checkpoint state is copied into the respective registers after the new `trereclaim.` instruction is executed. This process allows privileged code to save or modify the values. The checkpoint state is copied back into the speculative registers from the respective user-accessible registers after the execution of the new `trerechkpt.` instruction.

The following TM SPRs are added to the machine state for the processor:

Name	Title	Description	Privileged mtspr	Privileged mfspr	Size (bits)	SPR
FSCR	Facility Status and Control Register	Controls the available facilities in problem state and indicates the cause of a Facility unavailable interrupt.	yes	yes	64	153
TEXASR	Transaction Exception And Summary Register	Contains the transaction level and summary information that is used by the transaction failure handlers. The 0:31 bits contain the cause of the failure.	no	no	64	130
TFHAR	Transaction Failure Handler Address Register	Records the EA of the software failure handler. The TFHAR register is always set to the NIA for the <code>tbegin.</code> instruction that initiated the transaction.	no	no	64	128
TFIAR	Transaction Failure Instruction Address Register	Set to the exact EA of the instruction that causes the failure, when possible. The accuracy of the TFIAR register is recorded in the Exact field (bit 37) of the TEXASR register.	no	no	64	129
TEXASRU	Transaction Exception and Summary Register (Upper Half)	High-order half of TEXASR register.	no	no	32	131

The new TEXASR register contains information related to the state of a transaction and the cause of a transaction failure. The following table describes the fields included in the TEXASR register:

Field	Value-Meaning	Bits
Failure Code (Note: bit 7 is referred to as the Failure persistent field)	Transaction Failure Codes	0:7
Disallowed	0b1 - The access type instruction is not allowed	8
Nesting Overflow	0b1 - The maximum transaction level was exceeded.	9
Footprint Overflow	0b1 - The tracking limit for transactional storage accesses was exceeded.	10
Self-Induced Conflict	0b1 - A self-induced conflict occurred in suspended state.	11
Non-Transactional Conflict	0b1 - A conflict occurred with a non-transactional access by another processor.	12

Field	Value-Meaning	Bits
Transaction Conflict	Ob1 - A conflict occurred with another transaction.	13
Translation Invalidation Conflict	Ob1 - A conflict occurred with a TLB invalidation.	14
Implementation Specific	Ob1 - An implementation-specific condition caused the transaction to fail.	15
Instruction Fetch Conflict	Ob1 - An instruction fetch by the thread or another thread that was performed from a block that was previously written transactionally.	16
Reserved for future failure cases		17:30
Abort	Ob1 - An abort was caused by the execution of a particular TM instruction.	31
Suspended	Ob1 - The failure was recorded in Suspended State.	32
Reserved		33
Privilege	The thread was in the privilege state ([MSRHV PR]) at the time of failure recording.	34:35
Failure Summary (FS)	Ob1 - A failure was detected and recorded.	36
TFIAR Exact	Ob0 - The value in the TFIAR field is an approximate value. Ob1 - The value in the TFIAR field is an exact value.	37
ROT	Set to 0b0 when a non-ROT tbegin. instruction is executed. Set to 0b1 when a ROT is initiated.	38
Reserved		39:51
Transaction Level (TL)	Transaction level (nesting depth + 1) for the active transaction has the following values: <ul style="list-style-type: none"> • 0 if the most recent transaction completed successfully. • The transaction level at which the most recent transaction failed, if the transaction did not complete successfully. Note: A value of 1 corresponds to an <i>outer transaction</i> . A value greater than 1 corresponds to a <i>nested transaction</i> .	52:63

Notes:

- Exactly 1 bit of the 8-31 bits of the TEXASR register is set when the transaction failure is recorded. The single bit that is set indicates that the particular instruction or event caused failure.
- A Rollback Only Transaction (ROT) is a sequence of instructions that is executed either as a unit or the instructions are not executed. This construct allows for the speculative execution of a bulk of instructions with minimal cost. A ROT does not have the full atomic nature as a normal transaction or its synchronization and serialization properties. Therefore, ROTs must not be used to manipulate shared data.

Software failure handler

When a transaction fails, the machine hardware redirects control to the failure handler that is associated with the outermost transaction. When a transaction fails, the control is redirected to the instruction that follows the tbegin. instruction, CRO is set to either

```
0b101 || 0
```

or

```
0b010 || 0
```

Therefore, the instruction after the `tbegin.` instruction must be a branch instruction predicated on bit 2 of CRO. For example, after the `tbegin.` instruction is executed, the `beq` branch instruction predicated on bit 2 of CRO. The target of the branch must be a section of code that handles transaction failures. When the `tbegin.` instruction is successfully run at the start of the transaction, CRO is set to either

```
0b000 || 0 or 0b010 || 0
```

Note: The bits 0:31 of TEXASR reports the cause of the failure. The failure code (FC) field in bits 0-7, is used for the following scenarios:

- Privileged supervisor or hypervisor code causes the failure by using the `treclaim.` instruction.
- Problem-state code causes the failure by using a form of the `tabort.` instruction.

A value of 1 in bit 7 of TEXASR indicates that the failure is persistent and the transaction is bound to fail when the transaction is attempted again. The failure codes reserved by the AIX operating system indicate the cause of the failure that are defined in the `/usr/include/sys/machine.h.` directory.

A sample transaction

The following assembler code example shows a simple transaction that writes the value in GPR 5 into the address in GPR 4, which is assumed to be shared among multiple threads of execution. If the transaction fails due to a persistent cause, the code falls back to another code path at the `lock_based_update` label. The code for the alternate path is not shown.

```
trans_entry:
    tbegin                # Start transaction
    beq                   failure_hdlr    # Handle transaction failure
    stw                   r5, 0(r4)      # Write to memory pointed to by r4.
    tend.                # End transaction
    b                     trans_exit
failure_hdlr:
    mfspr                 r4, TEXASRU    # Read high-order half of TEXASR
    andis.                r5, r4, 0x0100 # Is the failure persistent?
    bne                   lock_based_update # If persistent, acquire lock and
                                           # then perform the write.
    b                     trans_entry    # If transient, try again.

lock_based_update:
trans_exit:
```

Runtime determination of Transactional Memory capability

A program can determine whether a system supports the TM category of the POWER ISA by reading the `SC_TM_VER` system variable using the `getsystemcfg` subroutine. A `__power_tm()` macro is provided in the `/usr/include/sys/systemcfg.h` file to determine the TM capability within a program. This macro is useful for software that conditionally uses the TM capability when it is present, or uses the functionally equivalent to lock-based code paths when the TM capability is not present.

Extended context structure

The earlier versions of the AIX operating system introduced support for extended context structures to support the vector state and user keys. The existing extended context structure support is further extended to support machine state that is required by TM.

An extended context is allocated and pinned for each transactional process-thread when it first uses TM. If the extended context area cannot be allocated and pinned, then the process receives a `SIGSEGV` signal that results in termination of the process.

The machine-context information is included in the `sigcontext` structure that is provided to signal handlers. When a signal handler returns, the machine context present in the `sigcontext` structure is activated. The `sigcontext` structure is actually a subset of the larger `ucontext` structure. The two

structures are identical up to `sizeof(struct sigcontext)`. When the AIX operating system builds a signal context to be passed to a signal handler, a `ucontext` structure is built on the stack of the signal handler. The machine-context portion of a signal context must contain all of the active machine state, including the volatile and nonvolatile state for the involuntarily interrupted context. The `ucontext` structure contains an indicator to determine whether extended context information is available.

The `__extctx` field in the `ucontext` structure is the address of an extended context structure is defined in the `/usr/include/sys/context.h` file. The `__extctx_magic` field in the `ucontext` structure indicates whether the extended context information is valid when the value of `__extctx_magic` field is equal to `__EXTCTX_MAGIC`. The additional machine state for a thread that uses the TM capability is restored and saved as a member of the context extension, to the `ucontext` structure as a part of the signal delivery and return.

If an application chooses to explicitly enable the use of Transactional Memory, it takes an extended size `ucontext` structure that already has space for the `__extctx` field that is included by the implicit definition of `__EXTABI__` by the compiler. The extended `ucontext` structure can also be picked up by an explicit definition of `__AIXEXTABI`.

The `getcontext()`, `setcontext()`, `makecontext()`, and `swapcontext()` subroutines of `libc` are not supported while in transactional or suspended state. When the subroutines are called within transaction, the `getcontext()`, `setcontext()`, `makecontext()` subroutines result in a persistent transaction failure of `TM_LIBC` type, which is defined in the `/usr/include/sys/machine.h` file.

When a `swapcontext()` subroutine is called within a transaction, it results in the following behavior:

- When a `swapcontext()` subroutine is in transactional state, it results in a persistent transaction failure of `TM_LIBC` type.
- When a `swapcontext()` subroutine is in suspended state, it results in the transaction being doomed, the specified `ucontext` structure swapped in and execution of the program is resumed by the specified `ucontext` structure. The resulting state and subsequent behavior after the `swapcontext()` subroutine returns are undefined.

If the `getcontext()`, `setcontext()`, and `swapcontext()` subroutines are called in a non-transactional state, the subroutines do not retrieve or restore any extended TM context into or from the `ucontext` structure pointed to by the `ucp` or `oucp` parameters. No error is indicated when the `setcontext()` or `swapcontext()` subroutines are called with the extended TM context present.

See the `/usr/include/sys/context.h` header file for detailed information of the extended context.

Signal delivery

Asynchronous signals that are received by an application while in a transaction is delivered non-transactionally. When in transactional state, the delivery of synchronous signals is not allowed and instead results in a persistent transaction failure of `TM_SYNC_SIGNAL` type, as defined in the `/usr/include/sys/machine.h` file.

Alignment interrupts and program interrupts

In transactional state, the alignment interrupts and program interrupts are caused due to an illegal operation or an operation that requires emulation result in a persistent transaction failure of `TM_ALIGN_INT` type or `TM_INV_OP` type, as defined in the `/usr/include/sys/machine.h` file. When in suspended state, alignment and program interrupts are processed normally by using the non-speculative semantics.

System calls

It is suggested that the system calls are not invoked within a transaction. System calls are only supported within a transaction when the transaction is suspended through the `tsuspend` instruction.

When a system call is invoked while a processor or thread is transactional and the transaction is not suspended, the system call is not invoked by the AIX kernel and the associated transaction fails

persistently. When this error occurs, the FC field of the TEXASR register contains the TM_ILL_SC failure code, which is defined in the `/usr/sys/include/machine.h` file.

It is assumed that any operations performed under a suspended transaction when the application programmer has explicitly suspended the transaction are intended to be persistent. Any operations that are performed by a system call that is invoked while in suspended state is not rolled-back even if the transaction fails.

The AIX operating system does not support system calls to be made while in transactional state because there is no way to roll back any operations, including I/O, performed by AIX underneath a system call.

setjmp() and longjmp() subroutines

The `setjmp()` and `longjmp()` subroutines of `libc` are not supported in transactional or suspended state because of the effects of setting a jump buffer and jumping back to the buffer. Consider the following scenarios

1. If the `setjmp()` subroutine is called inside of a transaction and the corresponding `longjmp()` subroutine is called after the transaction ends, the jump is to a speculative state that is now invalid.
2. If the `setjmp()` subroutine is called before the transaction, a corresponding `longjmp()` subroutine goes to the state before the transaction started regardless of whether the transaction has ended, failed, or aborted.
3. If the `setjmp()` subroutine is called within a transaction and then the transaction is aborted, the updates that are made to the jump buffer by the `setjmp()` subroutine will not appear to have occurred.

When the `setjmp()` subroutine is called within a transaction, it results in a persistent transaction failure either of `TM_LIBC` type or `TM_ILL_SC` type that is defined in the `/usr/include/sys/machine.h` file.

When a `longjmp()` subroutine is called within a transaction, it results in the following behavior:

- When a `longjmp()` subroutine is in a transactional state, it results in a persistent transaction failure of `TM_LIBC` type or `TM_ILL_SC` type that is defined in the `/usr/include/sys/machine.h` file.
- When a `longjmp()` subroutine is in a suspended state, it results in the transaction being doomed, the specified jump buffer is restored, and execution of the program returns to the corresponding `setjmp()` subroutine. The resulting state and subsequent behavior are undefined after the `setjmp()` subroutine returns from the `longjmp()` subroutine.

Compilers

An AIX operating system compiler that supports the Transactional Memory must conform to the AIX ABI. When the TM is enabled, a C or C++ compiler must predefine `__EXTABI__`. Refer to the compiler documentation for detailed information.

Assembler

The AIX operating system assembler, in the `/usr/ccs/bin/as` supports the additional instruction set defined for the TM in the POWER ISA, and as implemented by the POWER8 processor. You can use the `-m pwr8` assembly mode or the `.machine pwr8` pseudo op within the source file to enable the assembly of TM instructions. For more information, refer to the assembly language reference.

Debugger

The `/usr/ccs/bin/dbx` debugger supports machine-level debugging of TM programs. This support includes the ability to disassemble the new Transactional Memory instructions and to view the TM SPRs: TEXASR, TEXASRU, TFIAR, and TFHAR registers.

Setting a breakpoint inside of a transaction causes the transaction to unconditionally fail whenever the breakpoint is encountered. Therefore, the suggested approach is to debug a transaction that is failing to

set a breakpoint on the transaction's failure handler and then view the TEXASR and TFIAR registers when the breakpoint is encountered to determine the cause and location of the failure.

In dbx, the TEXASR, TFIAR, and TFHAR registers can be viewed by using the print subcommand with the `$texasr`, `$tfiar`, or `$tfhar` parameter. The line of code that is associated with the address found in the TFIAR, and TFHAR registers can be viewed through the list subcommand, for example:

```
(dbx) list at $tfiar
```

The dbx `tm_status` subcommand is used to view and interpret the contents of the TEXASR register. This subcommand is used to determine the nature of a transaction failure.

Enablement for third-party debuggers is provided in the form of a new `PTT_READ_TM` `ptrace` operation for reading the TM state of a thread. Refer to the `ptrace` documentation for details.

Tracing support

The AIX trace facility is expanded to include a set of trace events for TM operations that are performed by the AIX operating system including pre-emption causing transaction failure and other various operations that can cause transaction failure. The trace event identifier 675 can be used as input to the `trace` and `trcrpt` commands to view TM-related trace events.

Core files

The AIX operating system also supports the inclusion of TM machine state as part of the core file for processes or threads that uses TM. If a process or thread is using or used TM, the TM machine state is included in the core image for that thread.

Note: The TM state is only supported in the current core file formats for the AIX operating system. You can use the `dbx` command to read and view the TM machine state of a TM-enabled core file.

AIX threads library

The use of Transactional Memory is not supported for applications that use M:N threads. Undefined behavior may occur in transactional threads in an environment where more than one thread shares a single kernel thread. Usage of Transactional Memory by an application that uses M:N threads may lead to a persistent transaction failure with the Failure Code of `TM_PTH_PREEMPTED` being set in the TEXASR register.

Notices

This information was developed for products and services offered in the US.

IBM may not offer the products, services, or features discussed in this document in other countries. Consult your local IBM representative for information on the products and services currently available in your area. Any reference to an IBM product, program, or service is not intended to state or imply that only that IBM product, program, or service may be used. Any functionally equivalent product, program, or service that does not infringe any IBM intellectual property right may be used instead. However, it is the user's responsibility to evaluate and verify the operation of any non-IBM product, program, or service.

IBM may have patents or pending patent applications covering subject matter described in this document. The furnishing of this document does not grant you any license to these patents. You can send license inquiries, in writing, to:

*IBM Director of Licensing
IBM Corporation
North Castle Drive, MD-NC119
Armonk, NY 10504-1785
US*

For license inquiries regarding double-byte character set (DBCS) information, contact the IBM Intellectual Property Department in your country or send inquiries, in writing, to:

*Intellectual Property Licensing
Legal and Intellectual Property Law
IBM Japan Ltd.
19-21, Nihonbashi-Hakozakicho, Chuo-ku
Tokyo 103-8510, Japan*

INTERNATIONAL BUSINESS MACHINES CORPORATION PROVIDES THIS PUBLICATION "AS IS" WITHOUT WARRANTY OF ANY KIND, EITHER EXPRESS OR IMPLIED, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF NON-INFRINGEMENT, MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE. Some jurisdictions do not allow disclaimer of express or implied warranties in certain transactions, therefore, this statement may not apply to you.

This information could include technical inaccuracies or typographical errors. Changes are periodically made to the information herein; these changes will be incorporated in new editions of the publication. IBM may make improvements and/or changes in the product(s) and/or the program(s) described in this publication at any time without notice.

Any references in this information to non-IBM websites are provided for convenience only and do not in any manner serve as an endorsement of those websites. The materials at those websites are not part of the materials for this IBM product and use of those websites is at your own risk.

IBM may use or distribute any of the information you provide in any way it believes appropriate without incurring any obligation to you.

Licensees of this program who wish to have information about it for the purpose of enabling: (i) the exchange of information between independently created programs and other programs (including this one) and (ii) the mutual use of the information which has been exchanged, should contact:

*IBM Director of Licensing
IBM Corporation
North Castle Drive, MD-NC119
Armonk, NY 10504-1785
US*

Such information may be available, subject to appropriate terms and conditions, including in some cases, payment of a fee.

The licensed program described in this document and all licensed material available for it are provided by IBM under terms of the IBM Customer Agreement, IBM International Program License Agreement or any equivalent agreement between us.

The performance data and client examples cited are presented for illustrative purposes only. Actual performance results may vary depending on specific configurations and operating conditions.

Information concerning non-IBM products was obtained from the suppliers of those products, their published announcements or other publicly available sources. IBM has not tested those products and cannot confirm the accuracy of performance, compatibility or any other claims related to non-IBM products. Questions on the capabilities of non-IBM products should be addressed to the suppliers of those products.

Statements regarding IBM's future direction or intent are subject to change or withdrawal without notice, and represent goals and objectives only.

All IBM prices shown are IBM's suggested retail prices, are current and are subject to change without notice. Dealer prices may vary.

This information is for planning purposes only. The information herein is subject to change before the products described become available.

This information contains examples of data and reports used in daily business operations. To illustrate them as completely as possible, the examples include the names of individuals, companies, brands, and products. All of these names are fictitious and any similarity to actual people or business enterprises is entirely coincidental.

COPYRIGHT LICENSE:

This information contains sample application programs in source language, which illustrate programming techniques on various operating platforms. You may copy, modify, and distribute these sample programs in any form without payment to IBM, for the purposes of developing, using, marketing or distributing application programs conforming to the application programming interface for the operating platform for which the sample programs are written. These examples have not been thoroughly tested under all conditions. IBM, therefore, cannot guarantee or imply reliability, serviceability, or function of these programs. The sample programs are provided "AS IS", without warranty of any kind. IBM shall not be liable for any damages arising out of your use of the sample programs.

Each copy or any portion of these sample programs or any derivative work must include a copyright notice as follows:

© (your company name) (year).

Portions of this code are derived from IBM Corp. Sample Programs.

© Copyright IBM Corp. _enter the year or years_.

Privacy policy considerations

IBM Software products, including software as a service solutions, (“Software Offerings”) may use cookies or other technologies to collect product usage information, to help improve the end user experience, to tailor interactions with the end user or for other purposes. In many cases no personally identifiable information is collected by the Software Offerings. Some of our Software Offerings can help enable you to collect personally identifiable information. If this Software Offering uses cookies to collect personally identifiable information, specific information about this offering’s use of cookies is set forth below.

This Software Offering does not use cookies or other technologies to collect personally identifiable information.

If the configurations deployed for this Software Offering provide you as the customer the ability to collect personally identifiable information from end users via cookies and other technologies, you should seek your own legal advice about any laws applicable to such data collection, including any requirements for notice and consent.

For more information about the use of various technologies, including cookies, for these purposes, see IBM's Privacy Policy at <http://www.ibm.com/privacy> and IBM's Online Privacy Statement at <http://www.ibm.com/privacy/details> the section entitled "Cookies, Web Beacons and Other Technologies" and the "IBM Software Products and Software-as-a-Service Privacy Statement" at <http://www.ibm.com/software/info/product-privacy>.

Trademarks

IBM, the IBM logo, and [ibm.com](http://www.ibm.com) are trademarks or registered trademarks of International Business Machines Corp., registered in many jurisdictions worldwide. Other product and service names might be trademarks of IBM or other companies. A current list of IBM trademarks is available on the web at [Copyright and trademark information at www.ibm.com/legal/copytrade.shtml](http://www.ibm.com/legal/copytrade.shtml).

Intel, Intel logo, Intel Inside, Intel Inside logo, Intel Centrino, Intel Centrino logo, Celeron, Intel Xeon, Intel SpeedStep, Itanium, and Pentium are trademarks or registered trademarks of Intel Corporation or its subsidiaries in the United States and other countries.

Java and all Java-based trademarks and logos are trademarks or registered trademarks of Oracle and/or its affiliates.

UNIX is a registered trademark of The Open Group in the United States and other countries.

Index

Special Characters

[_exit subroutine](#) 463
[_LARGE_FILES](#) 130
[_system_configuration.ncpus](#) 562
[_system_configuration.ncpus variable](#) 562

Numerics

216840 102
41Map203831 46

A

[access subroutine](#) 148
[adb debug program](#)
 [adb debug program](#)
 [list of operators](#) 46
 [address maps](#)
 [displaying](#) 39
 [addresses](#)
 [displaying](#) 39
 [finding current](#) 39
 [forming](#) 39
 [binary files](#)
 [patching](#) 39
 [breakpoints](#) 30
 [C stack backtrace](#)
 [displaying](#) 39
 [commands, combining](#) 35
 [computing numbers](#) 39
 [creating scripts](#) 35
 [customizing](#) 35
 [data](#)
 [displaying](#) 39
 [data formats](#)
 [choosing](#) 39
 [data formatting, example](#) 56
 [default input formats](#)
 [setting](#) 35
 [directory dumps](#)
 [example](#) 54
 [displaying text](#) 39
 [examples](#)
 [data formatting](#) 56
 [directory dumps](#) 54
 [i-node dumps](#) 54
 [tracing multiple functions](#) 58
 [exiting](#) 29
 [expressions](#)
 [list of](#) 46
 [using integers](#) 34
 [using operators](#) 34
 [using symbols](#) 34
 [external variables](#)
 [displaying](#) 39
 [files](#)
 [adb debug program \(continued\)](#)
 [files \(continued\)](#)
 [locating values in](#) 39
 [writing](#) 39
 [i-node dumps](#)
 [example](#) 54
 [instructions](#)
 [displaying](#) 39
 [integers](#)
 [using in expressions](#) 34
 [list of subcommands](#) 46
 [list of variables](#) 46
 [maps](#)
 [memory, changing](#) 39
 [maps, address](#)
 [displaying](#) 39
 [maximum offsets](#)
 [setting](#) 35
 [memory](#)
 [changing](#) 39
 [memory maps](#)
 [changing](#) 39
 [numbers, computing](#) 39
 [operators](#)
 [using in expressions](#) 34
 [operators, list of](#) 46
 [output widths](#)
 [setting](#) 35
 [program execution](#)
 [controlling](#) 30
 [programs](#)
 [continuing execution](#) 30
 [preparing for debugging](#) 30
 [running](#) 30
 [single-stepping](#) 30
 [stopping](#) 30
 [stopping with keys](#) 30
 [prompts, using](#) 29
 [sample programs](#) 52
 [scripts, creating](#) 35
 [shell commands, using](#) 29
 [source files](#)
 [displaying and manipulating](#) 39
 [starting](#) 29
 [stopping a program](#) 30
 [subcommands, list of](#) 46
 [symbols](#)
 [using in expressions](#) 34
 [text, displaying](#) 39
 [tracing multiple functions, example of](#) 58
 [using prompts](#) 29
 [using shell commands](#) 29
 [values](#)
 [locating in files](#) 39
 [variables](#)
 [displaying external](#) 39
 [list of](#) 46

- adb debug program (*continued*)
 - variables (*continued*)
 - using [39](#)
- adb debug program overview [29](#)
- addresses
 - disk [123](#)
 - memory [588](#)
 - overview [588](#)
 - program [588](#)
- AIX transactional memory programming
 - Assembler [795](#)
 - Checkpoint State [795](#)
 - Compiler [795](#)
 - Debugger [795](#)
 - Signal delivery [795](#)
 - Software Failure Handler [795](#)
 - Trace facility [795](#)
- AIX vector programming [602](#)
- alarm subroutine [459](#)
- aliases
 - SMIT [720](#)
- allocating JFS2 file space [126](#)
- allocation
 - compressed file system [123](#)
 - JFS [123](#)
- allocation bitmaps [127](#)
- allocation groups [127](#)
- allocation policy
 - default [609](#)
 - Watson2 allocation policy [609](#)
- and [562](#)
- application binding [562](#)
- async-cancel safety [416](#)
- attributes object [413](#), [424](#), [429](#)
- aware [558](#)

B

- back end program [149](#)
- backup command [104](#)
- benefits of threads [489](#)
- bind CPU ID [562](#)
- binding a thread [177](#)
- bindprocessor [562](#)
- bindprocessor command [562](#)
- bindprocessor system call [562](#)
- blocks
 - boot [127](#)
 - data [127](#)
 - full logical [123](#)
 - indirect [123](#), [142](#)
 - logical [123](#)
 - partial logical [123](#)
 - super [127](#)
- boot block [127](#)
- buffer size operations [783](#)

C

- cancelability [416](#)
- cancelation [416](#), [417](#)
- capacity upgrade on demand (CUoD) [558](#)
- cbreak mode [17](#)

- characters
 - adding characters to windows [10](#)
 - adding to windows [10](#)
 - allowing 8-bit character return [10](#)
 - converting control characters to printables [10](#)
 - current
 - curses [3](#)
 - deleting [10](#)
 - getting from standard input [10](#)
 - returning if no input is available [10](#)
- chclass [562](#)
- chdir subroutine [116](#)
- check phase [563](#), [570](#)
- chmod subroutine [148](#)
- chown subroutine [148](#)
- chroot subroutine [116](#)
- class definitions [562](#)
- cleanup handlers [421](#)
- close subroutine
 - closing files with [140](#)
- command
 - smtctl [556](#)
- commands
 - backup [104](#)
 - cron [104](#)
 - diag [96](#)
 - errclear [104](#), [112](#)
 - errdemon [112](#)
 - errlogger [112](#)
 - errpt [96](#), [104](#), [112](#)
 - errstop [104](#)
 - errupdate [107](#), [112](#)
 - ls [104](#)
 - mycmd [769](#)
 - SCCS
 - list of [686](#)
 - trace [764](#), [769](#)
 - trcrpt [769](#)
 - trcstop [769](#)
- compiling a multithreaded program [478](#)
- compressed file system [123](#)
- condition variable
 - attributes
 - creation [429](#)
 - destruction [429](#)
 - object [429](#)
 - process-shared [464](#)
 - creation [429](#)
 - definition [429](#)
 - destruction [429](#)
 - usage [429](#)
 - wait [433](#)
 - waiting [429](#)
- configuring using malloc thread cache [637](#)
- contention scope
 - global [447](#)
 - local [447](#)
 - process [447](#)
 - system [447](#)
- contention-scope attribute [447](#)
- control characters
 - converting to printables [10](#)
- controlling processor use [177](#)
- controlling terminal [780](#)

- converter modules [787](#)
- creat subroutine
 - creating files with [139](#)
- creating
 - destroying [452](#)
 - keys [452](#)
- creating a copy-on-write mapped data file with the shmat subroutine [598](#)
- creating a mapped data file with the shmat subroutine [597](#)
- creating a thread [413](#)
- creation and destruction [424](#), [429](#)
- cron command [104](#)
- CUoD [558](#)
- curses
 - adding characters to windows [10](#)
 - characters
 - current [3](#)
 - returning if no input is available [10](#)
 - control characters [10](#)
 - converting control characters to printables [10](#)
 - converting termcap to terminfo [17](#)
 - current character [3](#)
 - current line [3](#)
 - cursors
 - controlling placement after refresh [9](#)
 - logical [3](#)
 - moving the logical cursor [9](#)
 - moving the physical cursor [9](#)
 - physical [3](#)
 - deleting characters [10](#)
 - displays [3](#)
 - getting characters from standard input [10](#)
 - inserting
 - blank lines in windows [10](#)
 - lines
 - current [3](#)
 - deleting [10](#)
 - logical cursor [3](#)
 - macros [3](#)
 - moving
 - logical cursor [9](#)
 - physical cursor [9](#)
 - naming conventions [3](#)
 - pads
 - creating [7](#)
 - deleting [7](#)
 - refreshing [7](#)
 - physical cursor [3](#)
 - printing
 - formatted printf on windows [10](#)
 - refreshing
 - pads [7](#)
 - windows [7](#)
 - screen [3](#)
 - scrolling
 - windows [10](#)
 - setting options [23](#)
 - starting [5](#)
 - stopping [5](#)
 - strings
 - adding to windows [10](#)
 - sub-windows
 - creating [7](#)
 - termcap

- curses (*continued*)
 - termcap (*continued*)
 - converting to terminfo [17](#)
 - terminology [3](#)
 - windows
 - clearing [10](#)
 - copying [7](#)
 - creating [7](#)
 - deleting [7](#)
 - drawing boxes around [7](#)
 - moving [7](#)
 - overlapping [7](#)
 - refreshing [7](#)
 - screens [3](#)
 - scrolling [10](#)
- curses compatibility [26](#)
- cursors
 - controlling placement after refresh [9](#)
 - getting location of logical cursor [9](#)
 - logical
 - curses [3](#)
 - moving logical curses [9](#)
 - moving physical [9](#)
 - physical
 - curses [3](#)

D

- data blocks [123](#), [127](#)
- data management application programming interface DMAPI [791](#)
- data types
 - pthread_once_t [451](#)
- dbx
 - debugging multithreaded program using [481](#)
 - symbolic debug program [60](#)
- dbx command
 - print subcommand [68](#)
 - step subcommand [72](#)
 - thread subcommand [68](#)
- dbx plug-in
 - alias
 - callback routine [77](#)
 - event types [77](#)
 - example [77](#)
 - fds
 - callback routine [77](#)
 - file format [77](#)
 - file names [77](#)
 - get_thread_context
 - callback routine [77](#)
 - header file [77](#)
 - Interface [77](#)
 - loading [77](#)
 - locate_symbol
 - callback routine [77](#)
 - location [77](#)
 - modules
 - callback routine [77](#)
 - naming [77](#)
 - overview [77](#)
 - print
 - callback routine [77](#)
 - process

- dbx plug-in (*continued*)
 - process (*continued*)
 - callback routine [77](#)
 - pthreads
 - callback routine [77](#)
 - read_memory
 - callback routine [77](#)
 - regions
 - callback routine [77](#)
 - session
 - callback routine [77](#)
 - set_pthread_context
 - callback routine [77](#)
 - set_thread_context
 - callback routine [77](#)
 - threads
 - callback routine [77](#)
 - unloading [77](#)
 - version control [77](#)
 - what_function
 - callback routine [77](#)
 - write_memory
 - callback routine [77](#)
- dbx plug-in callback routines [77](#)
- dbx symbolic debug program
 - .dbxinit file [75](#)
 - aliases
 - dbx subcommand, creating [75](#)
 - breakpoints [60](#)
 - calling procedures [66](#)
 - changing the current file [64](#)
 - command line editing [60](#)
 - dbx subcommand aliases
 - creating [75](#)
 - debugging spinlocks [75](#)
 - expressions
 - modifiers and operators for [66](#)
 - type checking [66](#)
 - files
 - .dbxinit [75](#)
 - current, changing [64](#)
 - reading dbx subcommands from [75](#)
 - source, displaying [64](#)
 - folding variables to lowercase and uppercase [66](#)
 - handling signals [66](#)
 - list of subcommands [94](#)
 - machine level debugging [72](#)
 - machine level programs
 - debugging [72](#)
 - running [72](#)
 - machine registers [72](#)
 - memory addresses [72](#)
 - modifiers
 - for expressions [66](#)
 - multiple threads [64](#)
 - names, scoping [66](#)
 - new dbx prompt
 - defining [75](#)
 - operators
 - for expressions [66](#)
 - print output
 - changing variables [66](#)
 - procedures
 - calling [66](#)
- dbx symbolic debug program (*continued*)
 - procedures (*continued*)
 - current, changing [64](#)
 - program control [60](#)
 - programs
 - controlling [60](#)
 - machine level [72](#)
 - machine level, running [72](#)
 - multiple threads [64](#)
 - running [60](#)
 - separating output from dbx [60](#)
 - prompts
 - defining [75](#)
 - reading dbx subcommands from a file [75](#)
 - running programs [60](#)
 - running shell commands [60](#)
 - scoping names [66](#)
 - separating dbx output from program output [60](#)
 - signals, handling [66](#)
 - source directory path
 - changing [64](#)
 - source files
 - displaying and manipulating [64](#), [66](#), [75](#)
 - spinlocks
 - debugging [75](#)
 - stack trace, displaying [66](#)
 - starting [60](#)
 - subcommands, list of [94](#)
 - tracing execution [60](#)
 - type checking in expressions [66](#)
 - using [60](#)
 - variables
 - changing print output [66](#)
 - displaying and modifying [66](#)
 - folding, lowercase and uppercase [66](#)
- deadlock [427](#)
- daemons
 - srcmstr [749](#)
 - trace [764](#)
- debug malloc tool [623](#)
- debugging multithreaded programs
 - dbx, using [481](#)
 - kernel debug program, using [481](#)
- debugging programs [28](#)
- dependencies [561](#), [562](#)
- descriptors [136](#)
- destructor routine [452](#)
- detached state [441](#)
- detachstate attribute [413](#)
- developing multithreaded program debuggers [485](#)
- developing multithreaded programs to examine modify pthread library objects [482](#)
- diag command [96](#)
- dialogs
 - SMIT [715](#)
- directories
 - changing
 - current [116](#)
 - root [116](#)
 - linking [134](#)
 - overview [116](#)
 - status [134](#)
 - working with
 - overview [116](#)

- directories (*continued*)
 - working with (*continued*)
 - subroutines for [116](#)
- disk address format [123](#)
- disk fragments [127](#)
- disk i-nodes [120](#), [127](#)
- disk space allocation [123](#)
- DLPAR [558](#), [563](#)
- DR_FAIL [570](#)
- DR_resource_POST_ERROR [570](#)
- DR_SUCCESS [570](#)
- DR_WAIT [570](#)
- Draft 7 applications
 - porting [480](#)
- dri_cpu [570](#)
- dri_cpu structure [570](#)
- dri_mem [570](#)
- dri_mem structure [570](#)
- dynamic logical partitioning [558](#), [562](#), [563](#)
- dynamic logical partitioning [561](#)
- dynamic memory guarding [181](#)
- Dynamic Processor Deallocation [177](#)

E

- ECHO directive [497](#)
- entry point routine [413](#)
- environment variables
 - NUM_SPAREVP [478](#)
- errclear command [104](#)
- error log
 - transferring to another system [98](#)
- error log descriptors [96](#)
- error logging
 - cleaning an error log [104](#)
 - commands [112](#)
 - copying an error log [111](#)
 - example report [104](#)
 - files [112](#)
 - generating a report [104](#)
 - kernel services [112](#)
 - managing [98](#)
 - overview [96](#)
 - reading an error report [104](#)
 - stopping an error log [104](#)
 - subroutines [112](#)
- error logging facility [97](#)
- error notification [101](#)
- error report
 - detailed example [104](#)
 - generating [104](#)
 - summary example [104](#)
- errpt command [96](#), [104](#), [112](#)
- errstop command [104](#), [112](#)
- errupdate command [107](#), [112](#)
- event parameter [570](#)
- event_in_prog parameter [570](#)
- event_list parameter [570](#)
- example
 - dbx plug-in [77](#)
- example adb program adbsamp3 [53](#)
- example adb program: adbsamp2 [53](#)
- example program for the lex and yacc Programs [515](#)
- example programs

- example programs (*continued*)
 - manipulating characters
 - isalnum (ctype) routine [706](#)
 - isalpha (ctype) routine [706](#)
 - isascii (ctype) routine [706](#)
 - iscntrl (ctype) routine [706](#)
 - isdigit (ctype) routine [706](#)
 - islower (ctype) routine [706](#)
 - ispunct (ctype) routine [706](#)
 - isspace (ctype) routine [706](#)
 - isupper (ctype) routine [706](#)
- exec subroutine [463](#)
- extended regular expressions
 - lex command [493](#)

F

- fast paths
 - SMIT [720](#)
- fclear subroutine [141](#)
- fdpr
 - debugging reordered executables [72](#)
- FIFO (first-in, first-out)
 - understanding [140–143](#), [145](#), [146](#)
- file descriptor tables
 - definition [136](#)
- file descriptors
 - definition [136](#)
 - duplicating
 - dup subroutine [136](#)
 - fcntl subroutine [136](#)
 - fork subroutine [136](#)
 - managing [136](#)
 - preset values [136](#)
 - resource limit [136](#)
- file systems
 - bitmap [123](#)
 - compressed [123](#)
 - fragment map [123](#)
 - fragmented [123](#)
 - layout [127](#)
 - overview [113](#)
 - quotas [123](#)
 - types
 - creating [149](#)
- file types
 - overview [114](#)
- file-system helpers
 - execution syntax [150](#)
 - operations [149](#)
 - sample invocations [149](#)
- files
 - access modes [148](#)
 - allocating space to [123](#)
 - closing [140](#)
 - creating [139](#)
 - input and output (I/O) [140–143](#), [145](#), [146](#)
 - large
 - _LARGE_FILES [130](#)
 - 64-bit file system [130](#)
 - allocation in file systems [123](#)
 - common pitfalls [130](#)
 - implications for existing programs [130](#)
 - open protection [130](#)

- files (*continued*)
 - large (*continued*)
 - porting applications [130](#)
 - writing programs that access [130](#)
 - linking [134](#)
 - locking fields [121](#)
 - masks [139](#)
 - opening [139](#)
 - overview [114](#)
 - pipes [140–143](#), [145](#), [146](#)
 - reading [140–143](#), [145](#), [146](#)
 - SCCS
 - controlling [684](#)
 - creating [683](#)
 - detecting damage [685](#)
 - editing [683](#)
 - repairing damage [685](#)
 - tracking [684](#)
 - updating [683](#)
 - sharing open [136](#)
 - status [147](#)
 - termcap [17](#)
 - terminfo [17](#)
 - truncating [140–143](#), [145](#), [146](#)
 - working with
 - subroutines for [114](#)
 - writing [140–143](#), [145](#), [146](#)
- first-in first-out scheduling policy [444](#)
- floating-point exceptions
 - disabled and enabled processing [155](#)
 - subroutines [154](#)
- fork cleanup handlers [463](#)
- fork subroutine [463](#)
- fragment map [123](#)
- fragmented file system [123](#)
- fragments
 - disk [127](#)
 - map [123](#)
- front end program [149](#)
- ftruncate subroutine [140–143](#), [145](#), [146](#)
- full logical block [123](#)
- fullstat subroutine [147](#)

G

- generic trace channels [769](#)

H

- h_arg parameter [570](#)
- h_token parameter [570](#)
- handlers, extended [570](#)
- Hardware Management Console [558](#)
- heap
 - 64-bit applications [609](#)
- help
 - SMIT [726](#)
- helpers [149](#)
- hlpadb [46](#)

I

- i-nodes

- i-nodes (*continued*)
 - definition [114](#)
 - disk [127](#)
 - i-number byte offset [116](#)
 - JFS2 [121](#)
 - modifying [120](#)
 - timestamp
 - changing [147](#)
- I/O offset
 - absolute [140](#)
 - and read subroutine [140–143](#), [145](#), [146](#)
 - and write subroutine [140–143](#), [145](#), [146](#)
 - description [140–143](#), [145](#), [146](#)
 - end_relative [140](#)
 - manipulating [140–143](#), [145](#), [146](#)
 - relative [140](#)
- in-core i-nodes [121](#)
- index nodes [120](#)
- indirect blocks [123](#)
- inheritance protocol [448](#)
- inheritsched attribute [444](#)
- init command
 - SRC [749](#)
- initializations [451](#)
- input and output handling [163](#)
- input modes [17](#)
- inter-process communication IPC limits [594](#)
- introduction [588](#)
- IPC (interprocess channel) [114](#)

J

- J2_CFG_ASSIST ioctl operation [153](#)
- JFS
 - disk space allocation [123](#)
- JFS2
 - i-nodes [121](#)
- JFS2 file system layout [129](#)
- joining threads [441](#)

K

- kernel extensions [570](#)
- kernel programming
 - multiprocessor issues [181](#)
- kernel services [570](#)
- kernel thread [408–412](#)
- key destruction [452](#)
- keys
 - creating [452](#)
 - destroying [452](#)
 - destruction [452](#)
- kill subroutine [459](#)

L

- large files
 - common pitfalls
 - arithmetic overflows [130](#)
 - failure to include proper headers [130](#)
 - file size limits [130](#)
 - fseek/ftell [130](#)
 - imbedded file offsets [130](#)

- large files (*continued*)
 - common pitfalls (*continued*)
 - improper data types [130](#)
 - parameter mismatches [130](#)
 - string conversions [130](#)
 - open protection [130](#)
 - porting applications to [130](#)
 - using `_LARGE_FILES` [130](#)
 - using 64-bit file system [130](#)
 - writing programs that access [130](#)
- large program support [172](#)
- lazy loading [583](#)
- lex command
 - compiling the lexical analyzer [491](#)
 - defining substitution strings [496](#)
 - extended regular expressions [493](#)
 - lexical analyzer [497](#)
 - matched strings [497](#)
 - passing code to program [496](#)
 - start conditions [501](#)
- lex Command
 - end-of-file processing [497](#)
 - input/output subroutines [497](#)
 - lex library [497](#)
 - lexical analyzer [491](#)
 - operators [493](#)
- lex Library [497](#)
- lexical analyzer
 - parser program [492](#)
- libpthread_compat.a library [412](#)
- libpthread.a library [408–412](#)
- library model test [447](#)
- line discipline module `ldterm` [784](#)
- linking
 - run time [581](#)
- links
 - directory [134](#)
 - symbolic [134](#)
- list numerical manipulation subroutines [694](#)
- list of additional curses subroutines [27](#)
- list of character manipulation subroutines [689](#)
- list of executable program creation subroutines [691](#)
- list of long integer numerical manipulation subroutines [696](#)
- list of multithreaded programming subroutines [700](#)
- list of threads processes interactions subroutines [462](#)
- list processes subroutines [697](#)
- list programmer's workbench library subroutines [701](#)
- list scheduling subroutines [447](#)
- list security auditing subroutines [702](#)
- list string manipulation subroutines [705](#)
- list synchronization subroutines [450](#)
- list thread advanced feature subroutines [468](#)
- list_size parameter [570](#)
- loader domains
 - using [789](#)
- locking
 - creating user locking services [181](#)
- logical block [123](#)
- logical partitioning [558](#)
- logical volume manager
 - library subroutines [153](#)
- logical volume programming [152](#)
- long locks [456](#)
- longjmp subroutine [459](#)
- LPAR [558](#)
- ls command [104](#)
- lsclass [562](#)
- lseek subroutine [140–143, 145, 146](#)
- lsrset [562](#)

M

- m4
 - built-in macros [532](#)
 - changing quote characters [532](#)
 - checking for defined macros [532](#)
 - conditional expressions [532](#)
 - creating user-defined macros [532](#)
 - integer arithmetic [532](#)
 - macro processing with arguments [532](#)
 - manipulating files [532](#)
 - manipulating strings [532](#)
 - predefined macros [532](#)
 - printing names and definitions [532](#)
 - quote characters [532](#)
 - redirecting output [532](#)
 - removing macro definitions [532](#)
 - system programs [532](#)
 - unique names [532](#)
 - using the macro processor [532](#)
- macros
 - PTHREAD_ONCE_INIT [451](#)
- main Subroutine [497](#)
- malloc buckets [630](#)
- malloc detect [636](#)
- malloc disclaim [636](#)
- malloc log [634](#)
- malloc multiheap [629](#)
- malloc subsystem
 - alloca [609](#)
 - allocation policy [611](#)
 - allocation policy, default [609](#)
 - calloc [609](#)
 - free [609](#)
 - heap [609](#)
 - mallinfo [609](#)
 - malloc [609](#)
 - mallopt [609](#)
 - process heap [609](#)
 - realloc [609](#)
 - valloc [609](#)
- malloc trace [633](#)
- manipulating soft labels [26](#)
- masks [139](#)
- memory management
 - allocating memory [609](#)
 - listing of memory manipulation services [601](#)
 - program address space
 - overview [588](#)
- memory mapping
 - listing of memory mapping services [602](#)
 - mmap comparison with `shmat` [590](#)
 - mmap compatibility considerations [590](#)
 - overview [590](#)
 - semaphore subroutines overview [590](#)
- memory removal failure [561](#)
- menus

- menus (*continued*)
 - SMIT [715](#)
- minicurses
 - characters
 - adding characters to windows [10](#)
 - strings
 - adding to windows [10](#)
- mkfifo subroutine [139](#)
- mknod subroutine
 - creating regular files with [139](#)
- mount command [149](#)
- mount helpers
 - execution syntax [149](#)
 - overview [149](#)
- multi-threading simultaneous [556](#)
- multiprocessor programming [490](#)
- multithreaded program
 - Compiler Invocation [478](#)
 - compiling [478](#)
 - debugging [478](#)
- multithreaded programming [407](#)
- mutex
 - attributes
 - creation [424](#)
 - destruction [424](#)
 - object [424](#)
 - prioceiling [448](#)
 - process-shared [464](#)
 - creation [424](#)
 - definition [424](#)
 - destruction [424](#)
 - locking [424](#)
 - protocol attributes [448](#)
 - protocols [448](#)
 - unlocking [424](#)
 - usage [424](#)
- mycmd command [769](#)
- mycpu [558](#)

N

- name parameter [570](#)
- named shared library areas [584](#)
- notify object class (SRC)
 - creating a subsystem notification method [750](#)
 - removing a subsystem notification method [750](#)
- NUM_SPAREVP environment variable [478](#)

O

- O_DEFER [141](#)
- object classes
 - SMIT [717](#)
- objects [541](#)
- ODM (Object Data Manager)
 - descriptors
 - link [541](#)
 - method [541](#)
 - terminal [541](#)
 - example code
 - adding objects [553](#)
 - creating object classes [553](#)
 - list of commands [551](#)

- ODM (Object Data Manager) (*continued*)
 - list of subroutines [551](#)
 - object classes
 - adding objects [541](#)
 - creating [541](#)
 - definition [541](#)
 - locking [541](#)
 - unlocking [541](#)
 - objects
 - adding to object classes [541](#)
 - definition [541](#)
 - searching for [541](#)
 - predicates
 - comparison operators [541](#)
 - constants in [541](#)
 - descriptor names [541](#)
- offset [140–143](#), [145](#), [146](#)
- one-time initialization [451](#)
- open subroutine
 - creating files with [139](#)
 - opening a file with [139](#)
- operating system
 - libraries [713](#)
- operations [563](#)
- options (threads library) [464](#)
- overviews
 - make command
 - creating description files [519–526](#), [528–532](#), [574–579](#)
 - creating target files [519–526](#), [528–532](#), [574–579](#)
 - internal rules [519–526](#), [528–532](#), [574–579](#)
 - using with non-SCCS files [519–526](#), [528–532](#), [574–579](#)
 - using with SCCS files [519–526](#), [528–532](#), [574–579](#)

P

- pads
 - creating [7](#)
 - curses [3](#)
 - deleting [7](#)
 - refreshing [7](#)
- paging space programming requirements [600](#)
- parser
 - lexical analyzer [492](#)
- parser operation generated yacc command [511](#)
- partial logical block [123](#)
- pbsearchsort [709](#)
- PCI slot removal failure [561](#)
- pclose subroutine [140–143](#), [145](#), [146](#)
- permissions
 - directories [148](#)
 - files [148](#)
- pipe subroutine [140–143](#), [145](#), [146](#)
- pipes
 - child processes [140–143](#), [145](#), [146](#)
 - creating with mkfifo [139](#)
- plug-in, dbx
 - alias
 - callback routine [77](#)
 - event types [77](#)
 - example [77](#)
 - fds
 - callback routine [77](#)

- plug-in, dbx (*continued*)
 - file format [77](#)
 - file names [77](#)
 - get_thread_context
 - callback routine [77](#)
 - header file [77](#)
 - Interface [77](#)
 - loading [77](#)
 - locate_symbol
 - callback routine [77](#)
 - location [77](#)
 - modules
 - callback routine [77](#)
 - naming [77](#)
 - print
 - callback routine [77](#)
 - process
 - callback routine [77](#)
 - pthreads
 - callback routine [77](#)
 - read_memory
 - callback routine [77](#)
 - regions
 - callback routine [77](#)
 - session
 - callback routine [77](#)
 - set_pthread_context
 - callback routine [77](#)
 - set_thread_context
 - callback routine [77](#)
 - threads
 - callback routine [77](#)
 - unloading [77](#)
 - version control [77](#)
 - what_function
 - callback routine [77](#)
 - write_memory
 - callback routine [77](#)
- plug-in, dbx callback routines [77](#)
- plug-in,dbx
 - overview [77](#)
- pobeve dynamic tracing facility [184](#)
- popen subroutine [140–143](#), [145](#), [146](#)
- POSIX thread [412](#)
- post phase [563](#), [570](#)
- post-add phase [570](#)
- post-error phase [570](#)
- post-remove phase [570](#)
- pre phase [563](#)
- pre-add phase [570](#)
- pre-phase [570](#)
- pre-remove phase [570](#)
- prioceiling attribute [448](#)
- priority
 - inheritance protocol [448](#)
 - protection protocol [448](#)
- priority inversion [448](#)
- priority scheduling POSIX option [465](#)
- process
 - properties [408–412](#)
- process heap
 - 32-bit applications [609](#)
- process-shared attribute [464](#)
- processes (*continued*)
 - using pipes [140–143](#), [145](#), [146](#)
- processor
 - numbers [176](#)
 - ODM names [176](#)
 - states [176](#)
- processor bindings [562](#)
- processor numbering [562](#)
- processor removal failure [561](#)
- protect keys [167](#)
- protection protocol [448](#)
- protocol attribute [448](#)
- protocols
 - inheritance [448](#)
 - mutex [448](#)
 - priority inheritance [448](#)
 - priority protection [448](#)
 - protection [448](#)
- ps [562](#)
- ps command [562](#)
- pthread [412](#)
- pthread_atfork subroutine [463](#)
- pthread_attr_destroy subroutine [413](#)
- pthread_attr_getdetachstate subroutine [413](#)
- pthread_attr_getinheitsched subroutine [444](#)
- pthread_attr_getsatckaddr subroutine [464](#)
- pthread_attr_getschedparam [447](#)
- pthread_attr_getschedparam subroutine
 - schedparam attribute [444](#)
- pthread_attr_getscope subroutine [447](#)
- pthread_attr_getstacksize subroutine [464](#)
- pthread_attr_init subroutine [413](#)
- pthread_attr_setdetachstate subroutine [413](#)
- pthread_attr_setinheritsched subroutine [444](#)
- pthread_attr_setsatckaddr subroutine [464](#)
- pthread_attr_setschedparam [447](#)
- pthread_attr_setschedparam subroutine
 - schedparam attribute [444](#)
- pthread_attr_setschedpolicy subroutine [444](#)
- pthread_attr_setscope subroutine [447](#)
- pthread_attr_setstacksize subroutine [464](#)
- pthread_attr_t data type [413](#)
- pthread_cancel subroutine [417](#)
- pthread_cleanup_pop subroutine [416](#)
- pthread_cleanup_push subroutine [416](#)
- pthread_cond_broadcast subroutine [429](#)
- pthread_cond_destroy subroutine [429](#)
- pthread_cond_init subroutine [429](#)
- PTHREAD_COND_INITIALIZER macro [429](#)
- pthread_cond_signal subroutine [429](#)
- pthread_cond_t data type [429](#)
- pthread_cond_timedwait subroutine [429](#)
- pthread_cond_wait subroutine [429](#)
- pthread_condattr_destroy subroutine [429](#)
- pthread_condattr_init subroutine [429](#)
- pthread_condattr_t data type [429](#)
- pthread_create subroutine [413](#)
- pthread_equal subroutine [413](#)
- pthread_exit subroutine [416](#)
- pthread_getschedparam [447](#)
- pthread_getschedparam subroutine
 - schedparam attribute [444](#)
 - schedpolicy attribute [444](#)
- pthread_getspecific subroutine [452](#)

- pthread_join subroutine [441](#)
- pthread_key_create subroutine [452](#)
- pthread_key_delete subroutine [452](#)
- pthread_key_t data type [452](#)
- pthread_kill subroutine [459](#)
- pthread_mutex_destroy subroutine [424](#)
- pthread_mutex_init subroutine [424](#)
- pthread_mutex_lock subroutine [424](#)
- pthread_mutex_t data type [424](#)
- pthread_mutex_trylock subroutine [424](#)
- pthread_mutex_unlock subroutine [424](#)
- pthread_mutexattr_destroy subroutine [424](#)
- pthread_mutexattr_getprioceiling subroutine [448](#)
- pthread_mutexattr_getprotocol subroutine [448](#)
- pthread_mutexattr_init subroutine [424](#)
- pthread_mutexattr_setprioceiling subroutine [448](#)
- pthread_mutexattr_setprotocol subroutine [448](#)
- pthread_mutexattr_t data type [424](#)
- pthread_once subroutine [451](#)
- PTHREAD_ONCE_INIT macro [451](#)
- pthread_once_t data type [451](#)
- PTHREAD_PRIO_INHERIT [448](#)
- PTHREAD_PRIO_NONE [448](#)
- PTHREAD_PRIO_PROTECT [448](#)
- pthread_self subroutine [413](#)
- pthread_setcancelstate subroutine [416](#)
- pthread_setcanceltype subroutine [416](#)
- pthread_setschedparam subroutine [447](#)
- pthread_setspecific subroutine [452](#)
- pthread_t data type [413](#)
- pthread_testcancel [416](#)
- pthread_yield subroutine [429](#)

Q

- quotas [123](#)

R

- race condition [424](#)
- raise subroutine [459](#)
- read subroutine [140–143](#), [145](#), [146](#)
- reconfig_handler_complete [570](#)
- reconfig_handler_complete kernel service [570](#)
- reconfig_handler_complete routine [570](#)
- reconfig_handler_complete subroutine [570](#)
- reconfig_register [570](#)
- reconfig_register function [570](#)
- reconfig_register subroutine [570](#)
- reconfig_register_list subroutine [570](#)
- reconfig_unregister function [570](#)
- reconfiguration handlers [570](#)
- registering handlers [570](#)
- registering kernel extension handlers [570](#)
- regular expressions
 - lex command [493](#)
- REJECT directive [497](#)
- remove subroutine [134](#)
- resource_info parameter [570](#)
- rmdir subroutine [134](#)
- round-robin scheduling policy [444](#)
- routines
 - destructor [452](#)

- run time linking [581](#)

S

- SAFE [558](#)
- SCCS
 - commands
 - list of [686](#)
 - files
 - controlling [684](#)
 - creating [683](#)
 - detecting damage [685](#)
 - editing [683](#)
 - repairing damage [685](#)
 - tracking [684](#)
 - updating [683](#)
 - flags and parameters [682](#)
- SCCS retrieval
 - suppress [519–526](#), [528–532](#), [574–579](#)
- sched_yield [447](#)
- sched_yield subroutine [444](#)
- scheduling
 - parameters [444](#)
 - policies [444](#)
 - priority [444](#)
 - synchronization of [448](#)
- screen types
 - SMIT [715](#)
- Scripts [563](#)
- scripts, using [563](#)
- searching and sorting
 - example program [709](#)
- sed command
 - starting the editor [574](#)
 - using string replacement [574](#)
 - using text in commands [574](#)
 - using the command summary [574](#)
- sed program information [574](#)
- selectors
 - SMIT [715](#)
- semaphores (from OSF/1) [456](#)
- setjmp subroutine [459](#)
- shared libraries
 - creating [586](#)
 - lazy loading [583](#)
- shared libraries and shared memory [580](#)
- shared memory
 - mmap comparison with shmat [590](#)
 - overview [590](#)
- shared objects
 - creating [581](#)
- sigaction subroutine [459](#)
- siglongjmp subroutine [459](#)
- signal
 - delivery [459](#)
 - generation [459](#)
 - handlers [459](#)
 - masks [459](#)
- signals [563](#)
- signals, using [563](#)
- sigprocmask subroutine [459](#)
- SIGRECONFIG [563](#)
- SIGRECONFIG signal [563](#)
- sigsetjmp subroutine [459](#)

- sigthreadmask subroutine [459](#)
- sigwait subroutine [429](#), [459](#)
- Simultaneous Multithreading [556](#)
- SMIT (System Management Interface Tool)
 - aliases [720](#)
 - dialogs
 - designing [715](#)
 - executing [723](#)
 - generating [723](#)
 - example program [739](#)
 - fast paths [720](#)
 - generating commands with [723](#)
 - help
 - creating [726](#)
 - understanding [726](#)
 - information command descriptors
 - cmd_to*_postfix [721](#)
 - cmd_to_discover [721](#)
 - understanding [721](#)
 - menus
 - designing [715](#)
 - name selects
 - designing [715](#)
 - object class
 - for aliases [727](#)
 - object classes
 - dialog [732](#)
 - dialog header [736](#)
 - menu [727](#)
 - selector header [729](#)
 - understanding [717](#)
 - screen types [715](#)
 - tasks
 - adding [725](#)
 - debugging [726](#)
- smit errclear command [104](#)
- smit errpt command [110](#)
- smit trace command [769](#)
- smtctl [556](#)
- software models
 - controller/worker [409](#)
 - divide-and-conquer [408–412](#)
 - master/slave [408–412](#)
 - producer/consumer [408–412](#)
- Software Partitioning [562](#)
- source code control system SCCS [681](#)
- SRC
 - capabilities [749](#)
 - defining subservers to the SRC object class [763](#)
 - defining subsystems to the SRC object class [763](#)
 - init command [749](#)
 - list of subroutines [764](#)
 - modifying subserver object definitions [763](#)
 - modifying subsystem object definitions [763](#)
 - operations [749](#)
 - removing subserver object definitions [763](#)
 - removing subsystem object definitions [763](#)
- SRC communication types
 - message queues (IPC)
 - overview [755](#)
 - programming requirements [757](#)
 - signals
 - programming requirements [757](#)
- SRC communication types (*continued*)
 - sockets
 - overview [755](#)
 - programming requirements [757](#)
- SRC object classes
 - descriptors [750](#)
 - notify object overview [750](#)
 - subserver type object overview [750](#)
 - subsystem environment object overview [750](#)
- SRC subsystem programming requirements
 - processing status requests [757](#)
 - receiving SRC request packets
 - using message queue (IPC) communication [757](#)
 - using signals communication [757](#)
 - using sockets communication [757](#)
 - returning continuation packets [757](#)
 - returning error packets [757](#)
 - returning subsystem reply packets [757](#)
- srcmstr [749](#)
- stack address POSIX option [465](#)
- stackaddr attribute [464](#)
- stacksize attribute [464](#)
- status
 - directories [134](#)
- status subroutines
 - overview [147](#)
- storage protect keys [167](#)
- string manipulation
 - using sed command [574](#)
- strings
 - adding to windows
 - curses [10](#)
 - minicurses [10](#)
- structures
 - posix_trace_event_info [772](#)
 - posix_trace_status_info [772](#)
- sub-windows
 - refreshing [7](#)
- subroutines
 - alloca [609](#)
 - calloc [609](#)
 - chmod [148](#)
 - chown [148](#)
 - controlling files and directories
 - list of [692](#)
 - errlog [112](#)
 - free [609](#)
 - mallinfo [609](#)
 - malloc [609](#)
 - mallopt [609](#)
 - pthread_getspecific [452](#)
 - pthread_setspecific [452](#)
 - realloc [609](#)
 - valloc [609](#)
- subroutines, example programs, and libraries [686](#)
- subsystems
 - using the system resource controller [757](#)
- super block [127](#)
- suppress
 - SCCS retrieval [519–526](#), [528–532](#), [574–579](#)
- symbolic debug program, dbx [60](#)
- synchronization
 - mutex [424](#)
- synchronization scheduling

- synchronization scheduling (*continued*)
 - definition [448](#)
- sysconf subroutine [464](#)
- system [588](#)
- system calls
 - bindprocessor [562](#)
- system environment
 - Dynamic Processor Deallocation [177](#)
- system file tables [136](#)
- system management interface tool SMIT [714](#)
- system resource controller [750](#), [755](#), [757](#), [763](#), [764](#)

T

- termcap file [17](#)
- terminal devices
 - tty subsystem overview [780](#)
- terminals
 - multiple [17](#)
 - setting input and output modes [17](#)
- terminating a thread [416](#)
- terminfo file [17](#)
- thread
 - attributes
 - contention-scope [447](#)
 - creating [413](#)
 - destroying [413](#)
 - detachstate [413](#)
 - inheritsched [444](#)
 - object [413](#)
 - schedparam [444](#)
 - schedpolicy [444](#)
 - stackaddr [464](#)
 - stacksize [464](#)
 - binding [177](#)
 - Compiler Invocation [478](#)
 - concurrency level [447](#)
 - contention scope [447](#)
 - data types [464](#)
 - default values [464](#)
 - definition [408–412](#)
 - detached state [441](#)
 - ID [413](#)
 - join [441](#)
 - kernel [408–412](#)
 - libpthreads_compat.a library [412](#)
 - libpthreads.a library [412](#)
 - library [408–412](#)
 - limits [464](#)
 - models [408–412](#)
 - POSIX thread [412](#), [478](#)
 - properties [408–412](#)
 - pthread [412](#)
 - supported [468](#)
 - user [408–412](#)
- thread-specific data
 - concurrent creation [453](#)
 - definition [452](#)
 - destroying data [455](#)
 - destructor [452](#)
 - destructor routine [452](#)
 - key [452](#)
 - setting values [452](#)
 - swapping data [452](#)

- thread-specific data (*continued*)
 - usage [452](#)
- threads
 - creating [413](#)
 - object-oriented interface [408–412](#)
 - synchronization [423](#)
 - terminating [416](#)
- threads library
 - cancelation [416](#)
 - cleanup [416](#)
 - condition variable
 - attributes creation and destruction [429](#)
 - synchronization point [429](#)
 - timed wait [429](#)
 - dynamic initialization
 - threadsafe [451](#)
 - exiting a thread [416](#)
 - freeing returned data [441](#)
 - join [441](#)
 - mutex
 - attributes creation and destruction [424](#)
 - read/write locks (from POSIX) [456](#)
 - thread creation [478](#)
- threadsafe
 - SRC subroutines [764](#)
- threadsafe code [473](#), [637](#)
- timestamp
 - changing [147](#)
- tools utilities [1](#)
- trace
 - configuring [769](#)
 - controlling [764](#)
 - generating reports [769](#)
 - recording trace event data [770](#)
 - starting [769](#)
 - starting a trace [769](#)
 - stopping [769](#)
 - using generic trace channels [769](#)
- trace command
 - configuring [769](#)
- trace facility
 - overview of [764](#)
- trace hook identifiers [769](#)
- tracing
 - configuring [769](#)
 - definition [776](#)
 - starting [769](#)
 - trace event type [776](#)
 - trace subroutines [777](#)
 - tracing stream attributes [776](#)
- tracing data
 - tracing user application [772](#)
- tracing data structures
 - tracing user application [772](#)
- trcrpt command [769](#)
- trcstop command [769](#)
- troubleshooting [561](#)
- truncate subroutine [140–143](#), [145](#), [146](#)
- tty
 - current modes [17](#)
 - restoring modes [17](#)
 - saving modes [17](#)
- tty (terminal type)
 - definition [780](#)

tty (terminal type) *(continued)*
 examples [780](#)
tty drivers [788](#)
tty subsystem [780](#)

yyreject Subroutine [497](#)
yywlang external variable [497](#)
yywrap Subroutine [497](#)

U

unlink subroutine [134](#)
unmount command [149](#)
user defined malloc replacement [619](#)
user thread [408–412](#)
using read-write locks [434](#)
using yacc grammar file [504](#)
utimes subroutine [147](#)

V

variables [562](#)
video attributes
 setting [23](#)
virtual memory
 addressing
 overview [588](#)
virtual processor [408–412](#)
vital product data (VPD) [104](#)
VP [408–412](#)
VPD (vital product data) [104](#)
vue functions [374](#)

W

windows
 clearing [10](#)
 copying [7](#)
 curses [3](#)
 drawing box around [7](#)
 moving [7](#)
 overlapping [7](#)
 refreshing [7](#)
 refreshing multiple [7](#)
 scrolling [10](#)
winsize structure [783](#)
wlmctrl [562](#)
working color [23](#)
working JFS2 directories [118](#)
Workload Manager [562](#)
write subroutine [140–143](#), [145](#), [146](#)

Y

yacc command
 actions [509](#)
 ambiguous rules [513](#)
 creating a parser [502](#)
 declarations [505](#)
 error handling [510](#)
 grammar file [502](#)
 rules [507](#)
 turning on debug mode [515](#)
yacc program [492](#)
yyless external variable [497](#)
yyless Subroutine [497](#)
yyless Subroutine [497](#)

