

AIX Version 6.1



# Kernel Extensions and Device Support Programming Concepts



AIX Version 6.1



# Kernel Extensions and Device Support Programming Concepts

**Note**

Before using this information and the product it supports, read the information in “Notices,” on page 391.

**Third Edition (October 2009)**

This edition applies to AIX Version 6.1 and to all subsequent releases of this product until otherwise indicated in new editions.

A reader's comment form is provided at the back of this publication. If the form has been removed, address comments to Information Development, Department 04XA-905-6B013, 11501 Burnet Road, Austin, Texas 78758-3400. To send comments electronically, use this commercial Internet address: [pserinfo@us.ibm.com](mailto:pserinfo@us.ibm.com). Any information that you supply may be used without incurring any obligation to you.

© **Copyright International Business Machines Corporation 1997, 2009.**

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

---

# Contents

<b>About This Book</b> . . . . .	vii
How to Use This Book . . . . .	vii
Highlighting . . . . .	vii
Case-Sensitivity in AIX. . . . .	vii
ISO 9000 . . . . .	vii
Related Publications . . . . .	vii
<b>Chapter 1. Kernel Environment</b> . . . . .	1
Understanding Kernel Extension Symbol Resolution . . . . .	1
Understanding Execution Environments . . . . .	5
Understanding Kernel Threads . . . . .	6
Using Kernel Processes . . . . .	8
Accessing User-Mode Data While in Kernel Mode . . . . .	11
Understanding Locking . . . . .	12
Understanding Exception Handling . . . . .	14
Using Kernel Extensions for 64-bit Processes . . . . .	18
64-bit Kernel Extension Programming Environment . . . . .	19
Related Information. . . . .	20
<b>Chapter 2. System Calls</b> . . . . .	23
Differences Between a System Call and a User Function . . . . .	23
Understanding Protection Domains . . . . .	23
Understanding System Call Execution . . . . .	24
Accessing Kernel Data While in a System Call. . . . .	24
Passing Parameters to System Calls . . . . .	25
Preempting a System Call . . . . .	29
Handling Signals While in a System Call . . . . .	30
Handling Exceptions While in a System Call . . . . .	31
Understanding Nesting and Kernel-Mode Use of System Calls . . . . .	31
Page Faulting within System Calls . . . . .	31
Returning Error Information from System Calls. . . . .	32
System Calls Available to Kernel Extensions . . . . .	32
Related Information. . . . .	33
<b>Chapter 3. Virtual File Systems.</b> . . . . .	35
Logical File System Overview . . . . .	35
Virtual File System Overview . . . . .	36
Understanding Data Structures and Header Files for Virtual File Systems . . . . .	38
Configuring a Virtual File System. . . . .	39
Related Information. . . . .	39
<b>Chapter 4. Kernel Services</b> . . . . .	41
Categories of Kernel Services . . . . .	41
I/O Kernel Services. . . . .	41
Block I/O Buffer Cache Kernel Services: Overview . . . . .	48
Understanding Interrupts . . . . .	49
Understanding DMA Transfers. . . . .	50
Kernel Extension and Device Driver Management Services . . . . .	58
Locking Kernel Services . . . . .	60
File Descriptor Management Services . . . . .	63
Logical File System Kernel Services . . . . .	63
Programmed I/O (PIO) Kernel Services . . . . .	64
Memory Kernel Services . . . . .	64

Understanding Virtual Memory Manager Interfaces . . . . .	68
Message Queue Kernel Services . . . . .	72
Network Kernel Services . . . . .	72
Process and Exception Management Kernel Services . . . . .	74
RAS Kernel Services . . . . .	77
Security Kernel Services . . . . .	77
Timer and Time-of-Day Kernel Services . . . . .	78
Using Fine Granularity Timer Services and Structures . . . . .	79
Using Multiprocessor-Safe Timer Services . . . . .	81
Virtual File System (VFS) Kernel Services . . . . .	82
Related Information . . . . .	82
<b>Chapter 5. Asynchronous I/O Subsystem . . . . .</b>	<b>83</b>
How Do I Know if I Need to Use AIO? . . . . .	84
Functions of Asynchronous I/O . . . . .	86
Asynchronous I/O Subroutines . . . . .	87
Subroutines Affected by Asynchronous I/O . . . . .	88
64-bit Enhancements . . . . .	88
LEGACY AIO Extended Functionality . . . . .	88
Related Information . . . . .	91
<b>Chapter 6. Device Configuration Subsystem . . . . .</b>	<b>93</b>
Scope of Device Configuration Support . . . . .	93
Device Configuration Subsystem Overview . . . . .	93
General Structure of the Device Configuration Subsystem . . . . .	94
Device Configuration Database Overview . . . . .	95
Basic Device Configuration Procedures Overview . . . . .	95
Device Configuration Manager Overview . . . . .	96
Device Classes, Subclasses, and Types Overview . . . . .	97
Writing a Device Method . . . . .	98
Understanding Device Methods Interfaces . . . . .	98
Understanding Device States . . . . .	99
Adding an Unsupported Device to the System . . . . .	100
Understanding Device Dependencies and Child Devices . . . . .	101
Accessing Device Attributes . . . . .	102
Device Dependent Structure (DDS) Overview . . . . .	103
List of Device Configuration Commands . . . . .	105
List of Device Configuration Subroutines . . . . .	105
Related Information . . . . .	106
<b>Chapter 7. Communications I/O Subsystem . . . . .</b>	<b>107</b>
User-Mode Interface to a Communications PDH . . . . .	107
Kernel-Mode Interface to a Communications PDH . . . . .	107
CDLI Device Drivers . . . . .	108
Communications Physical Device Handler Model Overview . . . . .	108
Status Blocks for Communications Device Handlers Overview . . . . .	109
MPQP Device Handler Interface Overview for the ARTIC960Hx PCI Adapter . . . . .	111
Serial Optical Link Device Handler Overview . . . . .	112
Configuring the Serial Optical Link Device Driver . . . . .	113
Forum-Compliant ATM LAN Emulation Device Driver . . . . .	114
Fiber Distributed Data Interface (FDDI) Device Driver . . . . .	127
High-Performance (8fc8) Token-Ring Device Driver . . . . .	131
High-Performance (8fa2) Token-Ring Device Driver . . . . .	139
PCI Token-Ring Device Drivers . . . . .	146
Ethernet Device Drivers . . . . .	155
Related Information . . . . .	196

<b>Chapter 8. Graphic Input Devices Subsystem</b>	197
open and close Subroutines	197
read and write Subroutines	197
ioctl Subroutines	197
Input Ring	199
<b>Chapter 9. Low Function Terminal Subsystem</b>	203
Low Function Terminal Interface Functional Description	203
Components Affected by the Low Function Terminal Interface	204
Accented Characters	206
Related Information	207
<b>Chapter 10. Logical Volume Subsystem</b>	209
Direct Access Storage Devices (DASDs)	209
Physical Volumes	209
Understanding the Logical Volume Device Driver	212
Understanding Logical Volumes and Bad Blocks	215
Related Information	216
<b>Chapter 11. Printer Addition Management Subsystem</b>	219
Printer Types Currently Supported	219
Printer Types Currently Unsupported	219
Adding a New Printer Type to Your System	219
Adding a Printer Definition	220
Adding a Printer Formatter to the Printer Backend	221
Understanding Embedded References in Printer Attribute Strings	221
Related Information	221
<b>Chapter 12. Small Computer System Interface Subsystem (Parallel SCSI)</b>	223
SCSI Subsystem Overview	223
Understanding SCSI Asynchronous Event Handling	224
SCSI Error Recovery	226
A Typical Initiator-Mode SCSI Driver Transaction Sequence	229
Understanding SCSI Device Driver Internal Commands	229
Understanding the Execution of Initiator I/O Requests	230
SCSI Command Tag Queuing	232
Understanding the sc_buf Structure	232
Other SCSI Design Considerations	237
SCSI Target-Mode Overview	242
Required SCSI Adapter Device Driver ioctl Commands	247
Related Information	253
<b>Chapter 13. SCSI Architectural Model Subsystem</b>	255
Programming SAM Device Drivers	255
SAM Subsystem Overview	277
SAM Asynchronous Event Handling	283
SAM Error Recovery	285
SAM Initiator-Mode Recovery When Not Command Tag Queuing	285
Initiator-Mode Recovery During Command Tag Queuing	286
A Typical Initiator-Mode SAM Driver Transaction Sequence	287
Understanding the Execution of SAM Initiator I/O Requests	288
SAM Command Tag Queuing	289
Understanding the scsi_buf Structure	290
Other SAM Design Considerations	296
SAM Adapter Device Driver ioctl Commands	301
Related Information	301

<b>Chapter 14. Integrated Device Electronics (IDE) Subsystem</b>	303
Responsibilities of the IDE Adapter Device Driver	303
Responsibilities of the IDE Device Driver	303
Communication Between IDE Device Drivers and IDE Adapter Device Drivers	303
IDE Error Recovery	304
A Typical IDE Driver Transaction Sequence	304
IDE Device Driver Internal Commands	305
Execution of I/O Requests	305
ataide_buf Structure	306
Other IDE Design Considerations	309
Required IDE Adapter Driver ioctl Commands	310
Related Information	312
<b>Chapter 15. Serial Direct Access Storage Device Subsystem</b>	313
DASD Device Block Level Description	313
<b>Chapter 16. Debug Facilities</b>	315
System Dump Facility	315
Live Dump Facility	326
Component Trace Facility	338
Error Logging	344
Debug and Performance Tracing	347
Memory Overlay Detection System (MODS)	367
Related Information	368
<b>Chapter 17. Loadable Authentication Module Programming Interface</b>	371
Overview	371
Load Module Interfaces	371
Authentication Interfaces	372
Identification Interfaces	374
Support Interfaces	378
Configuration Files	381
Compound Load Modules	382
<b>Chapter 18. Kernel Storage-Protection Keys</b>	385
Kernel Keys and Kernel Key Sets	386
Protection Gates	387
Making a Kernel Extension Key Safe	388
Designing the Key Protection in a Key-protected Kernel Extension	388
Related Information	390
<b>Appendix. Notices</b>	391
Trademarks	392
<b>Index</b>	393



---

## About This Book

This book provides system programmers with complete information about kernel programming and the kernel environment for the AIX operating system. Programmers can use this book to gain knowledge of device drivers, kernel services, debugging tools, and kernel subsystems. Topics include the kernel environment, system calls, the kernel debug program, system dump, and virtual file systems. Each subsystem is described in detail.

---

## How to Use This Book

This book provides two types of information: (1) an overview of the kernel programming environment and information a programmer needs to write kernel extensions, and (2) information about existing kernel subsystems.

---

## Highlighting

The following highlighting conventions are used in this book:

<b>Bold</b>	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.

---

## Related Publications

The following books contain additional information on kernel extension programming and the existing kernel subsystems:

- *Printers and printing*
- *Keyboard Technical Reference*
- *Operating system and device management*
- *AIX Version 6.1 Technical Reference: Kernel and Subsystems Volume 1*
- *AIX Version 6.1 Technical Reference: Kernel and Subsystems Volume 2*



---

## Chapter 1. Kernel Environment

The kernel is dynamically extendable and can be expanded by adding routines that belong to any of the following functional classes:

- System calls
- Virtual file systems
- Kernel Extension and Device Driver Management Kernel Services
- Device Drivers

The term *kernel extension* applies to all routines added to the kernel, independent of their purpose. Kernel extensions can be added at any time by a user with the appropriate privilege.

Kernel extensions run in the same mode as the kernel. That is, when the 64-bit kernel is used, kernel extensions run in 64-bit mode. Therefore, AIX 6.1 supports 64-bit kernel extensions only.

The following kernel-environment programming information is provided to assist you in programming kernel extensions:

- “Understanding Kernel Extension Symbol Resolution”
- “Understanding Execution Environments” on page 5
- “Understanding Kernel Threads” on page 6
- “Using Kernel Processes” on page 8
- “Accessing User-Mode Data While in Kernel Mode” on page 11
- “Understanding Locking” on page 12
- “Understanding Exception Handling” on page 14
- “Using Kernel Extensions for 64-bit Processes” on page 18

A process executing in user mode can customize the kernel by using the **sysconfig** subroutine, if the process has appropriate privilege. In this way, a user-mode process can load, unload, initialize, or terminate kernel routines. Kernel configuration can also be altered by changing tunable system parameters.

Kernel extensions can also customize the kernel by using kernel services to load, unload, initialize, and terminate dynamically loaded kernel routines; to create and initialize kernel processes; and to define interrupt handlers.

**Note:** Private kernel routines (or kernel services) execute in a privileged protection domain and can affect the operation and integrity of the whole system. See “Kernel Protection Domain” on page 23 for more information.

---

### Understanding Kernel Extension Symbol Resolution

The following information is provided to assist you in understanding kernel extension symbol resolution:

- “Exporting Kernel Services and System Calls” on page 2
- “Using Kernel Services” on page 2
- “Using System Calls with Kernel Extensions” on page 2
- “Using Private Routines” on page 3
- “Using Libraries” on page 4

## Exporting Kernel Services and System Calls

A kernel extension provides additional kernel services and system calls by specifying an export file when it is link-edited. An export file contains a list of symbols to be added to the kernel name space. In addition, symbols can be identified as system calls for 32-bit processes, 64-bit processes, or both.

In an export file, symbols are listed one per line. These system calls are available to both 32- and 64-bit processes. System calls are identified by using one of the **syscall32**, **syscall64** or **syscall3264** keywords after the symbol name. Use **syscall32** to make a system call available to 32-bit processes, **syscall64** to make a system call available to 64-bit processes, and **syscall3264** to make a system call available to both 32- and 64-bit processes. For more information about export files, see **ld** Command in *AIX Version 6.1 Commands Reference, Volume 3*.

When a new kernel extension is loaded by the **sysconfig** or **kmod\_load** subroutine, any symbols exported by the kernel extension are added to the kernel name space, and are available to all subsequently loaded kernel extensions. Similarly, system calls exported by a kernel extension are available to all user programs or shared objects subsequently loaded.

## Using Kernel Services

The kernel provides a set of base kernel services to be used by kernel extensions. Kernel extensions can export new kernel services, which can then be used by subsequently loaded kernel extensions. Base kernel services, which are described in the services documentation, are made available to a kernel extension by specifying the **/usr/lib/kernex.imp** import file during the link-edit of the extension.

**Note:** Link-editing of a kernel extension should always be performed by using the **ld** command. Do not use the compiler to create a kernel extension.

If a kernel extension depends on kernel services provided by other kernel extensions, an additional import file must be specified when link-editing. An import file lists additional kernel services, with each service listed on its own line. An import file must contain the line **#!/unix** before any services are listed. The same file can be used both as an import file and an export file. The **#!/unix** line is ignored when a file is used as an export file. For more information on import files, see **ld command** in *AIX Version 6.1 Commands Reference, Volume 3*.

## Using System Calls with Kernel Extensions

A restricted set of system calls can be used by kernel extensions. A kernel process can use a larger set of system calls than a user process in kernel mode. “System Calls Available to Kernel Extensions” on page 32 specifies which system calls can be used by either type of process. User-mode processes in kernel mode can only use system calls that have all parameters passed by value. Kernel routines running under user-mode processes cannot directly use a system call having parameters passed by reference.

The second restriction is imposed because, when they access a caller’s data, system calls with parameters passed by reference access storage across a protection domain. The cross-domain memory services performing these cross-memory operations support kernel processes as if they, too, accessed storage across a protection domain. However, these services have no way to determine that the caller is in the same protection domain when the caller is a user-mode process in kernel mode. For more information on cross-domain memory services, see “Cross-Memory Kernel Services” on page 67.

**Note:** System calls must not be used by kernel extensions executing in the interrupt handler environment.

System calls available to kernel extensions are listed in **/usr/lib/kernex.imp**, along with other kernel services.

## Loading and Unloading Kernel Extensions

Kernel extensions can be loaded and unloaded by calling the **sysconfig** function from user applications. A kernel extension can load another kernel extension by using the **kmod\_load** kernel service, and kernel extensions can be unloaded by using the **kmod\_unload** kernel service.

**Loading Kernel Extensions:** Normally, kernel extensions that provide new system calls or kernel services only need to be loaded once. For these kernel extensions, loading should be performed by specifying `SYS_SINGLELOAD` when calling the **sysconfig** function, or `LD_SINGLELOAD` when calling the **kmod\_load** function. If the specified kernel extension is already loaded, a second copy is not loaded. Instead, a reference to the existing kernel extension is returned. The loader uses the specified pathname to determine whether a kernel extension is already loaded. If multiple pathnames refer to the same kernel extension, multiple copies can be loaded into the kernel.

If a kernel extension can support multiple instances of itself (particularly its data), it can be loaded multiple times, by specifying `SYS_KLOAD` when calling the **sysconfig** function, or by not specifying `LD_SINGLELOAD` when calling the **kmod\_load** function. Either of these operations loads a new copy of the kernel extension, even when one or more copies are already loaded. When this operation is used, currently loaded routines bound to the old copy of the kernel extension continue to use the old copy. Subsequently loaded routines use the most recently loaded copy of the kernel extension.

**Unloading Kernel Extensions:** Kernel extensions can be unloaded. For each kernel extension, the loader maintains a use count and a load count. The use count indicates how many other object files have referenced some exported symbol provided by the kernel extension. The load count indicates how many explicit load requests have been made for each kernel extension.

When an explicit unload of a kernel extension is requested, the load count is decremented. If the load count and the use count are both equal to 0, the kernel extension is unloaded, and the memory used by the text and data of the kernel extension is freed.

If either the load count or use count is not equal to 0, the kernel extension is not unloaded. As processes exit or other kernel extensions are unloaded, the use counts for referenced kernel extensions are decremented. Even if the load and use counts become 0, the kernel extension may not be unloaded immediately. In this case, the kernel extension's exported symbols are still available for load-time binding unless another kernel extension is unloaded or the **slibclean** command is executed. At this time, the loader unloads all modules that have both load and use counts of 0.

## Using Private Routines

So far, symbol resolution for kernel extensions has been concerned with importing and exporting symbols *from* and *to* the kernel name space. Exported symbols are global in the kernel, and can be referenced by any subsequently loaded kernel extension.

Kernel extensions can also consist of several separately link-edited modules. This is particularly useful for device drivers, where a kernel extension contains the top (pageable) half of the driver and a dependent module contains the bottom (pinned) half of the driver. Using a dependent module also makes sense when several kernel extensions use common routines. In both cases, the symbols exported by the dependent modules are not added to the global kernel name space. Instead, these symbols are only available to the kernel extension being loaded.

When link-editing a kernel extension that depends on another module, an import file should be specified listing the symbols exported by the dependent module. Before any symbols are listed, the import file should contain one of the following lines:

```
#! path/file
```

or

```
#! path/file(member)
```

**Note:** This import file can also be used as an export file when building the dependent module. Dependent modules can be found in an archive file. In this case, the member name must be specified in the `#!` line.

While a kernel extension is being loaded, any dependent modules are only loaded a single time. This allows modules to depend on each other in a complicated way, without causing multiple instances of a module to be loaded.

**Note:** The loader uses the pathname of a module to determine whether it has already been loaded. Another copy of the module can be loaded if different path names are used for the same module.

The symbols exported by dependent modules are not added to the kernel name space. These symbols can only be used by a kernel extension and its other dependent modules. If another kernel extension is loaded that uses the same dependent modules, these dependent modules will be loaded a second time.

## Using Libraries

The operating system provides the following two libraries that can be used by kernel extensions:

- **libcsys.a**
- **libsys.a**

### libcsys Library

The **libcsys.a** library contains a subset of subroutines found in the user-mode **libc.a** library that can be used by kernel extensions. When using any of these routines, the header file **/usr/include/sys/libcsys.h** should be included to obtain function prototypes, instead of the application header files, such as **/usr/include/string.h** or **/usr/include/stdio.h**. The following routines are included in **libcsys.a**:

- **atoi**
- **bcmp**
- **bcopy**
- **bzero**
- **memccpy**
- **memchr**
- **memcmp**
- **memcpy**
- **memmove**
- **memset**
- **ovbcopy**
- **strcat**
- **strchr**
- **strcmp**
- **strcpy**
- **strcspn**
- **strlen**
- **strncat**
- **strncmp**
- **strncpy**
- **strpbrk**
- **strrchr**
- **strspn**
- **strstr**

- **strtok**

**Note:** In addition to these explicit subroutines, some additional functions are defined in **libcsys.a**. All kernel extensions should be linked with **libcsys.a** by specifying **-lcsys** at link-edit time. The library **libc.a** is intended for user-level code only. Do not link-edit kernel extensions with the **-lc** flag.

## libsys Library

The **libsys.a** library provides the following set of kernel services:

- **d\_align**
- **d\_roundup**
- **timeout**
- **timeoutcf**
- **untimeout**

When using these services, specify the **-lsys** flag at link-edit time.

## User-provided Libraries

To simplify the development of kernel extensions, you can choose to split a kernel extension into separately loadable modules. These modules can be used when linking kernel extensions in the same way that they are used when developing user-level applications and shared objects. In particular, a kernel module can be created as a shared object by linking with the **-bM:SRE** flag. The shared object can then be used as an input file when linking a kernel extension. In addition, shared objects can be put into an archive file, and the archive file can be listed on the command line when linking a kernel extension. In both cases, the shared object will be loaded as a dependent module when the kernel extension is loaded.

---

## Understanding Execution Environments

There are two major environments under which a kernel extension can run:

- Process environment
- Interrupt environment

A kernel extension runs in the *process environment* when invoked either by a user process in kernel mode or by a kernel process. A kernel extension is executing in the *interrupt environment* when invoked as part of an interrupt handler.

A kernel extension can determine in which environment it is called to run by calling the **getpid** or **thread\_self** kernel service. These services respectively return the process or thread identifier of the current process or thread, or a value of -1 if called in the interrupt environment. Some kernel services can be called in both environments, whereas others can only be called in the process environment.

**Note:** No floating-point functions can be used in the kernel.

## Process Environment

A routine runs in the process environment when it is called by a user-mode process or by a kernel process. Routines running in the process environment are executed at an interrupt priority of INTBASE (the least favored priority). A kernel extension running in this environment can cause page faults by accessing pageable code or data. It can also be replaced by another process of equal or higher process priority.

A routine running in the process environment can sleep or be interrupted by routines executing in the interrupt environment. A kernel routine that runs on behalf of a user-mode process can only invoke system calls that have no parameters passed by reference. A kernel process, however, can use all system calls listed in the System Calls Available to Kernel Extensions if necessary.

## Interrupt Environment

A routine runs in the interrupt environment when called on behalf of an interrupt handler. A kernel routine executing in this environment cannot request data that has been paged out of memory and therefore cannot cause page faults by accessing pageable code or data. In addition, the kernel routine has a stack of limited size, is not subject to replacement by another process, and cannot perform any function that would cause it to sleep.

A routine in this environment is only interruptible either by interrupts that have priority more favored than the current priority or by exceptions. These routines cannot use system calls and can use only kernel services available in both the process and interrupt environments.

A process in kernel mode can also put *itself* into an environment similar to the interrupt environment. This action, occurring when the interrupt priority is changed to a priority more favored than INTBASE, can be accomplished by calling the **i\_disable** or **disable\_lock** kernel service. A kernel-mode process is sometimes required to do this to serialize access to a resource shared by a routine executing in the interrupt environment. When this is the case, the process operates under most of the same restrictions as a routine executing in the interrupt environment. However, the **e\_sleep**, **e\_wait**, **e\_sleepl**, **et\_wait**, **lockl**, and **unlockl** process can sleep, wait, and use locking kernel services if the event word or lock word is pinned.

Routines executed in this environment can adversely affect system real-time performance and are therefore limited to a specific maximum path length. Guidelines for the maximum path length are determined by the interrupt priority at which the routines are executed. Understanding Interrupts provides more information.

---

## Understanding Kernel Threads

A *thread* is an independent flow of control that operates within the same address space as other independent flows of control within a process.

One process can have multiple threads, with each thread executing different code concurrently, while sharing data and synchronizing much more easily than cooperating processes. Threads require fewer system resources than processes, and can start more quickly.

Although threads can be scheduled, they exist in the context of their process. The following list indicates what is managed at process level and shared among all threads within a process:

- Address space
- System resources, like files or terminals
- Signal list of actions.

The process remains the swappable entity. Only a few resources are managed at thread level, as indicated in the following list:

- State
- Stack
- Signal masks.

## Kernel Threads, Kernel Only Threads, and User Threads

There are three kinds of threads:

- Kernel threads
- Kernel-only threads
- 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 in user mode environment when executing user functions or library calls; it switches to kernel mode environment when executing system calls.

A *kernel-only thread* is a kernel thread that executes only in kernel mode environment. Kernel-only threads are controlled by the kernel mode environment programmer through kernel services.

User mode programs can access *user threads* through a library (such as the **libpthreads.a** threads library). User threads are part of a portable programming model. User threads are mapped to kernel threads by the threads library, in an implementation dependent manner. The threads library uses a proprietary interface to handle kernel threads. See *Understanding Threads in AIX Version 6.1 General Programming Concepts: Writing and Debugging Programs* to get detailed information about the user threads library and their implementation.

All threads discussed in this article are kernel threads; and the information applies only to the kernel mode environment. Kernel threads cannot be accessed from the user mode environment, except through the threads library.

## Kernel Data Structures

The kernel maintains thread- and process-related information in two types of structures:

- The **user** structure contains process-related information
- The **uthread** structure contains thread-related information.

These structures cannot be accessed directly by kernel extensions and device drivers. They are encapsulated for portability reasons. Many fields that were previously in the **user** structure are now in the **uthread** structure.

## Thread Creation, Execution, and Termination

A process is always created with one thread, called the *initial thread*. The initial thread provides compatibility with previous single-threaded processes. The initial thread's stack is the process stack. See "Kernel Process Creation, Execution, and Termination" on page 10 to get more information about kernel process creation.

Other threads can be created, using a two-step procedure. The **thread\_create** kernel service allocates and initializes a new thread, and sets its state to idle. The **kthread\_start** kernel service then starts the thread, using the specified entry point routine.

A thread is terminated when it executes a return from its entry point, or when it calls the **thread\_terminate** kernel service. Its resources are automatically freed. If it is the last thread in the process, the process ends.

## Thread Scheduling

Threads are scheduled using one of the following scheduling policies:

- First-in first-out (FIFO) scheduling policy, with fixed priority. Using the FIFO policy with high favored priorities might lead to bad system performance.
- Round-robin (RR) scheduling policy, quantum based and with fixed priority.
- Default scheduling policy, a non-quantum based round-robin scheduling with fluctuating priority. Priority is modified according to the CPU usage of the thread.

Scheduling parameters can be changed using the **thread\_setsched** kernel service. The process-oriented **setpri** system call sets the priority of all the threads within a process. The process-oriented **getpri** system

call gets the priority of a thread in the process. The scheduling policy and priority of an individual thread can be retrieved from the `ti_policy` and `ti_pri` fields of the `thrdsinfo` structure returned by the `getthrds` system call.

## Thread Signal Handling

The signal handling concepts are the following:

- A signal mask is associated with each thread.
- The list of actions associated with each signal number is shared among all threads in the process.
- If the signal action specifies termination, stop, or continue, the entire process, thus including all its threads, is respectively terminated, stopped, or continued.
- Synchronous signals attributable to a particular thread (such as a hardware fault) are delivered to the thread that caused the signal to be generated.
- Signals can be directed to a particular thread. If the target thread has blocked the signal from delivery, the signal remains pending on the thread until the thread unblocks the signal from delivery, or the action associated with the signal is set to ignore by any thread within the process.

The signal mask of a thread is handled by the `limit_sigs` and `sigsetmask` kernel services. The `kthread_kill` kernel service can be used to direct a signal to a particular thread.

In the kernel environment, when a signal is received, no action is taken (no termination or handler invocation), even for the **SIGKILL** signal. A thread in kernel environment, especially kernel-only threads, must *poll* for signals so that signals can be delivered. Polling ensures the proper kernel-mode serialization. For example, **SIGKILL** will not be delivered to a kernel-only thread that does not poll for signals. Therefore, **SIGKILL** is not necessarily an effective means for terminating a kernel-only thread.

Signals whose actions are applied at generation time (rather than delivery time) have the same effect regardless of whether the target is in kernel or user mode. A kernel-only thread can poll for unmasked signals that are waiting to be delivered by calling the `sig_chk` kernel service. This service returns the signal number of a pending signal that was not blocked or ignored. The thread then uses the signal number to determine which action should be taken. The kernel does not automatically call signal handlers for a thread in kernel mode as it does for user mode.

See “Kernel Process Signal and Exception Handling” on page 11 for more information about signal handling at process level.

---

## Using Kernel Processes

A kernel process is a process that is created in the kernel protection domain and always executes in the kernel protection domain. Kernel processes can be used in subsystems, by complex device drivers, and by system calls. They can also be used by interrupt handlers to perform asynchronous processing not available in the interrupt environment. Kernel processes can also be used as device managers where asynchronous input/output (I/O) and device management is required.

## Introduction to Kernel Processes

A kernel process (`kproc`) exists only in the kernel protection domain and differs from a user process in the following ways:

- It is created using the `creatp` and `initp` kernel services.
- It executes only within the kernel protection domain and has all security privileges.
- It can call a restricted set of system calls and all applicable kernel services. For more information, see “System Calls Available to Kernel Extensions” on page 32.
- It has access to the global kernel address space (including the kernel pinned and pageable heaps), kernel code, and static data areas.

- It must poll for signals and can choose to ignore any signal delivered, including a **kill** signal.
- Its text and data areas come from the global kernel heap.
- It cannot use application libraries.
- It has a process-private region containing only the **u-block** (user block) structure and possibly the kernel stack.
- Its *parent process* is the process that issued the **creatp** kernel service to create the process.
- It can change its parent process to the **init** process and can use interrupt disable functions for serialization.
- It can use locking to serialize process-time access to critical data structures.
- It can only be a 64-bit process in the 64-bit kernel.

A kernel process controls directly the kernel threads. Because kernel processes are always in the kernel protection domain, threads within a kernel process are kernel-only threads. For more information on kernel threads, see “Understanding Kernel Threads” on page 6.

A kernel process inherits the environment of its parent process (the one calling the **creatp** kernel service to create it), but with some exceptions. The kernel process does not have a root directory or a current directory when initialized. All uses of the file system functions must specify absolute path names.

Kernel processes created during phase 1 of system boot must not keep any long-term opens on files until phase 2 of system boot or run time has been reached. This is because Base Operating System changes root file systems between phase 1 and phase 2 of system boot. As a result, the system crashes if any files are open at root file system transition time.

## Accessing Data from a Kernel Process

Because kernel processes execute in the more privileged kernel protection domain, a kernel process can access data that user processes cannot. This applies to all kernel data, of which there are three general categories:

- The **user block** data structure

The **u-block** (or **u-area**) structure exists for kernel processes and contains roughly the same information for kernel processes as for user-mode processes. A kernel process must use kernel services to query or manipulate data from the **u-area** to maintain modularity and increase portability of code to other platforms.

- The stack for a kernel process

To ensure binary compatibility with older applications, each kernel process has a stack called the *process stack*. This stack is used by the process initial thread.

The location of the stack for a kernel process is implementation-dependent. This stack can be located in global memory or in the process-private segment of the kernel process. A kernel process must not assume automatically that its stack is located in global memory.

- Global kernel memory

A kernel process can also access global kernel memory as well as allocate and de-allocate memory from the kernel heaps. Because it runs in the kernel protection domain, a kernel process can access any valid memory location within the global kernel address space. Memory dynamically allocated from the kernel heaps by the kernel process must be freed by the kernel process before exiting. Unlike user-mode processes, memory that is dynamically allocated by a kernel process is not freed automatically upon process exit.

## Cross-Memory Services

Kernel processes must be provided with a valid cross-memory descriptor to access address regions outside the kernel global address space or kernel process address space. For example, if a kernel process is to access data from a user-mode process, the system call using the process must obtain a

cross-memory descriptor for the user-mode region to be accessed. Calling the **xmattach** or **xmattach64** kernel service provides a descriptor that can then be made available to the kernel process.

The kernel process should then call the **xmemin** and **xmemout** kernel services to access the targeted cross-memory data area. When the kernel process has completed its operation on the memory area, the cross-memory descriptor must be detached by using the **xmdetach** kernel service.

## Kernel Process Creation, Execution, and Termination

A kernel process is created by a kernel-mode routine by calling the **creatp** kernel service. This service allocates and initializes a process block for the process and sets the new process state to idle. This new kernel process does not run until it is initialized by the **initp** kernel service, which must be called in the same process that created the new kernel process (with the **creatp** service). The **creatp** kernel service returns the process identifier for the new kernel process.

The process is created with one kernel-only thread, called the *initial thread*. See Understanding Kernel Threads to get more information about threads.

After the **initp** kernel service has completed the process initialization, the initial thread is placed on the run queue. On the first dispatch of the newly initialized kernel process, it begins execution at the entry point previously supplied to the **initp** kernel service. The initialization parameters were previously specified in the call to the **initp** kernel service.

A kernel process terminates when it executes a return from its main entry routine. A process should never exit without both freeing all dynamically allocated storage and releasing all locks owned by the kernel process.

When kernel processes exit, the parent process (the one calling the **creatp** and **initp** kernel services to create the kernel process) receives the **SIGCHLD** signal, which indicates the end of a child process. However, it is sometimes undesirable for the parent process to receive the **SIGCHLD** signal due to ending a process. In this case, the kproc can call the **setpinit** kernel service to designate the **init** process as its parent. The **init** process cleans up the state of all its child processes that have become zombie processes. A kernel process can also issue the **setsid** subroutine call to change its session. Signals and job control affecting the parent process session do not affect the kernel process.

## Kernel Process Preemption

A kernel process is initially created with the same process priority as its parent. It can therefore be replaced by a more favored kernel or user process. It does not have higher priority just because it is a kernel process. Kernel processes can use the **setpri** subroutine to modify their execution priority.

The kernel process can use the locking kernel services to serialize access to critical data structures. This use of locks does not guarantee that the process will not be replaced, but it does ensure that another process trying to acquire the lock waits until the kernel process owning the lock has released it.

Using locks, however, does not provide serialization if a kernel routine can access the critical data while executing in the interrupt environment. Serialization with interrupt handlers must be handled by using locking together with interrupt control. The **disable\_lock** and **unlock\_enable** kernel services should be used to serialize with interrupt handlers.

Kernel processes must ensure that their maximum path lengths adhere to the specifications for interrupt handlers when executing at an interrupt priority more favored than INTBASE. This ensures that system real-time performance is not degraded.

## Kernel Process Signal and Exception Handling

Signals are delivered to exactly one thread within the process which has not blocked the signal from delivery. If all threads within the target process have blocked the signal from delivery, the signal remains pending on the process until a thread unblocks the signal from delivery, or the action associated with the signal is set to ignore by any thread within the process. See “Thread Signal Handling” on page 8 for more information on signal handling by threads.

Signals whose action is applied at generation time (rather than delivery time) have the same effect regardless of whether the target is a kernel or user process. A kernel process can poll for unmasked signals that are waiting to be delivered by calling the **sig\_chk** kernel service. This service returns the signal number of a pending signal that was not blocked or ignored. The kernel process then uses the signal number to determine which action should be taken. The kernel does not automatically call signal handlers for a kernel process as it does for user processes.

A kernel process should also use the exception-catching facilities (**setjmpx**, and **clrjmpx**) available in kernel mode to handle exceptions that can be caused during run time of the kernel process. Exceptions received during the execution of a kernel process are handled the same as exceptions that occur in any kernel-mode routine.

Unhandled exceptions that occur in kernel mode (in any user process while in kernel mode, in an interrupt handler, or in a kernel process) result in a system crash. To avoid crashing the system due to unhandled exceptions, kernel routines should use the **setjmpx**, **clrjmpx**, and **longjmpx** kernel services to handle exceptions that might possibly occur during run time. See “Understanding Exception Handling” on page 14 for more details on handling exceptions.

## Kernel Process Use of System Calls

System calls made by kernel processes do not result in a change of protection domain because the kernel process is already within the kernel protection domain. Routines in the kernel (including routines executing in a kernel process) are bound by the loader to the system call function and not to the system call handler. When system calls use kernel services to access user-mode data, these kernel services recognize that the system call is running within a kernel process instead of a user process and correctly handle the data accesses.

However, the error information returned from a kernel process system call must be accessed differently than for a user process. A kernel process must use the **getuerror** kernel service to retrieve the system call error information normally provided in the **errno** global variable for user-mode processes. In addition, the kernel process can use the **setuerror** kernel service to set the error information to 0 before calling the system call. The return code from the system call is handled the same for all processes.

Kernel processes can use only a restricted set of the base system calls. “System Calls Available to Kernel Extensions” on page 32 lists system calls available to kernel processes.

---

## Accessing User-Mode Data While in Kernel Mode

Kernel extensions must use a set of kernel services to access data that is in the user-mode protection domain. These services ensure that the caller has the authority to perform the desired operation at the time of data access and also prevent system crashes in a system call when accessing user-mode data. These services can be called only when running in the process environment of the process that contains the user-mode data. For more information on user-mode protection, see “User Protection Domain” on page 23. For more information on the process environment, see “Process Environment” on page 5.

## Data Transfer Services

The following list shows user-mode data access kernel services (primitives):

Kernel Service	Purpose
<b>suword</b>	Stores a word of data in user memory.
<b>fubyte</b>	Fetches, or retrieves, a byte of data from user memory.
<b>fuword</b>	Fetches, or retrieves, a word of data from user memory.
<b>copyin</b>	Copies data between user and kernel memory.
<b>copyout</b>	Copies data between user and kernel memory.
<b>copyinstr</b>	Copies a character string (including the terminating null character) from user to kernel space.

Additional kernel services allow data transfer between user mode and kernel mode when a **uio** structure is used, thereby describing the user-mode data area to be accessed. Following is a list of services that typically are used between the file system and device drivers to perform device I/O:

Kernel Service	Purpose
<b>uioMOVE</b>	Moves a block of data between kernel space and a space defined by a <b>uio</b> structure.
<b>ureadc</b>	Writes a character to a buffer described by a <b>uio</b> structure.
<b>uwritec</b>	Retrieves a character from a buffer described by a <b>uio</b> structure.

Kernel services ending in “64”, such as **suword64**, **copyin64** and so on, are deprecated in AIX 6.1. To maintain binary compatibility of applications, macros in the **sys/uio.h** header file redefine these services to their counterparts when compiling in 64-bit mode.

## Using Cross-Memory Kernel Services

Occasionally, access to user-mode data is required when not in the environment of the user-mode process that has addressability to the data. Such cases occur when the data is to be accessed asynchronously. Examples of asynchronous accessing include:

- Direct memory access to the user data by I/O devices
- Data access by interrupt handlers
- Data access by a kernel process

In these circumstances, the kernel cross-memory services are required to provide the necessary access. The **xmattach** kernel service allows a cross-memory descriptor to be obtained for the data area to be accessed. These services must be called in the process environment of the process containing the data area.

After a cross-memory descriptor has been obtained, the **xmemin** and **xmemout** kernel services can be used to access the data area outside the process environment containing the data. When access to the data area is no longer required, the access must be removed by calling the **xmdetach** kernel service. Kernel extensions should use these services only when absolutely necessary. Because of the machine dependencies of cross-memory operations, using them increases the difficulty of porting the kernel extension to other machine platforms.

---

## Understanding Locking

The following information is provided to assist you in understanding locking.

### Lockl Locks

The *lockl locks* (previously called *conventional locks*) are provided for compatibility only and should not be used in new code: simple or complex locks should be used instead. These locks are used to protect a

critical section of code which accesses a resource such as a data structure or device, serializing access to the resource. Every thread which accesses the resource must acquire the lock first, and release the lock when finished.

A conventional lock has two states: locked or unlocked. In the *locked* state, a thread is currently executing code in the critical section, and accessing the resource associated with the conventional lock. The thread is considered to be the owner of the conventional lock. No other thread can lock the conventional lock (and therefore enter the critical section) until the owner unlocks it; any thread attempting to do so must wait until the lock is free. In the *unlocked* state, there are no threads accessing the resource or owning the conventional lock.

*lockl* locks are recursive and, unlike simple and complex locks, can be awakened by a signal.

## Simple Locks

A *simple* lock provides exclusive-write access to a resource such as a data structure or device. Simple locks are not recursive and have only two states: locked or unlocked.

## Complex Locks

A *complex* lock can provide either shared or exclusive access to a resource such as a data structure or device. Complex locks are not recursive by default (but can be made recursive) and have three states: exclusive-write, shared-read, or unlocked.

If several threads perform read operations on the resource, they must first acquire the corresponding lock in shared-read mode. Because no threads are updating the resource, it is safe for all to read it. Any thread which writes to the resource must first acquire the lock in exclusive-write mode. This guarantees that no other thread will read or write the resource while it is being updated.

## Types of Critical Sections

There are two types of critical sections which must be protected from concurrent execution in order to serialize access to a resource:

<b>thread-thread</b>	These critical sections must be protected (by using the locking kernel services) from concurrent execution by multiple processes or threads.
<b>thread-interrupt</b>	These critical sections must be protected (by using the <b>disable_lock</b> and <b>unlock_enable</b> kernel services) from concurrent execution by an interrupt handler and a thread or process.

## Priority Promotion

When a lower priority thread owns a lock which a higher-priority thread is attempting to acquire, the owner has its priority promoted to that of the most favored thread waiting for the lock. When the owner releases the lock, its priority is restored to its normal value. Priority promotion ensures that the lock owner can run and release its lock, so that higher priority processes or threads do not remain blocked on the lock.

## Locking Strategy in Kernel Mode

**Attention:** A kernel extension should not attempt to acquire the kernel lock if it owns any other lock. Doing so can cause unpredictable results or system failure.

A hierarchy of locks exists. This hierarchy is imposed by software convention, but is not enforced by the system. The lock **kernel\_lock** variable, which is the global kernel lock, has the the coarsest granularity. Other types of locks have finer granularity. The following list shows the ordering of locks based on granularity:

- The **kernel\_lock** global kernel lock

**Note:** Avoid using the **kernel\_lock** global kernel lock variable in new code. It is only included for compatibility purposes.

- File system locks (private to file systems)
- Device driver locks (private to device drivers)
- Private fine-granularity locks

Locks should generally be released in the reverse order from which they were acquired; all locks must be released before a kernel process or thread exits. Kernel mode processes do not receive any signals while they hold any lock.

---

## Understanding Exception Handling

Exception handling involves a basic distinction between *interrupts* and *exceptions*:

- An interrupt is an asynchronous event and is not associated with the instruction that is executing when the interrupt occurs.
- An exception is a synchronous event and is directly caused by the instruction that is executing when the exception occurs.

The computer hardware generally uses the same mechanism to report both interrupts and exceptions. The machine saves and modifies some of the event's state and forces a branch to a particular location. When decoding the reason for the machine interrupt, the interrupt handler determines whether the event is an interrupt or an exception, then processes the event accordingly.

## Exception Processing

When an exception occurs, the current instruction stream cannot continue. If you ignore the exception, the results of executing the instruction may become undefined. Further execution of the program may cause unpredictable results. The kernel provides a default exception-handling mechanism by which an instruction stream (a process- or interrupt-level program) can specify what action is to be taken when an exception occurs. Exceptions are handled differently depending on whether they occurred while executing in kernel mode or user mode.

### Default Exception-Handling Mechanism

If no exception handler is currently defined when an exception occurs, typically one of two things happens:

- If the exception occurs while a process is executing in user mode, the process is sent a signal relevant to the type of exception.
- If the exception occurs while in kernel mode, the system halts.

## Kernel-Mode Exception Handling

Exception handling in kernel mode extends the **setjump** and **longjump** subroutines context-save-and-restore mechanism by providing **setjmpx** and **longjmpx** kernel services to handle exceptions. The traditional system mechanism is extended by allowing these exception handlers (or context-save checkpoints) to be stacked on a per-process or per-interrupt handler basis.

This stacking mechanism allows the execution point and context of a process or interrupt handler to be restored to a point in the process or interrupt handler, *at the point of return from the **setjmpx** kernel service*. When execution returns to this point, the return code from **setjmpx** kernel service indicates the type of exception that occurred so that the process or interrupt handler state can be fully restored. Appropriate retry or recovery operations are then invoked by the software performing the operation.

When an exception occurs, the kernel first-level exception handler gets control. The first-level exception handler determines what type of exception has occurred and saves information necessary for handling the specific type of exception. For an I/O exception, the first-level handler also enables again the programmed I/O operations.



The first-level exception handler then modifies the saved context of the interrupted process or interrupt handler. It does so to execute the **longjmpx** kernel service when the first-level exception handler returns to the interrupted process or interrupt handler.

The **longjmpx** kernel service executes in the environment of the code that caused the exception and restores the current context from the topmost jump buffer on the stack of saved contexts. As a result, the state of the process or interrupt handler that caused the exception is restored to the point of the return from the **setjmpx** kernel service. (The return code, nevertheless, indicates that an exception has occurred.)

The process or interrupt handler software should then check the return code and invoke exception handling code to restore fully the state of the process or interrupt handler. Additional information about the exception can be obtained by using the **getexcept** kernel service.

## User-Defined Exception Handling

A typical exception handler should do the following:

- Perform any necessary clean-up such as freeing storage or segment registers and releasing other resources.
- If the exception is recognized by the current handler and can be handled entirely within the routine, the handler should establish itself again by calling the **setjmpx** kernel service. This allows normal processing to continue.
- If the exception is not recognized by the current handler, it must be passed to the previously stacked exception handler. The exception is passed by calling the **longjmpx** kernel service, which either calls the previous handler (if any) or takes the system's default exception-handling mechanism.
- If the exception is recognized by the current handler but cannot be handled, it is treated as though it is unrecognized. The **longjmpx** kernel service is called, which either passes the exception along to the previous handler (if any) or takes the system default exception-handling mechanism.

When a kernel routine that has established an exception handler completes normally, it must remove its exception handler from the stack (by using the **clrjmpx** kernel service) before returning to its caller.

**Note:** When the **longjmpx** kernel service invokes an exception handler, that handler's entry is automatically removed from the stack.

## Implementing Kernel Exception Handlers

The following information is provided to assist you in implementing kernel exception handlers.

### **setjmpx, longjmpx, and clrjmpx Kernel Services**

The **setjmpx** kernel service provides a way to save the following portions of the program state at the point of a call:

- Nonvolatile general registers
- Stack pointer
- TOC pointer
- Interrupt priority number (**intpri**)
- Ownership of kernel-mode lock

This state can be restored later by calling the **longjmpx** kernel service, which accomplishes the following tasks:

- Reloads the registers (including TOC and stack pointers)
- Enables or disables to the correct interrupt level
- Conditionally acquires or releases the kernel-mode lock
- Forces a branch back to the point of original return from the **setjmpx** kernel service

The **setjmpx** kernel service takes the address of a jump buffer (a **label\_t** structure) as an explicit parameter. This structure can be defined anywhere including on the stack (as an automatic variable). After noting the state data in the jump buffer, the **setjmpx** kernel service pushes the buffer onto the top of a stack that is maintained in the machine-state save structure.

The **longjmpx** kernel service is used to return to the point in the code at which the **setjmpx** kernel service was called. Specifically, the **longjmpx** kernel service returns to the most recently created jump buffer, as indicated by the top of the stack anchored in the machine-state save structure.

The parameter to the **longjmpx** kernel service is an exception code that is passed to the resumed program as the return code from the **setjmp** kernel service. The resumed program tests this code to determine the conditions under which the **setjmpx** kernel service is returning. If the **setjmpx** kernel service has just saved its jump buffer, the return code is 0. If an exception *has* occurred, the program is entered by a call to the **longjmpx** kernel service, which passes along a return code that is *not* equal to 0.

**Note:** Only the resources listed here are saved by the **setjmpx** kernel service and restored by the **longjmpx** kernel service. Other resources, in particular segment registers, are not restored. A call to the **longjmpx** kernel service, by definition, returns to an earlier point in the program. The program code must free any resources that are allocated between the call to the **setjmpx** kernel service and the call to the **longjmpx** kernel service.

If the exception handler stack is empty when the **longjmpx** kernel service is issued, there is no place to jump to and the system default exception-handling mechanism is used. If the stack is not empty, the context that is defined by the topmost jump buffer is reloaded and resumed. The topmost buffer is then removed from the stack.

The **clrjmpx** kernel service removes the top element from the stack as placed there by the **setjmpx** kernel service. The caller to the **clrjmpx** kernel service is expected to know exactly which jump buffer is being removed. This should have been established earlier in the code by a call to the **setjmpx** kernel service. Accordingly, the address of the buffer is required as a parameter to the **clrjmpx** kernel service. It can then perform consistency checking by asserting that the address passed is indeed the address of the top stack element.

## Exception Handler Environment

The stacked exception handlers run in the environment in which the exception occurs. That is, an exception occurring in a process environment causes the next dispatch of the process to run the exception handler on the top of the stack of exception handlers for that process. An exception occurring in an interrupt handler causes the interrupt handler to return to the context saved by the last call to the **setjmpx** kernel service made by the interrupt handler.

**Note:** An interrupt handler context is newly created each time the interrupt handler is invoked. As a result, exception handlers for interrupt handlers must be registered (by calling the **setjmpx** kernel service) each time the interrupt handler is invoked. Otherwise, an exception detected during execution of the interrupt handler will be handled by the default handler.

## Restrictions on Using the setjmpx Kernel Service

Process and interrupt handler routines registering exception handlers with the **setjmpx** kernel service must not return to their caller before removing the saved jump buffer or buffers from the list of jump buffers. A saved jump buffer can be removed by invoking the **clrjmpx** kernel service in the reverse order of the **setjmpx** calls. The saved jump buffer must be removed before return because the routine's context no longer exists once the routine has returned to its caller.

If, on the other hand, an exception does occur (that is, the return code from **setjmpx** kernel service is nonzero), the jump buffer is automatically removed from the list of jump buffers. In this case, a call to the **clrjmpx** kernel service for the jump buffer must not be performed.

Care must also be taken in defining variables that are used after the context save (the call to the **setjmpx** service), and then again by the exception handler. Sensitive variables of this nature must be restored to their correct value by the exception handler when an exception occurs.

**Note:** If the last value of the variable is desired at exception time, the variable data type must be declared as "volatile."

Exception handling is concluded in one of two ways. Either a registered exception handler handles the exception and continues from the saved context, or the default exception handler is reached by exhausting the stack of jump buffers.

## Exception Codes

The `/usr/include/sys/except.h` file contains a list of code numbers corresponding to the various types of hardware exceptions. When an exception handler is invoked (the return from the **setjmpx** kernel service is not equal to 0), it is the responsibility of the handler to test the code to ensure that the exception is one the routine can handle. If it is not an expected code, the exception handler must:

- Release any resources that would not otherwise be freed (buffers, segment registers, storage acquired using the **xmalloc** routines)
- Call the **longjmpx** kernel service, passing it the exception code as a parameter

Thus, when an exception handler does not recognize the exception for which it has been invoked, it passes the exception on to the next most recent exception handler. This continues until an exception handler is reached that recognizes the code and can handle it. Eventually, if no exception handler can handle the exception, the stack is exhausted and the system default action is taken.

In this manner, a component can allocate resources (after calling the **setjmpx** kernel service to establish an exception handler) and be assured that the resources will later be released. This ensures the exception handler gets a chance to release those resources regardless of what events occur before the instruction stream (a process- or interrupt-level code) is terminated.

By coding the exception handler to recognize what exception codes it can process rather than encoding this knowledge in the stack entries, a powerful and simple-to-use mechanism is created. Each handler need only investigate the exception code that it receives rather than just assuming that it was invoked because a particular exception has occurred to implement this scheme. The set of exception codes used cannot have duplicates.

Exceptions generated by hardware use one of the codes in the `/usr/include/sys/except.h` file. However, the **longjmpx** kernel service can be invoked by any kernel component, and any integer can serve as the exception code. A mechanism similar to the old-style **setjmp** and **longjmp** kernel services can be implemented on top of the **setjmpx/longjmpx** stack by using exception codes outside the range of those used for hardware exceptions.

To implement this old-style mechanism, a unique set of exception codes is needed. These codes must not conflict with either the pre-assigned hardware codes or codes used by any other component. A simple way to get such codes is to use the addresses of unique objects as code values.

For example, a program that establishes an exception handler might compare the exception code to the address of its own entry point. Later on in the calling sequence, after any number of intervening calls to the **setjmpx** kernel service by other programs, a program can issue a call to the **longjmpx** kernel service and pass the address of the agreed-on function descriptor as the code. This code is only recognized by a single exception handler. All the intervening ones just clean up their resources and pass the code to the **longjmpx** kernel service again.

Addresses of functions are not the only possibilities for unique code numbers. For example, addresses of external variables can also be used. By using unique, system-wide addresses, the problem of code-space

collision is transformed into a problem of external-name collision. This problem is easier to solve, and is routinely solved whenever the system is built. By comparison, pre-assigning exception numbers by using **#define** statements in a header file is a much more cumbersome and error-prone method.

## Hardware Detection of Exceptions

Each of the exception types results in a hardware interrupt. For each such interrupt, a first-level interrupt handler (FLIH) saves the state of the executing program and calls a second-level handler (SLIH). The SLIH is passed a pointer to the machine-state save structure and a code indicating the cause of the interrupt.

When a SLIH determines that a hardware interrupt should actually be considered a synchronous exception, it sets up the machine-state save to invoke the **longjmpx** kernel service, and then returns. The FLIH then resumes the instruction stream at the entry to the **longjmpx** service.

The **longjmpx** service then invokes the top exception handler on the stack or takes the system default action as previously described.

## User-Mode Exception Handling

Exceptions that occur in a user-mode process and are not automatically handled by the kernel cause the user-mode process to be signaled. If the process is in a state in which it cannot take the signal, it is terminated and the information logged. Kernel routines can install user-mode exception handlers that catch exceptions before they are signaled to the user-mode process.

The **uexadd** and **uexdel** kernel services allow system-wide user-mode exception handlers to be added and removed.

The most recently registered exception handler is the first called. If it cannot handle the exception, the next most recent handler on the list is called, and this second handler attempts to handle the exception. If this attempt fails, successive handlers are tried, until the default handler is called, which generates the signal.

Additional information about the exception can be obtained by using the **getexcept** kernel service.

---

## Using Kernel Extensions for 64-bit Processes

Kernel extensions in a 64-bit kernel run in 64-bit mode. Therefore, only 64-bit kernel extensions can run on AIX 6.1. You can program kernel extensions for both 32-bit and 64-bit applications.

System calls can be made available to 32- or 64-bit processes, selectively. If an application invokes a system call that is not exported to processes running in the current mode, the call fails.

Because only 64-bit kernel extensions can run on AIX 6.1, the interaction between kernel extensions and user address space is simplified. To examine and manipulate user address space, use kernel services such as the **as\_att64**, **as\_det64**, **as\_puth64**, **as\_seth64**, and **as\_getsrval64** kernel services.

Address space remapping is no longer necessary for kernel extensions running on AIX 6.1 and later. 64-bit pointers or 64-bit data values of the **long** type can be used directly by the kernel without manipulation because the kernel always runs in 64-bit mode. The 64-bit kernel can also handle 64-bit addresses without mapping them to 32-bit values.

**32-bit:** Do not use the **as\_att**, **as\_det**, **as\_seth**, **as\_geth**, **as\_puth**, and **as\_getsrval** kernel services; they are obsolete. The **as\_remap64**, **as\_unremap64**, **get64bitparm**, and **saveretval64** kernel services are also obsolete.

---

## 64-bit Kernel Extension Programming Environment

### C Language Data Model

The 64-bit kernel uses the LP64 (Long Pointer 64-bit) C language data model and requires kernel extensions to do the same. The LP64 data model defines **pointers**, **long**, and **long long** types as 64 bits, **int** as 32 bits, **short** as 16 bits, and **char** as 8 bits.

In the ILP32 data model, **long** and **pointer** types are 32 bits. In order to port an existing 32-bit kernel extension to the 64-bit kernel environment, source code must be modified to be type-safe under LP64. This means ensuring that data types are used in a consistent fashion. Source code is incorrect for the 64-bit environment if it assumes that pointers, **long**, and **int** are all the same size.

In addition, the use of system-derived types must be examined whenever values are passed from an application to the kernel. For example, **size\_t** is a system-derived type whose size depends on the compilation mode, and **key\_t** is a system-derived type that is 64 bits in the 64-bit kernel environment.

### Kernel Data Structures

Several global, exported kernel data structures have been changed in the 64-bit kernel, in order to support scalability and future functionality. These changes include larger structure sizes as a result of being compiled under the LP64 data model. In porting a kernel extension to the 64-bit kernel environment, these data structure changes must be considered.

### Function Prototypes

Function prototypes are more important in the 64-bit programming environment than the 32-bit programming environment, because the default return value of an undeclared function is **int**. If a function prototype is missing for a function returning a pointer, the compiler will convert the returned value to an **int** by setting the high-order word to 0, corrupting the value. In addition, function prototypes allow the compiler to do more type checking, regardless of the compilation mode.

When compiled in 64-bit mode, system header files define full function prototypes for all kernel services provided by the 64-bit kernel. By defining the **\_\_FULL\_PROTO** macro, function prototypes are provided in 32-bit mode as well. It is recommended that function prototypes be provided by including the system header files, instead of providing a prototype in a source file.

### Compiler Options

To compile a kernel extension in 64-bit mode, the **-q64** flag must be used. To check for missing function prototypes, **-qinfo=pro** can be specified. To compile in ANSI mode, use the **-qlanglvl=ansi** flag. When this flag is used, additional error checking will be performed by the compiler. To link-edit a kernel extension, the **-b64** option must be used with the **ld** command.

**Note:** Do not link kernel extensions using the **cc** command.

### Conditional Compilation

When compiling in 64-bit mode, the compiler automatically defines the macro **\_\_64BIT\_\_**. Kernel extensions should always be compiled with the **\_KERNEL** macro defined, and if **sys/types.h** is included, the macro **\_\_64BIT\_KERNEL** will be defined for kernel extensions being compiled in 64-bit mode. The **\_\_64BIT\_KERNEL** macro can be used to provide for conditional compilation when compiling kernel extensions from common source code.

Kernel extensions should not be compiled with the **\_KERNSYS** macro defined. If this macro is defined, the resulting kernel extension will not be supported, and binary compatibility will not be assured with future releases.

## Kernel Extension Libraries

The **libcsys.a** and **libsys.a** libraries are supported for kernel extensions. Function prototypes for all the functions in **libcsys.a** are found in **sys/libcsys.h**.

## Kernel Execution Mode

Within the 64-bit kernel, all kernel mode subsystems, including kernel extensions, run exclusively in 64-bit processor mode and are capable of accessing data or executing instructions at any location within the kernel's 64-bit address space, including those found above the first 4GBs of this address space. This availability of the full 64-bit address space extends to all kernel entities, including kprocs and interrupt handlers, and enables the potential for software resource scalability through the introduction of an enormous kernel address space.

## Kernel Address Space

The 64-bit kernel provides a common user and kernel 64-bit address space.

## Kernel Address Space Organization

The organization of kernel space differs between hardware systems. Therefore, kernel extensions must not have any dependencies on the locations, relative or absolute, of the kernel text, kernel global data, kernel heap data, and kernel stack values, and must appropriately type variables used to hold kernel addresses.

## Temporary Attachment

The 64-bit kernel provides kernel extensions with the capability to temporarily attach virtual memory segments to the kernel space for the current thread of kernel mode execution. This capability is provided through the **vm\_att** and **vm\_det** services.

A total of four concurrent temporary attaches will be supported under a single thread of execution.

## Global Regions

The 64-bit kernel provides kernel extensions with the capability to create global regions within the kernel address space. Once created, a region is globally accessible to all kernel code until it is destroyed. Regions may be created with unique characteristics, for example, page protection, that suit kernel extension requirements and are different from the global virtual memory allocated from the **kernel\_heap**.

Global regions are also useful for kernel extensions that in the past have organized their data around virtual memory segments and require sizes and alignments that are inappropriate for the kernel heap. Under the 64-bit kernel, this memory can be provided through global regions rather than separate virtual memory segments, thus avoiding the complexity and performance cost of temporarily attaching virtual memory segments.

Global regions are created and destroyed with the **vm\_galloc** and **vm\_gfree** kernel services.

---

## Related Information

Chapter 15, "Serial Direct Access Storage Device Subsystem," on page 313

"Locking Kernel Services" on page 60

"Handling Signals While in a System Call" on page 30

"System Calls Available to Kernel Extensions" on page 32

## Subroutine References

The **setpri** subroutine, **sysconfig** subroutine in *AIX Version 6.1 Technical Reference: Base Operating System and Extensions Volume 2*.

## Commands References

The **ar** command in *AIX Version 6.1 Commands Reference, Volume 1*.

The **ld** command in *AIX Version 6.1 Commands Reference, Volume 3*.

## Technical References

The **clrjmpx** kernel service, **copyin** kernel service, **copyinstr** kernel service, **copyout** kernel service, **creatp** kernel service, **disable\_lock** kernel service, **e\_sleep** kernel service, **e\_sleepl** kernel service, **e\_wait** kernel service, **et\_wait** kernel service, **fubyte** kernel service, **fuword** kernel service, **getexcept** kernel service, **i\_disable** kernel service, **i\_enable** kernel service, **i\_init** kernel service, **initp** kernel service, **lockl** kernel service, **longjmpx** kernel service, **setjmpx** kernel service, **setpinit** kernel service, **sig\_chk** kernel service, **subyte** kernel service, **suword** kernel service, **uiomove** kernel service, **unlockl** kernel service, **ureadc** kernel service, **uwritec** kernel service, **uexadd** kernel service, **uexdel** kernel service, **xmalloc** kernel service, **xmattach** kernel service, **xmdetach** kernel service, **xmemin** kernel service, **xmemout** kernel service in *AIX Version 6.1 Technical Reference: Kernel and Subsystems Volume 1*.

The **uio** structure in *AIX Version 6.1 Technical Reference: Kernel and Subsystems Volume 1*.





---

## Chapter 2. System Calls

A system call is a routine that allows a user application to request actions that require special privileges. Adding system calls is one of several ways to extend the functions provided by the kernel.

The distinction between a system call and an ordinary function call is only important in the kernel programming environment. User-mode application programs are not usually aware of this distinction.

Operating system functions are made available to the application program in the form of *programming libraries*. A set of library functions found in a library such as **libc.a** can have functions that perform some user-mode processing and then internally start a system call. In other cases, the system call can be directly exported by the library without any user-space code. For more information on programming libraries, see “Using Libraries” on page 4.

Operating system functions available to application programs can be split or moved between user-mode functions and kernel-mode functions as required for different releases or machine platforms. Such movement does not affect the application program. Chapter 1, “Kernel Environment,” on page 1 provides more information on how to use system calls in the kernel environment.

---

### Differences Between a System Call and a User Function

A system call differs from a user function in several key ways:

- A system call has more privilege than a normal subroutine. A system call runs with kernel-mode privilege in the kernel protection domain.
- System call code and data are located in global kernel memory.
- System call routines can create and use kernel processes to perform asynchronous processing.
- System calls cannot use shared libraries or any symbols not found in the kernel protection domain.

---

### Understanding Protection Domains

There are two protection domains in the operating system: the *user protection domain* and the *kernel mode protection domain*.

#### User Protection Domain

Application programs run in the *user protection domain*, which provides:

- Read and write access to the data region of the process
- Read access to the text and shared text regions of the process
- Access to shared data regions using the shared memory functions.

When a program is running in the user protection domain, the processor executes instructions in the problem state, and the program does not have direct access to kernel data.

#### Kernel Protection Domain

The code in the kernel and kernel extensions run in the *kernel protection domain*. This code includes interrupt handlers, kernel processes, device drivers, system calls, and file system code. The processor is in the kernel protection domain when it executes instructions in the privileged state, which provides:

- Read and write access to the global kernel address space
- Read and write access to the thread’s **uthread** block and u-block, except when an interrupt handler is running.

Code running in the kernel protection domain can affect the execution environments of all processes because it:

- Can access global system data
- Can use all kernel services
- Is exempt from all security constraints.

Programming errors in the code running in the kernel protection domain can cause the operating system to fail. In particular, a process's user data cannot be accessed directly, but must be accessed using the **copyin** and **copyout** kernel services, or their variants. These routines protect the kernel from improperly supplied user data addresses.

Application programs can gain controlled access to kernel data by making system calls. Access to functions that directly or indirectly invoke system calls is typically provided by programming libraries, providing access to operating system functions.

---

## Understanding System Call Execution

When a user program invokes a system call, a system call instruction is executed, which causes the processor to begin executing the system call handler in the kernel protection domain. This system call handler performs the following actions:

1. Sets the `ut_error` field in the **uthread** structure to 0
2. Switches to a kernel stack associated with the calling thread
3. Calls the function that implements the requested system call.

The system loader maintains a table of the functions that are used for each system call.

The system call runs within the calling thread, but with more privilege because system calls run in the kernel protection domain. After the function implementing the system call has performed the requested action, control returns to the system call handler. If the `ut_error` field in the **uthread** structure has a non-zero value, the value is copied to the application's thread-specific **errno** variable. If a signal is pending, signal processing takes place, which can result in an application's signal handler being invoked. If no signals are pending, the system call handler restores the state of the calling thread, which is resumed in the user protection domain. For more information on protection domains, see "Understanding Protection Domains" on page 23.

---

## Accessing Kernel Data While in a System Call

A system call can access data that the calling thread cannot access because system calls execute in the kernel protection domain. The following are the general categories of kernel data:

- The **ublock** or **u-block** (user block data) structure:

System calls should use the kernel services to read or modify data traditionally found in the **ublock** or **uthread** structures. For example, the system call handler uses the value of the thread's `ut_error` field to update the thread-specific **errno** variable before returning to user mode. This field can be read or set by using the **getuerror** and **setuerror** kernel services. The current process ID can be obtained by using the **getpid** kernel service, and the current thread ID can be obtained by using the **thread\_self** kernel service.

- Global memory

System calls can also access global memory such as the kernel and kernel data regions. These regions contain the code and static data for the system call as well as the rest of the kernel.

- The stack for a system call:

A system call routine runs on a protected stack associated with a calling thread, which allows a system call to execute properly even when the stack pointer to the calling thread is invalid. In addition, privileged data can be saved on the stack without danger of exposing the data to the calling thread.

**Attention:** Incorrectly modifying fields in kernel or user block structures can cause unpredictable results or system crashes.

---

## Passing Parameters to System Calls

Parameters are passed to system calls in the same way that parameters are passed to other functions, but some additional calling conventions and limitations apply.

First, system calls cannot have floating-point parameters. In fact, the operating system does not preserve the contents of floating-point registers when a system call is preempted by another thread, so system calls cannot use any floating-point operations.

Second, because a system call runs on its own stack, the number of arguments that can be passed to a system call is limited. The operating system linkage conventions specify that up to eight general purpose registers are used for parameter passing. If more parameters exist than will fit in eight registers, the remaining parameters are passed in the stack. Because a system call does not have direct access to the application's stack, all parameters for system calls must fit in eight registers.

Third, some parameters are passed in multiple registers. For example, 32-bit applications pass **long long** parameters in two registers, and structures passed by value can require multiple registers, depending on the structure size. The writer of a system call should be familiar with the way parameters are passed by the compiler and ensure that the 8-register limit is not exceeded. For more information on parameter calling conventions, see Subroutine Linkage Convention in *Assembler Language Reference*.

Finally, because 32-bit applications are supported by the 64-bit kernel, the data model used by the kernel does not always match the data model used by the application. When the data models do not match, the system call might have to perform extra processing before parameters can be used.

The **IS64U** macro can be used to determine if the caller of a system call is a 64-bit process. For more information on the IS64U macro, see IS64U Kernel Service in *AIX Version 6.1 Technical Reference: Kernel and Subsystems Volume 1*.

The ILP32 and LP64 data models differ in the way that pointers and **long** and **long long** parameters are treated when used in structures or passed as functional parameters. The following tables summarize the differences.

Type	Size	Used as Parameter
<b>long</b>	32 bits	One register
<b>pointer</b>	32 bits	One register
<b>long long</b>	64 bits	Two registers

Type	Size	Used as Parameter
<b>long</b>	64 bits	One register
<b>pointer</b>	64 bits	One register
<b>long long</b>	64 bits	One register

System calls using these types must take the differing data models into account. The treatment of these types depends on whether they are used as parameters or in structures passed as parameters by value or by reference.

## Passing Scalar Parameters to System Calls

*Scalar* parameters (pointers and integral values) are passed in registers. The combinations of kernel and application modes are:

- 32-bit application support on the 64-bit kernel
- 64-bit application support on the 64-bit kernel

### 32-bit Application Support on the 64-bit Kernel

When a 32-bit application makes a system call to the 64-bit kernel, the system call handler zeros the high-order word of each parameter register. This allows 64-bit system calls to use pointers and unsigned **long** parameters directly. Signed and unsigned integer parameters can also be used directly by 64-bit system calls. This is because in 64-bit mode, the compiler generates code that sign extends or zero fills integers passed as parameters. Similar processing is performed for **char** and **short** parameters, so these types do not require any special handling either. Only signed **long** and **long long** parameters need additional processing.

**Signed long Parameters:** To convert a 32-bit signed **long** parameter to a 64-bit value, the 32-bit value must be sign extended. The **LONG32TOLONG64** macro is provided for this operation. It converts a 32-bit signed value into a 64-bit signed value, as shown in this example:

```
syscall1(long incr)
{
    /* If the caller is a 32-bit process, convert
     * 'incr' to a signed, 64-bit value.
     */
    if (!IS64U)
        incr = LONG32TOLONG64(incr);
    .
    .
    .
}
```

If a parameter can be either a pointer or a symbolic constant, special handling is needed. For example, if -1 is passed as a pointer argument to indicate a special case, comparing the pointer to -1 will fail, as will unconditionally sign-extending the parameter value. Code similar to the following should be used:

```
syscall2(void *ptr)
{
    /* If caller is a 32-bit process,
     * check for special parameter value.
     */
    if (!IS64U && (LONG32TOLONG64(ptr) == -1))
        ptr = (void *)-1;

    if (ptr == (void *)-1)
        special_handling();
    else {
        .
        .
        .
    }
}
```

Similar treatment is required when an unsigned long parameter is interpreted as a signed value.

**long long Parameters:** A 32-bit application passes a **long long** parameter in two registers, while a 64-bit kernel system call uses a single register for a **long long** parameter value.

The system call function prototype cannot match the function prototype used by the application. Instead, each **long long** parameter should be replaced by a pair of **uintptr\_t** parameters. Subsequent parameters should be replaced with **uintptr\_t** parameters as well. When the caller is a 32-bit process, a single 64-bit

value will be constructed from two consecutive parameters. This operation can be performed using the **INTSTOLLONG** macro. For a 64-bit caller, a single parameter is used directly.

For example, suppose the application function prototype is:

```
syscall3(void *ptr, long long len1, long long len2, int size);
```

The corresponding system call code should be similar to:

```
syscall3(void *ptr, uintptr_t L1,
         uintptr_t L2, uintptr_t L3,
         uintptr_t L4, uintptr_t L5)
{
    long len1;
    long len2;
    int size;

    /* If caller is a 32-bit application, len1
     * and len2 must be constructed from pairs of
     * parameters. Otherwise, a single parameter
     * can be used for each length.
     */

    if (!IS64U) {
        len1 = INTSTOLLONG(L1, L2);
        len2 = INTSTOLLONG(L3, L4);
        size = (int)L5;
    }
    else {
        len1 = (long)L1
        len2 = (long)L2
        size = (int)L3;
    }
    .
    .
    .
}
```

## 64-bit Application Support on the 64-bit Kernel

For the most part, system call parameters from a 64-bit application can be used directly by 64-bit system calls. The system call handler does not modify the parameter registers, so the system call sees the same values that were passed by the application. The only exceptions are the **pid\_t** and **key\_t** types, which are 32-bit signed types in 64-bit applications, but are 64-bit signed types in 64-bit system calls. Before these two types can be used, the 32-bit parameter values must be sign extended using the **LONG32TOLONG64** macro.

## Passing Structure Parameters to System Calls

When structures are passed to or from system calls, whether by value or by reference, the layout of the structure in the application might not match the layout of the same structure in the system call. There are two ways that system calls can process structures passed from or to applications: structure reshaping and dual implementation.

### Structure Reshaping

Structure reshaping allows system calls to support both 32- and 64-bit applications using a single system call interface and using code that is predominately common to both application types.

Structure reshaping requires defining more than one version of a structure. One version of the structure is used internally by the system call to process the request. The other version should use size-invariant types, so that the layout of the structure fields matches the application's view of the structures. When a structure is copied in from user space, the application-view structure definition is used. The structure is

reshaped by copying each field of the application's structure to the kernel's structure, converting the fields as required. A similar conversion is performed on structures that are being returned to the caller.

Structure reshaping is used for structures whose size and layout as seen by an application differ from the size and layout as seen by the system call. If the system call uses a structure definition with fields big enough for both 32- and 64-bit applications, the system call can use this structure, independent of the mode of the caller.

While reshaping requires two versions of a structure, only one version is public and visible to the end user. This version is the natural structure, which can also be used by the system call if reshaping is not needed. The private version should only be defined in the source file that performs the reshaping. The following example demonstrates the techniques for passing structures to system calls that are running in the 64-bit kernel and how a structure can be reshaped:

```
/* Public definition */
struct foo {
    int a;
    long b;
};

/* Private definition--matches 32-bit
 * application's view of the data structure. */
struct foo32 {
    int a;
    int b;
}

syscall7(struct foo *f)
{
    struct foo f1;
    struct foo32 f2;

    if (IS64U()) {
        copyin(&f1, f, sizeof(f1));
    }
    else {
        copyin(&f2, f, sizeof(f2));
        f1.a = f2.a;
        f1.b = f2.b;
    }
    /* Common structure f1 used from now on. */
    .
    .
    .
}
```

**Dual Implementation:** The dual implementation approach involves separate code paths for calls from 32-bit applications and calls from 64-bit applications. Similar to reshaping, the system call code defines a private view of the application's structure. With dual implementations, the function *syscall7* could be rewritten as:

```
syscall8(struct foo *f)
{
    struct foo f1;
    struct foo32 f2;

    if (IS64U()) {
        copyin(&f1, f, sizeof(f1));
        /* Code for 64-bit process uses f1 */
        .
        .
        .
    }
    else {
        copyin(&f2, f, sizeof(f2));
    }
}
```

```

        /* Code for 32-bit process uses f2 */
        .
        .
    }
}

```

Dual implementation is most appropriate when the structures are so large that the overhead of reshaping would affect the performance of the system call.

**Passing Structures by Value:** When structures are passed by value, the structure is loaded into as many parameter registers as are needed. When the data model of an application and the data model of the kernel extension differ, the values in the registers cannot be used directly. Instead, the registers must be stored in a temporary variable. For example:

**Note:** This example builds upon the structure definitions defined in “Dual Implementation” on page 28.

```

/* Application prototype: syscall19(struct foo f); */

syscall19(unsigned long a1, unsigned long a1)
{
    union {
        struct foo  f1;      /* Structure for 64-bit caller. */
        struct foo32 f2;    /* Structure for 32-bit caller. */
        unsigned long p64[2]; /* Overlay for parameter registers
                               * when caller is 64-bit program
                               */
        unsigned int  p32[2]; /* Overlay for parameter registers
                               * when caller is 32-bit program
                               */
    } uarg;
    if (IS64U()) {
        uarg.p64[0] = a1;
        uarg.p64[1] = a2;
        /* Now uarg.f1 can be used */
        .
        .
    }
    else {
        uarg.p32[0] = a1;
        uarg.p32[1] = a2;
        /* Now uarg.f2 can be used */
        .
        .
    }
}

```

### Comparisons to AIX 4.3

In AIX 4.3, the conventions for passing parameters from a 64-bit application to a system call required user-space library code to perform some of the parameter reshaping and address mapping. In AIX 5.1 and later, all parameter reshaping and address mapping should be performed by the system call, eliminating the need for kernel-specific library code. In fact, user-space address mapping is no longer supported. In most cases, system calls can be implemented without any application-specific library code.

---

## Preempting a System Call

The kernel allows a thread to be preempted by a more favored thread, even when a system call is executing. This capability provides better system responsiveness for large multi-user systems.

Because system calls can be preempted, access to global data must be serialized. Kernel locking services, such as **simple\_lock** and **simple\_unlock**, are frequently used to serialize access to kernel data.

A thread can be preempted even when it owns a lock. If multiple locks are obtained by system calls, a technique must be used to prevent multiple threads from deadlocking. One technique is to define a lock hierarchy. A system call must never return while holding a lock. For more information on locking, see “Understanding Locking” on page 12.

---

## Handling Signals While in a System Call

Signals can be generated asynchronously or synchronously with respect to the thread that receives the signal. An asynchronously generated signal is one that results from some action external to a thread. It is not directly related to the current instruction stream of that thread. Generally these are generated by other threads or by device drivers.

A synchronously generated signal is one that results from the current instruction stream of the thread. These signals cause interrupts. Examples of such cases are the execution of an illegal instruction, or an attempted data access to nonexistent address space.

## Delivery of Signals to a System Call

Delivery of signals to a thread only takes place when a user application is about to be resumed in the user protection domain. Signals cannot be delivered to a thread if the thread is in the middle of a system call. For more information on signal delivery for kernel processes, see “Using Kernel Processes” on page 8.

## Asynchronous Signals and Wait Termination

An asynchronous signal can alter the operation of a system call or kernel extension by terminating a long wait. Kernel services such as **e\_block\_thread**, **e\_sleep\_thread**, and **et\_wait** are affected by signals. The following options are provided when a signal is posted to a thread:

- Return from the kernel service with a return code indicating that the call was interrupted by a signal
- Call the **longjmpx** kernel service to resume execution at a previously saved context in the event of a signal
- Ignore the signal using the **short-wait** option, allowing the kernel service to return normally.

The **sleep** kernel service, provided for compatibility, also supports the **PCATCH** and **SWAKEONSIG** options to control the response to a signal during the **sleep** function.

Previously, the kernel automatically saved context on entry to the system call handler. As a result, any long (interruptible) sleep not specifying the **PCATCH** option returned control to the saved context when a signal interrupted the wait. The system call handler then set the **errno** global variable to **EINTR** and returned a return code of -1 from the system call.

The kernel, however, requires each system call that can directly or indirectly issue a **sleep** call without the **PCATCH** option to set up a saved context using the **setjmpx** kernel service. This is done to avoid overhead for system calls that handle waits terminated by signals. Using the **setjmpx** service, the system can set up a saved context, which sets the system call return code to a -1 and the **ut\_error** field to **EINTR**, if a signal interrupts a long wait not specifying **return-from-signal**.

It is probably faster and more robust to specify **return-from-signal** on all long waits and use the return code to control the system call return.

## Stacking Saved Contexts for Nested setjmpx Calls

The kernel supports nested calls to the **setjmpx** kernel service. It implements the stack of saved contexts by maintaining a linked list of context information anchored in the machine state save area. This area is in the user block structure for a process. Interrupt handlers have special machine state save areas.

An initial context is set up for each process by the **initp** kernel service for kernel processes and by the **fork** subroutine for user processes. The process terminates if that context is resumed.



---

## Handling Exceptions While in a System Call

Exceptions are interrupts detected by the processor as a result of the current instruction stream. They therefore take effect synchronously with respect to the current thread.

The default exception handler generates a signal if the process is in a state where signals can be delivered immediately. Otherwise, the default exception handler generates a system dump.

## Alternative Exception Handling Using the `setjmpx` Kernel Service

For certain types of exceptions, a system call can specify unique exception-handler routines through calls to the `setjmpx` service. The exception handler routine is saved as part of the stacked saved context. Each exception handler is passed the exception type as a parameter.

The exception handler returns a value that can specify any of the following:

- The process should resume with the instruction that caused the exception.
- The process should return to the saved context that is on the top of the stack of contexts.
- The exception handler did not handle the exception.

If the exception handler did not handle the exception, then the next exception handler in the stack of contexts is called. If none of the stacked exception handlers handle the exception, the kernel performs default exception handling. The `setjmpx` and `longjmpx` kernel services help implement exception handlers.

---

## Understanding Nesting and Kernel-Mode Use of System Calls

The operating system supports nested system calls with some restrictions. System calls (and any other kernel-mode routines running under the process environment of a user-mode process) can use system calls that pass all parameters by value. System calls and other kernel-mode routines must not start system calls that have one or more parameters passed by reference. Doing so can result in a system crash. This is because system calls with reference parameters assume that the referenced data area is in the user protection domain. As a result, these system calls must use special kernel services to access the data. However, these services are unsuccessful if the data area they are trying to access is not in the user protection domain.

This restriction does not apply to kernel processes. User-mode data access services can distinguish between kernel processes and user-mode processes in kernel mode. As a result, these services can access the referenced data areas accessed correctly when the caller is a kernel process.

Kernel processes cannot call the `fork` or `exec` system calls, among others. A list of the base operating system calls available to system calls or other routines in kernel mode is provided in “System Calls Available to Kernel Extensions” on page 32.

---

## Page Faulting within System Calls

**Attention:** A page fault that occurs while external interrupts are disabled results in a system crash. Therefore, a system call should be programmed to ensure that its code, data, and stack are pinned before it disables external interrupts.

Most data accessed by system calls is pageable by default. This includes the system call code, static data, dynamically allocated data, and stack. As a result, a system call can be preempted in two ways:

- By a more favored process, or by an equally favored process when a time slice has been exhausted
- By losing control of the processor when it page faults

In the latter case, even less-favored processes can run while the system call is waiting for the paging I/O to complete.

---

## Returning Error Information from System Calls

Error information returned by system calls differs from that returned by kernel services that are not system calls. System calls typically return a special value, such as -1 or NULL, to indicate that an error has occurred. When an error condition is to be returned, the `ut_error` field should be updated by the system call before returning from the system call function. The `ut_error` field is written using the **setuerror** kernel service.

Before actually calling the system call function, the system call handler sets the `ut_error` field to 0. Upon return from the system call function, the system call handler copies the value found in `ut_error` into the thread-specific **errno** variable if `ut_error` was nonzero. After setting the **errno** variable, the system call handler returns to user mode with the return code provided by the system call function.

Kernel-mode callers of system calls must be aware of this return code convention and use the **getuerror** kernel service to obtain the error value when an error indication is returned by the system call. When system calls are nested, the system call function called by the system call handler can return the error value provided by the nested system call function or can replace this value with a new one by using the **setuerror** kernel service.

---

## System Calls Available to Kernel Extensions

The following system calls are grouped according to which subroutines call them:

- System calls available to all kernel extensions
- System calls available to kernel processes only

**Note:** System calls are not available to interrupt handlers.

## System Calls Available to All Kernel Extensions

<b>gethostid</b>	Gets the unique identifier of the current host.
<b>getpgrp</b>	Gets the process ID, process group ID, and parent process ID.
<b>getppid</b>	Gets the process ID, process group ID, and parent process ID.
<b>getpri</b>	Returns the scheduling priority of a process.
<b>getpriority</b>	Gets or sets the <i>nice</i> value.
<b>semget</b>	Gets a set of semaphores.
<b>seteuid</b>	Sets the process user IDs.
<b>setgid</b>	Sets the process group IDs.
<b>sethostid</b>	Sets the unique identifier of the current host.
<b>setpgid</b>	Sets the process group IDs.
<b>setpgrp</b>	Sets the process group IDs.
<b>setpri</b>	Sets a process scheduling priority to a constant value.
<b>setpriority</b>	Gets or sets the <i>nice</i> value.
<b>setreuid</b>	Sets the process user IDs.
<b>setsid</b>	Creates a session and sets the process group ID.
<b>setuid</b>	Sets the process user IDs.
<b>ulimit</b>	Sets and gets user limits.
<b>umask</b>	Sets and gets the value of the file-creation mask.

## System Calls Available to Kernel Processes Only

<b>disclaim</b>	Disclaims the content of a memory address range.
-----------------	--

<b>getdomainname</b>	Gets the name of the current domain.
<b>getgroups</b>	Gets the concurrent group set of the current process.
<b>gethostname</b>	Gets the name of the local host.
<b>getpeername</b>	Gets the name of the peer socket.
<b>getrlimit</b>	Controls maximum system resource consumption.
<b>getrusage</b>	Displays information about resource use.
<b>getsockname</b>	Gets the socket name.
<b>getsockopt</b>	Gets options on sockets.
<b>gettimer</b>	Gets and sets the current value for the specified system-wide timer.
<b>resabs</b>	Manipulates the expiration time of interval timers.
<b>resinc</b>	Manipulates the expiration time of interval timers.
<b>restimer</b>	Gets and sets the current value for the specified system-wide timer.
<b>semctl</b>	Controls semaphore operations.
<b>semop</b>	Performs semaphore operations.
<b>setdomainname</b>	Sets the name of the current domain.
<b>setgroups</b>	Sets the concurrent group set of the current process.
<b>sethostname</b>	Sets the name of the current host.
<b>setrlimit</b>	Controls maximum system resource consumption.
<b>settimer</b>	Gets and sets the current value for the specified systemwide timer.
<b>shmat</b>	Attaches a shared memory segment or a mapped file to the current process.
<b>shmctl</b>	Controls shared memory operations.
<b>shmdt</b>	Detaches a shared memory segment.
<b>shmget</b>	Gets shared memory segments.
<b>sigaction</b>	Specifies the action to take upon delivery of a signal.
<b>sigprocmask</b>	Sets the current signal mask.
<b>sigstack</b>	Sets and gets signal stack context.
<b>sigsuspend</b>	Atomically changes the set of blocked signals and waits for a signal.
<b>sysconfig</b>	Provides a service for controlling system/kernel configuration.
<b>sys_parm</b>	Provides a service for examining or setting kernel run-time tunable parameters.
<b>times</b>	Displays information about resource use.
<b>uname</b>	Gets the name of the current system.
<b>unamex</b>	Gets the name of the current system.
<b>usrinfo</b>	Gets and sets user information about the owner of the current process.
<b>utimes</b>	Sets file access and modification times.

---

## Related Information

“Handling Signals While in a System Call” on page 30

“Understanding Protection Domains” on page 23

“Understanding Kernel Threads” on page 6

“Using Kernel Processes” on page 8

“Using Libraries” on page 4

“Understanding Locking” on page 12

“Locking Kernel Services” on page 60

“Understanding Interrupts” on page 49

## Subroutine References

The **fork** subroutine in *AIX Version 6.1 Technical Reference: Base Operating System and Extensions Volume 1*.

## Technical References

The **e\_sleep** kernel service, **e\_sleepl** kernel service, **et\_wait** kernel service, **getuerror** kernel service, **initp** kernel service, **lockl** kernel service, **longjmpx** kernel service, **setjmpx** kernel service, **setuerror** kernel service, **unlockl** kernel service in *AIX Version 6.1 Technical Reference: Kernel and Subsystems Volume 1*.

---

## Chapter 3. Virtual File Systems

The virtual file system (VFS) interface, also known as the v-node interface, provides a bridge between the physical and logical file systems. The information that follows discusses the virtual file system interface, its data structures, and its header files, and explains how to configure a virtual file system.

There are two essential components in the file system:

<b>Logical file system</b>	Provides support for the system call interface.
<b>Physical file system</b>	Manages permanent storage of data.

The interface between the physical and logical file systems is the *virtual file system interface*. This interface allows support for multiple concurrent instances of physical file systems, each of which is called a file system implementation. The file system implementation can support storing the file data in the local node or at a remote node. For more information on the virtual filesystem interface, see “Understanding the Virtual File System Interface” on page 37.

The virtual file system interface is usually referred to as the *v-node* interface. The v-node structure is the key element in communication between the virtual file system and the layers that call it. For more information on v-nodes, see “Understanding Virtual Nodes (V-nodes)” on page 36.

Both the virtual and logical file systems exist across all of this operating system family’s platforms.

---

### Logical File System Overview

The *logical file system* is the level of the file system at which users can request file operations by system call. This level of the file system provides the kernel with a consistent view of what might be multiple physical file systems and multiple file system implementations. As far as the logical file system is concerned, file system types, whether local, remote, or strictly logical, and regardless of implementation, are indistinguishable.

A consistent view of file system implementations is made possible by the virtual file system abstraction. This abstraction specifies the set of file system operations that an implementation must include in order to carry out logical file system requests. Physical file systems can differ in how they implement these predefined operations, but they must present a uniform interface to the logical file system. A list of file system operators can be found at “Requirements for a File System Implementation” on page 37. For more information on the virtual file system, see “Virtual File System Overview” on page 36.

Each set of predefined operations implemented constitutes a virtual file system. As such, a single physical file system can appear to the logical file system as one or more separate virtual file systems.

Virtual file system operations are available at the logical file system level through the *virtual file system switch*. This array contains one entry for each virtual file system, with each entry holding entry point addresses for separate operations. Each file system type has a set of entries in the virtual file system switch.

The logical file system and the virtual file system switch support other operating system file-system access semantics. This does not mean that only other operating system file systems can be supported. It does mean, however, that a file system implementation must be designed to fit into the logical file system model. Operations or information requested from a file system implementation need be performed only to the extent possible.

Logical file system can also refer to the tree of known path names in force while the system is running. A virtual file system that is mounted onto the logical file system tree itself becomes part of that tree. In fact, a

single virtual file system can be mounted onto the logical file system tree at multiple points, so that nodes in the virtual subtree have multiple names. Multiple mount points allow maximum flexibility when constructing the logical file system view.

## Component Structure of the Logical File System

The logical file system is divided into the following components:

- System calls

Implement services exported to users. System calls that carry out file system requests do the following:

- Map the user's parameters to a file system object. This requires that the system call component use the v-node (virtual node) component to follow the object's path name. In addition, the system call must resolve a file descriptor or establish implicit (mapped) references using the open file component.
- Verify that a requested operation is applicable to the type of the specified object.
- Dispatch a request to the file system implementation to perform operations.

- Logical file system file routines

Manage open file table entries and per-process file descriptors. An open file table entry records the authorization of a process's access to a file system object. A user can refer to an open file table entry through a file descriptor or by accessing the virtual memory to which the file was mapped. The logical file system routines are those kernel services, such as **fp\_ioctl** and **fp\_select**, that begin with the prefix **fp\_**.

- v-nodes

Provide system calls with a mechanism for local name resolution. Local name resolution allows the logical file system to access multiple file system implementations through a uniform name space.

---

## Virtual File System Overview

The virtual file system is an abstraction of a physical file system implementation. It provides a consistent interface to multiple file systems, both local and remote. This consistent interface allows the user to view the directory tree on the running system as a single entity even when the tree is made up of a number of diverse file system types. The interface also allows the logical file system code in the kernel to operate without regard to the type of file system being accessed. For more information on the logical file system, see "Logical File System Overview" on page 35.

A virtual file system can also be viewed as a subset of the logical file system tree, that part belonging to a single file system implementation. A virtual file system can be physical (the instantiation of a physical file system), remote, or strictly logical. In the latter case, for example, a virtual file system need not actually be a true file system or entail any underlying physical storage device.

A virtual file system mount point grafts a virtual file system subtree onto the logical file system tree. This mount point ties together a mounted-over v-node (virtual node) and the root of the virtual file system subtree. A mounted-over, or stub, v-node points to a virtual file system, and the mounted VFS points to the v-node it is mounted over.

## Understanding Virtual Nodes (V-nodes)

A *virtual node* (v-node) represents access to an object within a virtual file system. V-nodes are used only to translate a path name into a generic node (g-node). For more information on g-nodes, see "Understanding Generic I-nodes (G-nodes)" on page 37.

A v-node is either created or used again for every reference made to a file by path name. When a user attempts to open or create a file, if the VFS containing the file already has a v-node representing that file, a use count in the v-node is incremented and the existing v-node is used. Otherwise, a new v-node is created.

Every path name known to the logical file system can be associated with, at most, one file system object. However, each file system object can have several names. Multiple names appear in the following cases:

- The object can appear in multiple virtual file systems. This can happen if the object (or an ancestor) is mounted in different virtual file systems using a local file-over-file or directory-over-directory mount.
- The object does not have a unique name within the virtual file system. (The file system implementation can provide synonyms. For example, the use of links causes files to have more than one name. However, opens of synonymous paths do not cause multiple v-nodes to be created.)

## Understanding Generic I-nodes (G-nodes)

A *generic i-node* (g-node) is the representation of an object in a file system implementation. There is a one-to-one correspondence between a g-node and an object in a file system implementation. Each g-node represents an object owned by the file system implementation.

Each file system implementation is responsible for allocating and destroying g-nodes. The g-node then serves as the interface between the logical file system and the file system implementation. Calls to the file system implementation serve as requests to perform an operation on a specific g-node.

A g-node is needed, in addition to the file system i-node, because some file system implementations may not include the concept of an i-node. Thus the g-node structure substitutes for whatever structure the file system implementation may have used to uniquely identify a file system object.

The logical file system relies on the file system implementation to provide valid data for the following fields in the g-node:

<b>gn_type</b>	Identifies the type of object represented by the g-node.
<b>gn_ops</b>	Identifies the set of operations that can be performed on the object.

## Understanding the Virtual File System Interface

Operations that can be performed upon a virtual file system and its underlying objects are divided into two categories. Operations upon a file system implementation as a whole (not requiring the existence of an underlying file system object) are called **vfs** operations. Operations upon the underlying file system objects are called v-node (virtual node) operations. Before writing specific virtual file system operations, it is important to note the requirements for a file system implementation.

### Requirements for a File System Implementation

File system implementations differ in how they implement the predefined operations. However, the logical file system expects that a file system implementation meets the following criteria:

- All **vfs** and v-node operations must supply a return value:
  - A return value of 0 indicates the operation was successful.
  - A nonzero return value is interpreted as a valid error number (taken from the `/usr/include/sys/errno.h` file) and returned through the system call interface to the application program.
- All **vfs** operations must exist for each file system type, but can return an error upon startup. The following are the necessary **vfs** operations:
  - **vfs\_cntl**
  - **vfs\_mount**
  - **vfs\_root**
  - **vfs\_statfs**
  - **vfs\_sync**
  - **vfs\_unmount**
  - **vfs\_vget**
  - **vfs\_quotactl**

For a complete list of file system operations, see List of Virtual File System Operations in *AIX Version 6.1 Technical Reference: Kernel and Subsystems Volume 1*.

## Important Data Structures for a File System Implementation

There are two important data structures used to represent information about a virtual file system, the **vfs** structure and the v-node. Each virtual file system has a **vfs** structure in memory that describes its type, attributes, and position in the file tree hierarchy. Each file object within that virtual file system can be represented by a v-node.

The **vfs** structure contains the following fields:

<code>vfs_flag</code>	Contains the state flags: <b>VFS_DEVMOUNT</b> Indicates whether the virtual file system has a physical mount structure underlying it. <b>VFS_READONLY</b> Indicates whether the virtual file system is mounted read-only.
<code>vfs_type</code>	Identifies the type of file system implementation. Possible values for this field are described in the <code>/usr/include/sys/vmount.h</code> file.
<code>vfs_ops</code>	Points to the set of operations for the specified file system type.
<code>vfs_mntdover</code>	Points to the mounted-over v-node.
<code>vfs_data</code>	Points to the file system implementation data. The interpretation of this field is left to the discretion of the file system implementation. For example, the field could be used to point to data in the kernel extension segment or as an offset to another segment.
<code>vfs_mdata</code>	Records the user arguments to the <b>mount</b> call that created this virtual file system. This field has a time stamp. The user arguments are retained to implement the <b>mntctl</b> call, which replaces the <code>/etc/mnttab</code> table.

---

## Understanding Data Structures and Header Files for Virtual File Systems

These are the data structures used in implementing virtual file systems:

- The **vfs** structure contains information about a virtual file system as a single entity.
- The **vnode** structure contains information about a file system object in a virtual file system. There can be multiple v-nodes for a single file system object.
- The **gnode** structure contains information about a file system object in a physical file system. There is only a single g-node for a given file system object.
- The **gfs** structure contains information about a file system implementation. This is distinct from the **vfs** structure, which contains information about an instance of a virtual file system.

The header files contain the structure definitions for the key components of the virtual file system abstraction. Understanding the contents of these files and the relationships between them is essential to an understanding of virtual file systems. The following are the necessary header files:

- `sys/vfs.h`
- `sys/gfs.h`
- `sys/vnode.h`
- `sys/vmount.h`



---

## Configuring a Virtual File System

The kernel maintains a table of active file system types. A file system implementation must be registered with the kernel before a request to mount a virtual file system (VFS) of that type can be honored. Two kernel services, **gfsadd** and **gfsdel**, are supplied for adding a file system type to the **gfs** file system table.

These are the steps that must be followed to get a file system configured.

1. A user-level routine must call the **sysconfig** subroutine requesting that the code for the virtual file system be loaded.
2. The user-level routine must then request, again by calling the **sysconfig** subroutine, that the virtual file system be configured. The name of a VFS-specific configuration routine must be specified.
3. The virtual file system-specific configuration routine calls the **gfsadd** kernel service to have the new file system added to the **gfs** table. The **gfs** table that the configuration routine passes to the **gfsadd** kernel service contains a pointer to an initialization routine. This routine is then called to do any further virtual file system-specific initialization.
4. The file system is now operational.

---

## Related Information

“Logical File System Kernel Services” on page 63

“Understanding Data Structures and Header Files for Virtual File Systems” on page 38

“Configuring a Virtual File System”

“Understanding Protection Domains” on page 23

List of Virtual File System Operations in *AIX Version 6.1 Technical Reference: Kernel and Subsystems Volume 1*.

## Subroutine References

The **mntctl** subroutine, **mount** subroutine, **sysconfig** subroutine in *AIX Version 6.1 Technical Reference: Base Operating System and Extensions Volume 1*.

## Files References

The **vmount.h** file in *AIX Version 6.1 Files Reference*.

## Technical References

The **gfsadd** kernel service, **gfsdel** kernel service in *AIX Version 6.1 Technical Reference: Kernel and Subsystems Volume 1*.



---

## Chapter 4. Kernel Services

*Kernel services* are routines that provide the runtime kernel environment to programs executing in kernel mode. Kernel extensions call kernel services, which resemble library routines. In contrast, application programs call library routines.

Callers of kernel services execute in kernel mode. They therefore share with the kernel the responsibility for ensuring that system integrity is not compromised.

For a list of system calls that kernel extensions are allowed to use, see “System Calls Available to Kernel Extensions” on page 32.

---

### Categories of Kernel Services

Following are the categories of kernel services:

- “I/O Kernel Services”
- “Kernel Extension and Device Driver Management Services” on page 58
- “Locking Kernel Services” on page 60
- “Logical File System Kernel Services” on page 63
- “Memory Kernel Services” on page 64
- “Message Queue Kernel Services” on page 72
- “Network Kernel Services” on page 72
- “Process and Exception Management Kernel Services” on page 74
- “RAS Kernel Services” on page 77
- “Security Kernel Services” on page 77
- “Timer and Time-of-Day Kernel Services” on page 78
- “Virtual File System (VFS) Kernel Services” on page 82

---

### I/O Kernel Services

The I/O kernel services fall into the following categories:

- “Block I/O Kernel Services”
- “Buffer Cache Kernel Services” on page 42
- “Character I/O Kernel Services” on page 42
- “Interrupt Management Kernel Services” on page 42
- “Memory Buffer (mbuf) Kernel Services” on page 43
- “DMA Management Kernel Services” on page 43
- “Enhanced I/O Error Handling Kernel Services” on page 44

### Block I/O Kernel Services

The Block I/O kernel services are:

<b>iodone</b>	Performs block I/O completion processing.
<b>iowait</b>	Waits for block I/O completion.
<b>uphysio</b>	Performs character I/O for a block device using a <b>uio</b> structure.

## Buffer Cache Kernel Services

For information on how to manage the buffer cache with the Buffer Cache kernel services, see “Block I/O Buffer Cache Kernel Services: Overview” on page 48. The Buffer Cache kernel services are:

<b>bawrite</b>	Writes the specified buffer's data without waiting for I/O to complete.
<b>bdwrite</b>	Releases the specified buffer after marking it for delayed write.
<b>bflush</b>	Flushes all write-behind blocks on the specified device from the buffer cache.
<b>binval</b>	Invalidates all of the specified device's blocks in the buffer cache.
<b>blkflush</b>	Flushes the specified block if it is in the buffer cache.
<b>bread</b>	Reads the specified block's data into a buffer.
<b>breada</b>	Reads in the specified block and then starts I/O on the read-ahead block.
<b>brelease</b>	Frees the specified buffer.
<b>bwrite</b>	Writes the specified buffer's data.
<b>clrbuf</b>	Sets the memory for the specified buffer structure's buffer to all zeros.
<b>getblk</b>	Assigns a buffer to the specified block.
<b>getblk</b>	Allocates a free buffer.
<b>geterror</b>	Determines the completion status of the buffer.
<b>purblk</b>	Purges the specified block from the buffer cache.

## Character I/O Kernel Services

The Character I/O kernel services are:

<b>getc</b>	Retrieves a character from a character list.
<b>getcb</b>	Removes the first buffer from a character list and returns the address of the removed buffer.
<b>getcbp</b>	Retrieves multiple characters from a character buffer and places them at a designated address.
<b>getc</b>	Retrieves a free character buffer.
<b>getc</b>	Returns the character at the end of a designated list.
<b>pincf</b>	Manages the list of free character buffers.
<b>putc</b>	Places a character at the end of a character list.
<b>putcb</b>	Places a character buffer at the end of a character list.
<b>putcbp</b>	Places several characters at the end of a character list.
<b>putc</b>	Frees a specified buffer.
<b>putc</b>	Frees the specified list of buffers.
<b>putc</b>	Places a character on a character list.
<b>waitcfree</b>	Checks the availability of a free character buffer.

## Interrupt Management Kernel Services

The operating system provides the following set of kernel services for managing interrupts. See Understanding Interrupts for a description of these services:

<b>i_clear</b>	Removes an interrupt handler from the system.
<b>i_reset</b>	Resets a bus interrupt level.
<b>i_sched</b>	Schedules off-level processing.
<b>i_mask</b>	Disables an interrupt level.
<b>i_unmask</b>	Enables an interrupt level.
<b>i_disable</b>	Disables all of the interrupt levels at a particular interrupt priority and all interrupt levels at a less-favored interrupt priority.
<b>i_enable</b>	Enables all of the interrupt levels at a particular interrupt priority and all interrupt levels at a more-favored interrupt priority.

## Memory Buffer (mbuf) Kernel Services

The Memory Buffer (mbuf) kernel services provide functions to obtain, release, and manipulate memory buffers, or **mbufs**. These **mbuf** services provide the means to easily work with the **mbuf** data structure, which is defined in the `/usr/include/sys/mbuf.h` file. Data can be stored directly in an **mbuf**'s data portion or in an attached external cluster. **Mbufs** can also be chained together by using the `m_next` field in the **mbuf** structure. This is particularly useful for communications protocols that need to add and remove protocol headers.

The Memory Buffer (**mbuf**) kernel services are:

<b>m_adj</b>	Adjusts the size of an <b>mbuf</b> chain.
<b>m_clattach</b>	Allocates an <b>mbuf</b> structure and attaches an external cluster.
<b>m_cat</b>	Appends one <b>mbuf</b> chain to the end of another.
<b>m_clgetm</b>	Allocates and attaches an external buffer.
<b>m_collapse</b>	Guarantees that an <b>mbuf</b> chain contains no more than a given number of <b>mbuf</b> structures.
<b>m_copydata</b>	Copies data from an <b>mbuf</b> chain to a specified buffer.
<b>m_copym</b>	Creates a copy of all or part of a list of <b>mbuf</b> structures.
<b>m_dereg</b>	Deregisters expected <b>mbuf</b> structure usage.
<b>m_free</b>	Frees an <b>mbuf</b> structure and any associated external storage area.
<b>m_freem</b>	Frees an entire <b>mbuf</b> chain.
<b>m_get</b>	Allocates a memory buffer from the <b>mbuf</b> pool.
<b>m_getclr</b>	Allocates and zeros a memory buffer from the <b>mbuf</b> pool.
<b>m_getclustm</b>	Allocates an <b>mbuf</b> structure from the <b>mbuf</b> buffer pool and attaches a cluster of the specified size.
<b>m_gethdr</b>	Allocates a header memory buffer from the <b>mbuf</b> pool.
<b>m_pullup</b>	Adjusts an <b>mbuf</b> chain so that a given number of bytes is in contiguous memory in the data area of the head <b>mbuf</b> structure.
<b>m_reg</b>	Registers expected <b>mbuf</b> usage.

In addition to the **mbuf** kernel services, the following macros are available for use with **mbufs**:

<b>m_clget</b>	Allocates a page-sized <b>mbuf</b> structure cluster.
<b>m_copy</b>	Creates a copy of all or part of a list of <b>mbuf</b> structures.
<b>m_getclust</b>	Allocates an <b>mbuf</b> structure from the <b>mbuf</b> buffer pool and attaches a page-sized cluster.
<b>M_HASCL</b>	Determines if an <b>mbuf</b> structure has an attached cluster.
<b>DTOM</b>	Converts an address anywhere within an <b>mbuf</b> structure to the head of that <b>mbuf</b> structure.
<b>MTOCL</b>	Converts a pointer to an <b>mbuf</b> structure to a pointer to the head of an attached cluster.
<b>MTOD</b>	Converts a pointer to an <b>mbuf</b> structure to a pointer to the data stored in that <b>mbuf</b> structure.
<b>M_XMEMD</b>	Returns the address of an <b>mbuf</b> cross-memory descriptor.

## DMA Management Kernel Services

The operating system kernel provides several services for managing direct memory access (DMA) channels and performing DMA operations. Understanding DMA Transfers provides additional kernel services information.

The services provided are:

<b>d_align</b>	Provides needed information to align a buffer with a processor cache line.
<b>d_cflush</b>	Flushes the processor and I/O controller (IOCC) data caches when using the long term <b>DMA_WRITE_ONLY</b> mapping of DMA buffers approach to the bus device DMA.
<b>d_map_clear</b>	Deallocates resources previously allocated on a <code>d_map_init</code> call.
<b>d_map_disable</b>	Disables DMA for the specified handle.
<b>d_map_enable</b>	Enables DMA for the specified handle.
<b>d_map_init</b>	Allocates and initializes resources for performing DMA with PCI and ISA devices.

<code>d_map_list</code>	Performs platform-specific DMA mapping for a list of virtual addresses.
<code>d_map_page</code>	Performs platform-specific DMA mapping for a single page.
<code>d_map_slave</code>	Accepts a list of virtual addresses and sizes and sets up the slave DMA controller.
<code>d_roundup</code>	Rounds the value length up to a given number of cache lines.
<code>d_unmap_list</code>	Deallocates resources previously allocated on a <code>d_map_list</code> call.
<code>d_unmap_page</code>	Deallocates resources previously allocated on a <code>d_map_page</code> call.
<code>d_unmap_slave</code>	Deallocates resources previously allocated on a <code>d_map_slave</code> call.

## Enhanced I/O Error Handling Kernel Services

Enhanced I/O Error Handling (EEH) kernel services is an error recovery strategy for errors that occur during I/O operations on a Peripheral Component Interconnect (PCI), Peripheral Component Interconnect Extended (PCI-X), or PCI Express (PCIe) bus. Bridges and switches, PCI-to-PCI or PCI-X-to-PCI-X PCIe switches, that enable each slot to be on its own bus provide a form of electrical and logical isolation of slots. The smallest PCI unit (device or function) that can be assigned to a logical partition (LPAR) is called a *Partitionable Endpoint (PE)*.

For example, if each port of a 4-port Ethernet adapter can be assigned to a different LPAR, and each port is a PCI function, then each port would be a PE. In this case, granularity of a PE is a PCI function. A PE is not necessarily the same as a PCIe Endpoint, and the two should not be confused. A bridge or switch above the PE where the EEH state is maintained forms a PE domain. This bridge or switch is called the *Top of PE domain*. EEH recovery is performed according to the PE domain and is carried out by the Top of PE domain as directed by the software (operating system and device drivers).

Several PCI functions in one or more adapters that belong to the same EEH recovery domain are called a *shared EEH domain*. This has been typically limited to a multifunction adapter, in which the functions on the adapter are recovered together. Because a shared EEH domain supports any number of PCI functions to be included in it, including the functions on different adapters, its function is more general than a multifunction adapter. For present purposes, the multifunction model are referred to as the *shared EEH model*.

In the LPAR environment, a PE domain is the same as a shared EEH domain and includes all PCI functions in the PE domain. In other words, if multiple functions belong to a shared EEH domain, they cannot be individual PEs because the EEH recovery can only affect the LPAR to which the PE belongs.

The types of adapters supported in the slot created by a PE domain are:

- Single-function adapter with or without a bridge or switch on the adapter.
- Multifunction adapter without a bridge or switch on the adapter
- Multifunction adapter with a bridge or switch on the adapter.

The bridges can be of different types, such as PCI-to-PCI or PCIX-to-PCI, and so on. A switch is a PCIe switch, which is logically a collection of bridges. The exact type of the bridge or switch is not important to this programming model. These details are handled by the hardware and firmware. A bridged single-function adapter is treated like a bridged multifunction adapter for the purposes of EEH programming model.

The device drivers for all these types of adapters use the same EEH kernel services to drive the error recovery except for the registration service. A nonbridged single-function adapter calls the `eeh_init()` registration service function. Adapters in a shared EEH domain call the `eeh_init_multifunc()` function. This includes any bridged or nonbridged multifunction adapters and bridged single-function adapters. Although a nonbridged single-function adapter typically calls `eeh_init()`, it can choose the shared EEH model and call `eeh_init_multifunc()` instead. Regardless of the number of functions and bridges, the device drivers should always use the shared EEH model and call `eeh_init_multifunc()`. The PCIe device drivers are required to use the shared EEH model. Although the same services are used by the single and shared EEH adapter drivers, the error recovery models are different. In addition, any time there are intermediate

bridges between the Top of PE domain and the PCI functions in the PE domain, those bridges have to be recovered as well. For example, a multifunction bridged adapter requires that the bridge on the adapter also be recovered.

The error recovery is performed by resetting the PCI bus between the Top of PE domain and the PE under it, and then reconfiguring any intermediate bridges. The basic steps in error detection and recovery are as follows:

- An adapter driver suspects an error on the card when it receives some invalid values from one or more locations in its I/O or memory spaces.
- The driver then confirms the existence of the error by calling EEH kernel services. After the error state is confirmed, the slot is declared frozen.
- After the slot is frozen, all further activities to the card are suspended until the error is recovered. For example, new read/write requests are blocked or failed.
- The driver attempts to recover the slot by toggling the reset line. After three attempts to recover, the driver declares the slot unusable (or dead). If the slot is reset successfully, normal operations resume.

The key difference in the single-function and shared EEH models is that in the shared EEH model, there is a need for coordination among different driver instances controlling the same PE domain. For example, a PE domain can include a physical device on a single slot. The driver instances controlling each function on the device require coordination. In addition, there are steps in the recovery process that need to be carried out only once. So among all registered drivers in a shared EEH domain, one is chosen to be the *master*. The drivers follow a state machine. The EEH kernel services are implemented so that they present an EEH recovery state machine to the device drivers. It is the master driver's responsibility to drive the state machine. The section titled Shared EEH Programming Model, which follows, contains more details on how a master driver is determined. Many details are hidden from the device drivers for simplicity. Because the shared EEH model is more flexible and extensible, it is recommended for the new device drivers.

In the single-function model, the device drivers are responsible for driving their own error recovery. In other words, they are responsible for implementing their own state machine. Every time EEH recovery is extended in some way at the hardware or firmware level, there is probably a code and testing impact on the single-function implementations. As previously described, a single-function adapter should still use the shared EEH model. In that case, all the messages from the EEH kernel services are sent to just one driver instance.

## Shared EEH Programming Model

For the shared EEH programming model, the EEH kernel services present the following state machine to the drivers:

1. A slot starts out in the NORMAL state.
2. When an EEH event happens, the driver receives all F's from an MMIO load. Because all F's might be a legal value for a driver, the driver must call **eeh\_read\_slot\_state()** to confirm the event.
3. If **eeh\_read\_slot\_state()** finds the slot to be frozen, it broadcasts an EEH\_DD\_SUSPEND message to all registered drivers, and the slot state moves to SUSPEND. The kernel messages like this one are broadcast by invoking the callback routine sequentially. The messages are broadcast at INTIODONE priority.
4. When the drivers receive the EEH\_DD\_SUSPEND message, they can do one of the following:
  - a. Gather some debug data from the adapter and proceed to reset the slot.

Gathering the debug data is really an optional step in the recovery process, where a driver can choose to read certain registers on the adapter in an attempt to understand what caused the EEH event in the first place.

To gather the debug data, the drivers must enable PIO to the adapter. PIO is frozen when an EEH event occurs. To enable PIO:

- 1) The master driver must call **eeh\_enable\_pio()**. The master driver is picked by the EEH kernel services. It has the EEH\_MASTER flag set on the callback routine and is the last driver called

in the callback chain. This ensures that all other drivers in the shared EEH domain have finished the last step of the recovery and that the master driver can now proceed to the next step (such as enabling PIO).

When **eeh\_enable\_pio()** is called, an EEH\_DD\_DEBUG message is sent to the drivers indicating that PIO is enabled, and the slot state moves to DEBUG.

- 2) The drivers then gather the data.

**eeh\_enable\_pio()** can be called multiple times. Each time it is called, another EEH\_DD\_DEBUG message is broadcast.

- 3) When the drivers receive EEH\_DD\_SUSPEND or EEH\_DD\_DEBUG messages, they call **eeh\_slot\_error()** to create an AIX error log entry with hardware debug data. This step is required to figure out the reason for the EEH event.

- 4) The master driver must call **eeh\_reset\_slot()** to reset the slot. Only one driver calls reset because it is not necessary to reset the slot multiple times.

b. Proceed directly to reset the slot.

5. The reset line on the PCI bus is toggled with 100 ms delay between activate and deactivate to reset the slot. The delay is hidden from the device drivers and is enforced by the **eeh\_reset\_slot()** kernel service internally. The slot internally moves through the ACTIVATE and the DEACTIVATE states.
6. If there are any intermediate bridges present (such as a bridge on the adapter), at the end of a successful reset, EEH kernel services configures the bridge using **eeh\_configure\_bridge()** service. Kernel services also enforces a certain amount of delay between the deactivation of the reset line and the configuration of bridge.  
The device drivers do not need to call **eeh\_configure\_bridge()** directly.
7. If everything goes well, the EEH\_DD\_RESUME message is sent to the drivers indicating that the slot recovery is complete.
8. At this point, most drivers would have to reinitialize their adapters before starting normal operations again. Reinitialization typically requires a partial restore of the config space (such as the BARs and Cache Line). Determining the config space registers to be restore depends on the device.

**Note:** This is the usual recovery sequence. If any of the services fail, the EEH\_DD\_DEAD message is broadcast asking the drivers to mark their adapters unavailable (for example, the drivers might have to perform some cleanup work and mark their internal states appropriately). The master driver must call **eeh\_slot\_error()** to create an AIX error log and mark the adapter permanently unavailable.

There are two special scenarios that a driver developer needs to be aware of:

1. If a driver receives either an EEH\_DD\_SUSPEND or an EEH\_DD\_DEAD message, it can return an EEH\_BUSY return code from its callback routine instead of an EEH\_SUCC return code. If EEH kernel services receives an EEH\_BUSY message, EEH kernel services waits for some time and then calls the same driver again. This process continues until EEH kernel services receive a different return code. This process is repeated because some drivers need more time to cleanup before recovery can continue. Cleanup would include such activities like killing a kproc or notifying a user level app.
2. If **eeh\_enable\_dma()** and **eeh\_enable\_pio()** cannot succeed due to the platform state restrictions, the service returns an EEH\_FAIL return code followed by an EEH\_DD\_DEAD message unless you take action. To avoid receiving an EEH\_FAIL return code, the driver must supply an **EEH\_ENABLE\_NO\_SUPPORT\_RC** flag when **eeh\_init\_multifunc()** kernel services is initiated. If an **EEH\_ENABLE\_NO\_SUPPORT\_RC** flag is supplied, **eeh\_enable\_pio()** and **eeh\_enable\_dma()** return the EEH\_NO\_SUPPORT return code that indicates to the drivers that they cannot collect debug data but can continue with the next step in recovery. For more information, see **eeh\_read\_slot\_state**.

The EEH kernel services that you can use are listed in the following table:

**Note:** **eeh\_init()** and **eeh\_init\_multifunc()** are the only exported kernel services. All other kernel services are called using function pointers in the **eeh\_handle** kernel service.



Kernel Service	Single Function	Shared EEH	Process Environment	Interrupt Environment
eeh_init	Y	N	Y	N
eeh_init_multifunc	N	Y	Y	N
eeh_clear	Y	Y	Y	N
eeh_read_slot_state	Y	Y	Y	Y
eeh_enable_pio	Y	Y	Y	Y
eeh_enable_dma	Y	Y	Y	Y
eeh_enable_slot	Y	N	Y	Y
eeh_disable_slot	Y	N	Y	Y
eeh_reset_slot	Y	Y	Y	Y
eeh_slot_error	Y	Y	Y	Y
eeh_broadcast	N	Y	Y	Y

## Callback Routine

The **(\*callback\_ptr)()** function prototype is defined as:

```
long eeh_callback(
    unsigned long long cmd, /* EEH messages */
    void *arg,             /* Pointer to dd defined argument */
    unsigned long flag)    /* DD defined flag */
```

- *cmd* – contains a kernel and driver message
- *arg* – is a cookie to a target device driver that is usually a pointer to the adapter structure
- *flag* argument can be either just EEH\_MASTER or EEH\_MASTER ORed with EEH\_DD\_PIO\_ENABLED

### EEH\_MASTER

Indicates that the target device driver is the EEH\_MASTER.

### EEH\_DD\_PIO\_ENABLED

Set only with the EEH\_DD\_DEBUG message to indicate that the PIO is enabled.

When **eeh\_init\_multifunc()** is called, the callback routines are registered. When **eeh\_clear()** is called the callback routines are unregistered. The callback routines are necessary for EEH kernel services recovery. They coordinate shared EEH domain driver instances. For more information on how this coordination is done, see “Enhanced I/O Error Handling Kernel Services” on page 44.

The shared EEH domain drivers are expected to handle the following EEH kernel services messages:

### EEH\_DD\_SUSPEND

Notifies all the device drivers on a slot that an EEH kernel services event occurred. The slot is either frozen or temporarily unavailable. Because an EEH kernel services event occurred, the device drivers suspend operations. Then, the EEH\_MASTER driver either enables PIO or resets the slot.

### EEH\_DD\_DEBUG

Notifies all drivers on a slot that they can now gather debug data from the devices. The device drivers log errors by calling the **eeh\_slot\_error()** function and passing in the gathered debug data. This message is sent when the EEH\_MASTER calls the **eeh\_enable\_pio()** function. On the callback routine, the flag argument is set to EEH\_DD\_PIO\_ENABLED.

### EEH\_DD\_DEAD

Notifies all drivers on a slot that the slot reached an unrecoverable state and the slot is no longer usable. This message is sent anytime EEH kernel services fail because of hardware or firmware

problems. This message is also broadcast when a driver calls the **eeh\_slot\_error()** function with the flag set to **EEH\_RESET\_PERM**. The device drivers usually perform necessary cleanup and mark the adapter as permanently unavailable.

#### **EEH\_DD\_RESUME**

Notifies all drivers on a slot that the EEH kernel services event was recovered successfully and that the callback routines can now resume normal operation. This message is sent at the end of a successful toggle of reset line and optional bridge configuration (for example, the bridge on the adapter). The device drivers must usually re-initialize their adapters before normal operation can begin again.

The device drivers define their own messages based on the contents of the `sys/eeh.h` file.

The `eeh_callback()` functions are scheduled to start sequentially at **INTIODONE** priority. They are not started in any specific order. For more information, see `eeh_broadcast`.

---

## **Block I/O Buffer Cache Kernel Services: Overview**

The Block I/O Buffer Cache services are provided to support user access to device drivers through block I/O special files. This access is required by the operating system file system for mounts and other limited activity, as well as for compatibility services required when other file systems are installed on these kinds of systems. These services are not used by the operating system's JFS (journal file system), NFS (Network File System), or CDRFS (CD-ROM file system) when processing standard file I/O data. Instead they use the virtual memory manager and pager to manage the system's memory pages as a buffer cache.

For compatibility support of other file systems and block special file support, the buffer cache services serve two important purposes:

- They ensure that multiple processes accessing the same block of the same device in multiprogrammed fashion maintain a consistent view of the data in the block.
- They increase the efficiency of the system by keeping in-memory copies of blocks that are frequently accessed.

The Buffer Cache services use the **buf** structure or buffer header as their main data-tracking mechanism. Each buffer header contains a pair of pointers that maintains a doubly-linked list of buffers associated with a particular block device. An additional pair of pointers maintain a doubly-linked list of blocks available for use again on another operation. Buffers that have I/O in progress or that are busy for other purposes do not appear in this available list.

Kernel buffers are discussed in more detail in Introduction to Kernel Buffers in *AIX Version 6.1 Technical Reference: Kernel and Subsystems Volume 1*.

See "Block I/O Kernel Services" on page 41 for a list of these services.

## **Managing the Buffer Cache**

Fourteen kernel services provide management of this block I/O buffer cache mechanism. The **getblk** kernel service allocates a buffer header and a free buffer from the buffer pool. Given a device and block number, the **getblk** and **bread** kernel services both return a pointer to a buffer header for the block. But the **bread** service is guaranteed to return a buffer actually containing a current data for the block. In contrast, the **getblk** service returns a buffer that contains the data in the block only if it is already in memory.

In either case, the buffer and the corresponding device block are made busy. Other processes attempting to access the buffer must wait until it becomes free. The **getblk** service is used when:

- A block is about to be rewritten totally.

- Its previous contents are not useful.
- No other processes should be allowed to access it until the new data has been placed into it.

The **breada** kernel service is used to perform read-ahead I/O and is similar to the **bread** service except that an additional parameter specifies the number of the block on the same device to be read asynchronously after the requested block is available. The **brelse** kernel service makes the specified buffer available again to other processes.

## Using the Buffer Cache write Services

There are three slightly different write routines. All of them take a buffer pointer as a parameter and all logically release the buffer by placing it on the free list. The **bwrite** service puts the buffer on the appropriate device queue by calling the device's strategy routine. The **bwrite** service then waits for I/O completion and sets the caller's error flag, if required. This service is used when the caller wants to be sure that I/O takes place synchronously, so that any errors can be handled immediately.

The **bawrite** service is an asynchronous version of the **bwrite** service and does not wait for I/O completion. This service is normally used when the overlap of processing and device I/O activity is desired.

The **bdwrite** service does not start any I/O operations, but marks the buffer as a delayed write and releases it to the free list. Later, when the buffer is obtained from the free list and found to contain data from some other block, the data is written out to the correct device before the buffer is used. The **bdwrite** service is used when it is undetermined if the write is needed immediately.

For example, the **bdwrite** service is called when the last byte of the write operation associated with a block special file falls short of the end of a block. The **bdwrite** service is called on the assumption that another write will soon occur that will use the same block again. On the other hand, as the end of a block is passed, the **bawrite** service is called, because it is assumed the block will not be accessed again soon. Therefore, the I/O processing can be started as soon as possible.

Note that the **getblk** and **bread** services dedicated the specified block to the caller while making other processes wait, whereas the **brelse**, **bwrite**, **bawrite**, or **bdwrite** services must eventually be called to free the block for use by other processes.

---

## Understanding Interrupts

Each hardware interrupt has an interrupt level, trigger, and interrupt priority.

### Interrupt Level

The interrupt level defines the source of the interrupt and is often referred to as the *interrupt source*. There are basically two types of interrupt levels: system and bus. The bus interrupts are generated by the devices on the buses (such as PCI, ISA, VDEVICE, and PCI-E). Examples of system interrupts are the timer and Environmental and Power Off Warning (EPOW).

The interrupt level of a bus or system interrupt is one of the resources managed by the respective configuration methods.

### Interrupt Trigger

There are two types of trigger mechanisms, level-triggered interrupts and edge-triggered interrupts. All ISA and VDEVICE interrupts are edge-triggered. PCI/PCIX and PCI-E buses define two types of interrupts, Level Signalled Interrupts (LSI) and Message Signalled Interrupts (MSI). LSIs are level-triggered, and MSIs are edge-triggered. PCI/PCIX device drivers in AIX must handle only level-triggered interrupts even though edge-triggered interrupts using MSIs are supported by PCIX. Similarly, even though PCI-E supports LSI

interrupts, all PCI-E interrupts are MSIs only in AIX, and are therefore edge-triggered. Consequently, all PCI-E device drivers for AIX are required to handle only edge-triggered interrupts.

A key difference between the edge-triggered and level-triggered interrupts is interrupt sharing. Level-triggered interrupts can be shared. Edge-triggered interrupts cannot be shared. Because they cannot be shared, edge-triggered interrupt handlers should pass the INTR\_EDGE flag on the `i_init()` kernel service.

Another difference between the edge-triggered and level-triggered interrupts is in issuing the End of Interrupt (EOI). For level-triggered interrupts, the AIX kernel issues the EOI. For ISA edge-triggered interrupts, the AIX kernel also issues EOI. However, for the VDEVICE and PCI-E edge-triggered interrupts, the device driver must issue the EOI before returning from the interrupt handler. The VDEVICE and PCI-E drivers can use the `i_eoi()` kernel service to accomplish this.

**Note:** During the processing of `i_eoi()`, there is a brief period of time in which a newly arrived interrupt could also be issued an EOI. Therefore, it is necessary to check for additional work between a call to `i_eoi()` and the return from the interrupt handler.

## Interrupt Priorities

The interrupt priority defines which of a set of pending interrupts is serviced first. INTMAX is the most favored interrupt priority and INTBASE is the least favored interrupt priority. The interrupt priorities for bus interrupts range from INTCLASS0 to INTCLASS3. The rest of the interrupt priorities are reserved for the base kernel. Interrupts that cannot be serviced within the time limits specified for bus interrupts qualify as off-level interrupts.

A device's interrupt priority is selected based on two criteria: its maximum interrupt *latency* requirements and the device driver's interrupt *execution time*. The interrupt latency requirement is the maximum time within which an interrupt must be serviced. (If it is not serviced in this time, some event is lost or performance is seriously degraded.) The interrupt execution time is the number of machine cycles required by the device driver to service the interrupt. Interrupts with a short interrupt latency time must have a short interrupt service time.

The general rule for interrupt service times is based on the following interrupt priority table:

Priority	Service Time (machine cycles)
INTCLASS0	200 cycles
INTCLASS1	400 cycles
INTCLASS2	600 cycles
INTCLASS3	800 cycles

The valid interrupt priorities are defined in the `/usr/include/sys/intr.h` file.

See "Interrupt Management Kernel Services" on page 42 for a list of these services.

---

## Understanding DMA Transfers

AIX DMA support deals with the issues of direct memory access by I/O devices to and from system memory. The programming framework supports common I/O buses (such as PCI, industry standard architecture (ISA), and PCIe) and is easily extensible to additional bus types in the future. The framework supports 64-bit addressability, and also allows for mappings from 32-bit devices to 64-bit addresses to be hidden from the devices and their drivers.

## DMA Programming Model

This is the basic DMA programming model. It is completely independent of:

- System hardware
- LPAR mode or non-LPAR mode
- 32-bit bus/devices or 64-bit bus or devices
- 32-bit kernel or 64-bit kernel

A device driver allocates and initializes DMA-related resources with the **d\_map\_init** service and the **d\_map\_init\_ext** service and frees the resources with the **d\_map\_clear** service. Each time a DMA mapping needs to be established, the driver calls **d\_map\_page** or **d\_map\_list** service.

The **d\_map\_page** and **d\_map\_list** services map DMA buffers in the bus memory. In other words, given a set of DMA buffer addresses, a corresponding set of bus addresses is returned to the driver. The driver programs its device with the bus addresses and sets it up to start the DMA. When the DMA is complete:

- The device generates an interrupt that is handled by the driver.
- If no more DMA will be done to the buffers, the driver unmaps the DMA buffers with the **d\_unmap\_page** or **d\_unmap\_list** services.

## Data Structures

The **d\_map\_init** service and the **d\_map\_init\_ext** service return a **d\_handle\_t** structure to the caller upon a successful completion. Only the **d\_map\_init** service and the **d\_map\_init\_ext** service are the exported kernel services. All other DMA kernel services are called through the function pointers in the **d\_handle\_t** structure (see the **sys/dma.h** system file).

```
d_handle {
    uint id; /* identifier for this device */
    uint flags; /* device capabilities */
#ifdef __64BIT_KERNEL
    /* pointer to d_map_page routine */
    int (*d_map_page)(d_handle_t,int,caddr_t, ulong *, struct xmem
*);
    /* pointer to d_unmap_page routine */
    void (*d_unmap_page)(d_handle_t, ulong *);
    /* pointer to d_map_list routine */
    int (*d_map_list)(d_handle_t, int, int, dio_t, dio_t);
    /* pointer to d_unmap_list routine */
    void (*d_unmap_list)(d_handle_t, dio_t);
    /* pointer to d_map_slave routine */
    int (*d_map_slave)(d_handle_t, int, int, dio_t, uint);
    /* pointer to d_unmap_slave routine */
    int (*d_unmap_slave)(d_handle_t);
    /* pointer to d_map_disable routine */
    int (*d_map_disable)(d_handle_t);
    /* pointer to d_map_enable routine */
    int (*d_map_enable)(d_handle_t);
    /* pointer to d_map_clear routine */
    void (*d_map_clear)(d_handle_t);
    /* pointer to d_sync_mem routine */
    int (*d_sync_mem)(d_handle_t, dio_t);
#else
    int (*d_map_page)(); /* pointer to d_map_page routine */
    void (*d_unmap_page)(); /* pointer to d_unmap_page routine */
    int (*d_map_list)(); /* pointer to d_map_list routine */
    void (*d_unmap_list)(); /* pointer to d_unmap_list routine */
    int (*d_map_slave)(); /* pointer to d_map_slave routine */
    int (*d_unmap_slave)(); /* pointer to d_unmap_slave routine */
    int (*d_map_disable)(); /* pointer to d_map_disable routine */
    int (*d_map_enable)(); /* pointer to d_map_enable routine */
    void (*d_map_clear)(); /* pointer to d_map_clear routine */
    int (*d_sync_mem)(); /* pointer to d_sync_mem routine */
#endif
};
```

```

#endif
        int bid; /* bus id passed to d_map_init or d_map_init_ext*/
        void *bus_sys_xlate_ptr; /* pointer to dma bus to system
        translation information */
        uint reserved1; /* padding */
        uint reserved2; /* padding */
        uint reserved3; /* padding */
};

```

The following **dio** and **d\_iovec** structures are used to define the scatter and gather lists used by the **d\_map\_list**, **d\_unmap\_list**, and **d\_map\_slave** services (see the **sys/dma.h** system file).

```

struct dio {
    int32long64_t total_iovecs; /* total available iovec entries */
    int32long64_t used_iovecs; /* number of used iovecs */
    int32long64_t bytes_done; /* count of bytes processed */
    int32long64_t resid_iov; /* number of iovec that couldn't be */
    /* fully mapped due to NORES,DIOFULL*/
    /* vec =&dvec [resid_iov] */
    struct d_iovec *dvec; /* pointer to list of d_iovecs */
};

struct d_iovec {
    caddr_t iov_base; /* base memory address */
    int32long64_t iov_len; /* length of transfer for this area */
    struct xmem *xmp; /* cross memory pointer for this address*/
};

```

The following **dio\_64** and **d\_iovec\_64** structures are used to define the scatter and gather lists used by the **d\_map\_list** and **d\_unmap\_list** services when the **DMA\_ENABLE\_64** flag is set on the **d\_map\_init** call or the **d\_map\_init\_ext** call. These are not used under the 64-bit kernel and kernel extension environment because the **dio** and **d\_iovec** data structures are naturally 64-bit capable in that environment. (For more information, see the **sys/dma.h** system file.)

```

struct dio_64 {
    int total_iovecs; /* total available iovec entries */
    int used_iovecs; /* number of used iovecs */
    int bytes_done; /* count of bytes processed */
    int resid_iov; /* number of iovec that couldn't be */
    /* fully mapped due to NORES,DIOFULL*/
    /* vec = &dvec [resid_iov] */
    struct d_iovec_64 *dvec; /* pointer to list of d_iovecs */
};

struct d_iovec_64 {
    unsigned long long iov_base; /* base memory address */
    int iov_len; /* length of transfer for this area */
    struct xmem *xmp; /* cross memory pointer for this address*/
}

```

The following macros are provided in **sys/dma.h** for device drivers in order to call the DMA kernel services cleanly:

```

#define D_MAP_INIT(bid, flags, bus_flags, channel) \
    d_map_init(bid, flags, bus_flags, channel)

#define D_MAP_INIT_EXT(dma_input, info_size, handle_ptr) \
    d_map_init_ext(dma_input, info_size, handle_ptr)
#define D_MAP_CLEAR(handle) \
    (handle->d_map_clear)(handle)

#define D_MAP_PAGE(handle, flags, baddr, busaddr, xmp) \
    (handle->d_map_page)(handle, flags, baddr, busaddr, xmp)

#define D_UNMAP_PAGE(handle, bus_addr) \
    if (handle->d_unmap_page != NULL) (handle->d_unmap_page)(handle, bus_addr)

#define D_MAP_LIST(handle, flags, minxfer, virt_list, bus_list) \

```

```

(handle->d_map_list)(handle, flags, minxfer, virt_list,bus_list)

#define D_UNMAP_LIST(handle, bus_list) \
    if (handle->d_unmap_list != NULL)(handle->d_unmap_list)(handle, bus_list)

#define D_MAP_SLAVE(handle, flags, minxfer, vlist, chan_flags) \
    (handle->d_map_slave)(handle, flags, minxfer, vlist, chan_flags)

#define D_UNMAP_SLAVE(handle) \
    (handle->d_unmap_slave != NULL) ? \
    (handle->d_unmap_slave)(handle) : DMA_SUCC

#define D_MAP_DISABLE(handle) (handle->d_map_disable)(handle)

#define D_MAP_ENABLE(handle) (handle->d_map_enable)(handle)

#define D_SYNC_MEM(handle, bus_list) \
    (handle->d_sync_mem != NULL) ? \
    (handle->d_sync_mem)(handle, bus_list) : DMA_SUCC
#define D_MAP_ATTR(handle, cmd, attr, attr_size) \
    (handle->d_map_attr)(handle, cmd, attr, attr_size)

```

## The d\_map Return Code Map

The following table describes the possible return codes and requirements for the **d\_map** interfaces that map memory for DMA:

Table 1. Return codes and requirements for the **d\_map** interfaces

Return Codes	d_map_page	d_map_list	d_map_slave
DMA_SUCC	Page mapped successfully, busaddr contains the mapped bus address. <b>d_unmap_page</b> must be called to free any resources associated with the mapping.	List mapped successfully. bus_list describes list of mapped bus addresses. <b>d_unmap_list</b> must be called to free any resources associated with the mapping.	List mapped successfully. Slave DMA Controller initialized for the required transfer. <b>d_unmap_slave</b> must be called to free any resources associated with the mapping.
DMA_NORES	Not enough resources to map the page. No mapping is performed. <b>d_unmap_page</b> must not be called.	Not enough resources to map the entire list. A partial mapping is possible. <b>d_unmap_list</b> must be called to free any resources associated with the mapping.	Not enough resources to map the entire list. A partial mapping is possible. <b>d_unmap_slave</b> must be called to free any resources associated with the mapping.
DMA_NOACC	No access to the page. No mapping is performed. <b>d_unmap_page</b> must not be called.	No access to a page in the list. No mapping is performed. <b>d_unmap_list</b> must not be called.	No access to a page in the list. No mapping is performed. <b>d_unmap_slave</b> must not be called.
DMA_DIOFULL	Does not apply.	bus_list is exhausted. Successful partial mapping as indicated. <b>d_unmap_list</b> must be called when finished with the partial mapping.	Does not apply.
DMA_BAD_MODE	Does not apply.	Does not apply.	Requested channel modes are not supported.

## Using the dio Structure

A device driver can use the **dio** structure in many ways. It can be used to:

- Pass a list of virtual addresses and lengths of buffers to the **d\_map\_list** and **d\_map\_slave** services.
- Receive the resulting list of bus addresses (**d\_map\_list** only) for use by the device in the data transfer.

**Note:** The driver does not need a **dio** bus list for calls to **d\_map\_slave** because the address generation for slaves is hidden.

Typically, a device driver provides a **dio** structure that contains only one virtual buffer and one length in the list. If the virtual buffer length spans many pages, the bus address list contains multiple entries that reflect the physical locations of the virtually contiguous buffer. The driver can provide multiple virtual buffers in the virtual list. The driver can then place many buffer requests in one I/O operation.

The device driver is responsible for allocating the storage for all the **dio** lists it needs. For more information, see the **DIO\_INIT** and **DIO\_FREE** macros in the `sys/dma.h` header file. The driver must have at least two **dio** structures. One is needed for passing in the virtual list. Another is needed to accept the resulting bus list. The driver can have many **dio** lists if it plans to have multiple outstanding I/O commands to its device. The length of each list is dependent on the use of the device and driver. The virtual list needs as many elements as the device could place in one operation at the same time. A formula for estimating how many elements the bus address list needs is the sum of each of the virtual buffers lengths divided by page size plus 2. Or,

```
sum [i=0 to n] ((vlist[i].length / PSIZE) + 2).
```

This formula handles a worst-case situation. For a contiguous virtual buffer that spans multiple pages, each physical page is discontiguous, and neither the starting nor ending addresses are page-aligned.

If the **d\_map\_list** service runs out of space when filling in the **dio** bus list, a `DMA_DIOFULL` error is returned to the device driver and the `bytes_done` field of the **dio** virtual list is set to the number of bytes successfully mapped in the bus list. This byte count is guaranteed to be a multiple of the `minxfer` field provided to the **d\_map\_list** or **d\_map\_slave** services. Also, the `resid_iov` field of the virtual list is set to the index of the first **d\_iovec** entry that represents the remainder of `iovecs` that could not be mapped.

The device driver can:

- Initiate a partial transfer on its device and leave the remainder on its device queue.  
If the driver chooses not to initiate the partial transfer, it must still make a call to **d\_unmap\_list** to undo the partial mapping.
- Make another call to the **d\_map\_list** with new **dio** lists for the remainder and setup its device for the full transfer that was originally intended.

If **d\_map\_list** or **d\_map\_slave** encounter an access violation on a page within the virtual list, then a `DMA_NOACC` error is returned to the device driver and the `bytes_done` field of the **dio** virtual list is set to the number of bytes that preceded the faulting `iovec`. In this case, the `resid_iov` field is set to the index of the **d\_iovec** entry that encountered the violation. From this information, the driver can determine which virtual buffer contained the faulting page and fail that request back to the originator.

**Note:** If the `DMA_NOACC` error is returned, the `bytes_done` count is not guaranteed to be a multiple of the `minxfer` field provided to the **d\_map\_list** or **d\_map\_slave** services, and no partial mapping is done. For slaves, setup of the address generation hardware is not done. For masters, the bus list is undefined. If the driver desires a partial transfer, it must make another call to the mapping service with the **dio** list adjusted to not include the faulting buffer.

If either the **d\_map\_list** or **d\_map\_slave** services run out of resources when mapping a transfer, a `DMA_NORES` error is returned to the device driver. In this case, the `bytes_done` field of the **dio** virtual list is set to the number of bytes that were successfully mapped in the bus list. This byte count is guaranteed



to be a multiple of the **minxfer** field provided to the **d\_map\_list** or **d\_map\_slave** services. Also, the **resid\_iov** field of the virtual list is set to the index of the first **d\_iovec** of the remaining iovecs that could not be mapped. The device driver can:

- Initiate a partial transfer on its device and leave the remainder on its device queue  
If the driver chooses not to initiate the partial transfer, it still must make a call to **d\_unmap\_list** or **d\_unmap\_slave** (for slaves) to undo the partial mapping.
- Choose to leave the entire request on its device queue and wait for resources to free up (for example, after a device interrupt from a previous operation).

**Note:** If the **DMA\_ENABLE\_64** flag is indicated on the **d\_map\_init** call or the **d\_map\_init\_ext** call, the programming model is the same with one exception. The **dio\_64** and **d\_iovec\_64** structures are used in addition to 64-bit address fields on **d\_map\_page** and **d\_unmap\_page** calls.

## Fields of dio

The only field of the bus list that a device driver modifies is the **total\_iovecs** field to indicate how many elements are available in the list. The device driver never changes any of the other fields in the bus list. The device driver uses the bus list to set up its device for the transfer. The bus list is provided to the **d\_unmap\_list** service to unmap the transfer. The **d\_map\_list** service sets the **used\_iovecs** field to indicate how many elements it filled out. The device driver sets up all of the fields in the virtual list except for the **bytes\_done** and **resid\_iov** fields. These fields are set by the mapping service.

## Using DMA\_CONTIGUOUS

The **DMA\_CONTIGUOUS** flag in the **d\_map\_init** service or the **d\_map\_init\_ext** service is the preferred way for the drivers to ask for contiguous bus addresses. The other way is the old model of drivers explicitly using **rmalloc()** to guarantee contiguous allocation during boot. However, with the advent of PCI Hot Plug devices, the **rmalloc** reservation does not add a device after boot. If a PowerPC driver determines that the device was dynamically added, the driver can use the **DMA\_CONTIGUOUS** flag to ensure that a contiguous list of bus addresses is generated because no prior resources were reserved with **rmalloc**.

## Using DMA\_NO\_ZERO\_ADDR

**DMA\_NO\_ZERO\_ADDR** is supplied on the **d\_map\_init** service or the **d\_map\_init\_ext** service to prevent **d\_map\_page** and **d\_map\_list** from giving out bus address zero to this **d\_handle**. Because many off-the-shelf PCI devices are not tested for bus address of zero, such devices might not work. Striking out bus address 0 causes a driver's mappable memory to shrink by one I/O page (4 KB). On some systems, using the flag would cause the **d\_map\_init** service or the **d\_map\_init\_ext** service to fail even if there is not an error condition. In such a case, the driver should call the **d\_map\_init** service or the **d\_map\_init\_ext** service without the flag and then check the bus address to see whether zero falls in its range of addresses. The driver can do this by mapping all of its range and checking for address 0. Such a check should be done at the driver initialization time. If bus address 0 is assigned to the driver, it can leave it mapped for the life of the driver and unmap all other addresses. This guarantees that address 0 is not assigned to it again.

## Using DMA\_MAXMIN\_MAPSPACE

The **DMA\_MAXMIN\_MAPSPACE** flag is supplied on the **d\_map\_init** service or the **d\_map\_init\_ext** service to prevent the **d\_map\_init** service or the **d\_map\_init\_ext** service from allocating more DMA space for the DMA handle than what the device driver requests.

The device drivers request a specific amount of DMA space by specifying a **DMA\_MAXMIN\_\*** flag on the **d\_map\_init** service or by specifying the **d\_info\_t.di\_max\_mapspace** flag on the **d\_map\_init\_ext** service. Because some device drivers support non-page-aligned DMA transfers of the specified maximum DMA space, by default the **d\_map\_init** service or the **d\_map\_init\_ext** service allocates at least one additional page of DMA space.

Device drivers that use the **DMA\_MAXMIN\_MAPSPACE** flag cannot support non-page-aligned DMA transfers of the specified maximum DMA space. The **DMA\_MAXMIN\_MAPSPACE** flag indicates that the value of the maximum DMA space represents the amount of mappable address space the device driver requires, rather than the maximum transfer value.

## Using DMA\_INIT\_STMAP\_SUPPORT

To indicate that the driver supports short-term mapping and long-term mapping, you can specify the *flags* parameter of the following kernel services with the values in the following table.

Table 2. Short-term mapping and long-term mapping support

Kernel service	The value of the <i>flags</i> parameter
<b>d_map_init_ext</b>	<b>DMA_INIT_STMAP_SUPPORT</b>
<b>d_map_page</b>	<b>DMA_STMAP</b>
<b>d_map_list</b>	<b>DMA_STMAP</b>

## Sample pseudo-code for the PCI drivers

dd\_initialization:

```
determine bus type for device from configuration information
determine 64 vs. 32-bit capabilities from configuration information
```

```
di_info.di_bid = bid
di_info.di_flags = DMA_MASTER | flags (except DMA_MAXMIN_* flag);
di_info.di_bus_flags = 0;
di_info.di_channel = 0;
di_info.di_min_mapspace = min_mapspace;
di_info.di_des_mapspace = DMA_MAXMIN_* value;
di_info.di_max_mapspace = max_dma_space ;
call "rc = D_MAP_INIT_EXT(&di_info, sizeof(di_info), handle)"
if handle == DMA_FAIL OR rc != DMA_SUCC
    could not configure
```

dd\_start\_io:

```
if single page or less transfer
    call "result = D_MAP_PAGE(handle, baddr,busaddr, xmem)"
    if result == DMA_NORES
        no resources, leave request on device queue
    else if result == DMA_NOACC
        no access to page, fail request
    else
        program device for transfer using busaddr
else
    create dio list of virtual addresses involved in transfer
    call "result = D_MAP_LIST(handle, flags, minxfer, list, blist)"
    if result == DMA_NORES
        not enough resource, either initiate partial transfer
        and leave remainder on queue or leave entire request
        on the queue and call d_unmap_list to unmap the
        partial transfer.
    else if result == DMA_NOACC
        use bytes_done to pinpoint failing buffer and
        fail corresponding request adjust virtual list and
        call d_map_list again
    else if result == DMA_DIOFULL
        ran out of space in blist. either initiate partial
        transfer and leave remainder on queue or leave entire
        request on the queue and call d_unmap_list to
        unmap the partial transfer.
else
```

```

        program device for scatter/gather transfer using blist

dd_finish_io:
    if single page or less transfer
        call "D_UNMAP_PAGE(handle, busaddr)"
    else
        call "D_UNMAP_LIST(handle, blist)"

dd_unconfigure:
    call "D_MAP_CLEAR(handle)"

```

## Sample Pseudo-code for the ISA Slave drivers

```

dd_initialization:
    determine bus type for device from configuration information
    call "handle = D_MAP_INIT(bid, DMA_SLAVE, bus_flags, channel)"
    if handle == DMA_FAIL
        could not configure
    else
        call "D_MAP_ENABLE(handle)" (if necessary)

dd_start_io:
    create dio list of virtual addresses involved in transfer
    call "result = D_MAP_SLAVE(handle, flags, minxfer, vlist,
chan_flags)"
    if result == DMA_NORES
        not enough resource, either initiate partial transfer
        and leave remainder on queue or leave entire request
        on the queue and call d_unmap_slave to unmap the
        partial transfer.
    else if result == DMA_NOACC
        use bytes_done to pinpoint failing buffer and
        fail corresponding request
        adjust virtual list and call d_map_slave again
    else
        program device to initiate transfer

dd_finish_io:
    call "error = D_UNMAP_SLAVE(handle)"
    if error
        log
        retry, or fail

dd_unconfigure:
    call "D_MAP_DISABLE(handle)" (if necessary)
    call "D_MAP_CLEAR(handle)"

```

## Page Protection Checking and Enforcement

Page protection checking is performed by the `d_map_page`, `d_map_list`, and `d_map_slave` services for each page of a requested transfer. If the intended direction of a transfer is from the device to the memory, the page access permissions must support writing to the page. If the intended direction of a transfer is from the memory to the device, the page access permissions only require reading from the page. In the case of a protection violation, a `DMA_NOACC` return code is returned from the services in the form of an error code and no mapping for the DMA transfer is performed.

The `DMA_BYPASS` flag lets a device driver bypass the access checking functionality of these services. This should only be used for global system buffers such as mbufs or other command, control, and status

buffers used by a device driver. Also, the DMA buffers *must* be pinned before the DMA transfer begins and can only be unpinned after the DMA transfer is complete.

## Short term mapping

Specifying the **DMA\_STMAP** flag on both the **d\_map\_page** service and the **d\_map\_list** service is the preferred method for a caller to indicate whether the mapping is for a short term. Indicating that a mapping is for a short term may allow the driver to map additional memory beyond the amount of I/O mappable memory that is reserved by the **d\_map\_init\_ext** kernel service. The caller must pass the **DMA\_INIT\_STMAP\_SUPPORT** flag on the **d\_map\_init\_ext** service to recognize the **DMA\_STMAP** flag.

## A comparison of PCI and ISA devices

The ISA bus has the following unique concepts that do not apply to the PCI bus:

- Enabling and disabling a DMA channel applies only to the ISA bus and devices. Therefore, **d\_map\_enable** and **d\_map\_disable** services cannot be used by PCI device drivers.
- Master and slave devices are not applicable to the PCI bus. On a PCI bus, every device acts as master.

Starting with AIX 5.2, only ISA slave devices are supported (ISA masters are not supported). For such ISA slave devices, the PCI-to-ISA bridge acts as the PCI master and initiates DMA on behalf of the ISA slave devices. Because the PCI devices are always master, **d\_map\_slave** and **d\_unmap\_slave** services cannot be used by PCI device drivers. By the same token, the **DMA\_SLAVE** flag cannot be supplied on the **d\_map\_init** service or the **d\_map\_init\_ext** service by a PCI device driver. If **DMA\_SLAVE** is used by a PCI driver, the **d\_map\_init** service or the **d\_map\_init\_ext** service returns **DMA\_FAIL**.

## d\_align and d\_roundup

The **d\_align** service (provided in **libsys.a**) returns the alignment value required for starting a buffer on a processor cache line boundary. The **d\_roundup** service (also provided in **libsys.a**) can be used to round the desired DMA buffer length up to a value that is an integer number of cache lines. These two services use buffers for DMA alignment on a cache line boundary. The buffers are also used to allocate DMA in whole multiples of the cache line size so that the buffer is not split across processor cache lines. This reduces the possibility of consistency problems because of DMA and also minimizes the number of cache lines that must be flushed or invalidated when used for DMA. For example, these services can be used to provide alignment as follows:

```
align = d_align();
buffer_length = d_roundup(required_length);
buf_ptr = xmalloc(buffer_length, align, kernel_heap);
```

---

## Kernel Extension and Device Driver Management Services

The kernel provides a set of program and device driver management services. These services include kernel extension loading and unloading services and device driver binding services. Services that allow kernel extensions to be notified of base kernel configuration changes, user-mode exceptions, and process state changes are also provided.

The following information is provided to assist you in learning more about kernel services:

- “Kernel Extension Loading and Unloading Services”
- “Other Kernel Extension and Device Driver Management Services” on page 59
- “List of Kernel Extension and Device Driver Management Kernel Services” on page 59

## Kernel Extension Loading and Unloading Services

The **kmod\_load**, **kmod\_unload**, and **kmod\_entrpt** services provide kernel extension loading, unloading, and query services. User-mode programs and kernel processes can use the **sysconfig** subroutine to invoke the **kmod\_load** and **kmod\_unload** services. The **kmod\_entrpt** service returns a pointer to a kernel extension’s entry point.

The **kmod\_load**, **kmod\_unload** services can be used to dynamically alter the set of routines loaded into the kernel based on system configuration and application demand. Subsystems and device drivers can use these services to load large, seldom-used routines on demand.

## Other Kernel Extension and Device Driver Management Services

The device driver binding services are **devswadd**, **devswdel**, **devswchg**, and **devswqry**. The **devswadd**, **devswdel**, and **devswchg** services are used to add, remove, or modify device driver entries in the dynamically-managed device switch table. The **devswqry** service is used to obtain information about a particular device switch table entry.

Some kernel extensions might be sensitive to the settings of base kernel runtime configurable parameters that are found in the **var** structure defined in the **/usr/include/sys/var.h** file. These parameters can be set automatically during system boot or at runtime by a privileged user. Kernel extensions can register or unregister a configuration notification routine with the **cfgnadd** and **cfgndel** kernel services. Each time the **sysconfig** subroutine is used to change base kernel tunable parameters found in the **var** structure, each registered configuration notification routine is called.

The **prochadd** and **prochdel** kernel services allow kernel extensions to be notified when any process in the system has a state transition, such as being created, exiting, or being swapped in or swapped out.

The **uexadd** and **uexdel** kernel services give kernel extensions the capability to intercept user-mode exceptions. A user-mode exception handler can use this capability to dynamically reassign access to single-use resources or to clean up after some particular user-mode error. The associated **uexblock** and **uexcLEAR** services can be used by these handlers to block and resume process execution when handling these exceptions.

The **pio\_assist** and **getexcept** kernel services are used by device drivers to obtain detailed information about exceptions that occur during I/O bus access. The **getexcept** service can also be used by any exception handler requiring more information about an exception that has occurred. The **selreg** kernel service is used by file select operations to register unsatisfied asynchronous poll or select event requests with the kernel. The **selnotify** kernel service provides the same functionality as the **selwakeUP** service found on other operating systems.

The **iostadd** and **iostdel** services are used by tty and disk device drivers to register device activity reporting structures to be used by the **iostat** and **vmstat** commands.

The **getuerror** and **setuerror** services allow kernel extensions to read or set the **ut\_error** field for the current thread. This field can be used to pass an error code from a system call function to an application program, because kernel extensions do not have direct access to the application's **errno** variable.

## List of Kernel Extension and Device Driver Management Kernel Services

The Kernel Program and Device Driver Management kernel services are:

<b>cfgnadd</b>	Registers a notification routine to be called when system-configurable variables are changed.
<b>cfgndel</b>	Removes a notification routine for receiving broadcasts of changes to system configurable variables.
<b>devdump</b>	Calls a device driver dump-to-device routine.
<b>devstrat</b>	Calls a block device driver's strategy routine.
<b>devswadd</b>	Adds a device entry to the device switch table.
<b>devswchg</b>	Alters a device switch entry point in the device switch table.
<b>devswdel</b>	Deletes a device driver entry from the device switch table.
<b>devswqry</b>	Checks the status of a device switch entry in the device switch table.
<b>getexcept</b>	Allows kernel exception handlers to retrieve additional exception information.
<b>getuerror</b>	Allows kernel extensions to read the <b>ut_error</b> field for the current thread.

<b>iostatt</b>	Registers an I/O statistics structure used for updating I/O statistics reported by the <b>iostatt</b> subroutine.
<b>iostdel</b>	Removes the registration of an I/O statistics structure used for maintaining I/O statistics on a particular device.
<b>kmod_entrypt</b>	Returns a function pointer to a kernel module's entry point.
<b>kmod_load</b>	Loads an object file into the kernel or queries for an object file already loaded.
<b>kmod_unload</b>	Unloads a kernel object file.
<b>pio_assist</b>	Provides a standardized programmed I/O exception handling mechanism for all routines performing programmed I/O.
<b>prochadd</b>	Adds a system wide process state-change notification routine.
<b>prochdel</b>	Deletes a process state change notification routine.
<b>selreg</b>	Registers an asynchronous poll or select request with the kernel.
<b>selnotify</b>	Wakes up processes waiting in a <b>poll</b> or <b>select</b> subroutine or the <b>fp_poll</b> kernel service.
<b>setuerror</b>	Allows kernel extensions to set the <code>ut_error</code> field for the current thread.
<b>uexadd</b>	Adds a system wide exception handler for catching user-mode process exceptions.
<b>uexblock</b>	Makes the currently active kernel thread not runnable when called from a user-mode exception handler.
<b>uexclear</b>	Makes a kernel thread blocked by the <b>uexblock</b> service runnable again.
<b>uexdel</b>	Deletes a previously added system-wide user-mode exception handler.

---

## Locking Kernel Services

The following information is provided to assist you in understanding the locking kernel services:

- Lock Allocation and Other Services
- Simple Locks
- Complex Locks
- LockI Locks
- Atomic Operations

### Lock Allocation and Other Services

The following lock allocation services allocate and free internal operating system memory for simple and complex locks, or check if the caller owns a lock:

<b>lock_alloc</b>	Allocates system memory for a simple or complex lock.
<b>lock_free</b>	Frees the system memory of a simple or complex lock.
<b>lock_mine</b>	Checks whether a simple or complex lock is owned by the caller.

### Simple Locks

Simple locks are exclusive-write, non-recursive locks that protect thread-thread or thread-interrupt critical sections. Simple locks are preemptable, meaning that a kernel thread can be preempted by another, higher priority kernel thread while it holds a simple lock. The simple lock kernel services are:

<b>simple_lock_init</b>	Initializes a simple lock.
<b>simple_lock, simple_lock_try</b>	Locks a simple lock.
<b>simple_unlock</b>	Unlocks a simple lock.

On a multiprocessor system, simple locks that protect thread-interrupt critical sections must be used in conjunction with interrupt control in order to serialize execution both within the executing processor and between different processors. On a uniprocessor system interrupt control is sufficient; there is no need to use locks. The following kernel services provide appropriate locking calls for the system on which they are executed:

**disable\_lock**                      Raises the interrupt priority, and locks a simple lock if necessary.  
**unlock\_enable**                    Unlocks a simple lock if necessary, and restores the interrupt priority.

Using the **disable\_lock** and **unlock\_enable** kernel services to protect thread-interrupt critical sections (instead of calling the underlying interrupt control and locking kernel services directly) ensures that multiprocessor-safe code does not make unnecessary locking calls on uniprocessor systems.

Simple locks are spin locks; a kernel thread that attempts to acquire a simple lock may spin (busy-wait: repeatedly execute instructions which do nothing) if the lock is not free. The table shows the behavior of kernel threads and interrupt handlers that attempt to acquire a busy simple lock.

Caller	Owner is Running	Owner is Sleeping
Thread (with interrupts enabled)	Caller spins initially; it sleeps if the maximum spin threshold is crossed.	Caller sleeps immediately.
Interrupt handler or thread (with interrupts disabled)	Caller spins until lock is acquired.	Caller spins until lock is freed (must not happen).

**Note:** On uniprocessor systems, the maximum spin threshold is set to one, meaning that that a kernel thread will never spin waiting for a lock.

A simple lock that protects a thread-interrupt critical section must never be held across a sleep, otherwise the interrupt could spin for the duration of the sleep, as shown in the table. This means that such a routine must not call any external services that might result in a sleep. In general, using any kernel service which is callable from process level may result in a sleep, as can accessing unpinned data. These restrictions do not apply to simple locks that protect thread-thread critical sections.

The lock word of a simple lock must be located in pinned memory if simple locking services are called with interrupts disabled.

## Complex Locks

Complex locks are read-write locks that protect thread-thread critical sections. Complex locks are preemptable, meaning that a kernel thread can be preempted by another, higher priority kernel thread while it holds a complex lock. The complex lock kernel services are:

<b>lock_init</b>	Initializes a complex lock.
<b>lock_islocked</b>	Tests whether a complex lock is locked.
<b>lock_done</b>	Unlocks a complex lock.
<b>lock_read, lock_try_read</b>	Locks a complex lock in shared-read mode.
<b>lock_read_to_write, lock_try_read_to_write</b>	Upgrades a complex lock from shared-read mode to exclusive-write mode.
<b>lock_write, lock_try_write</b>	Locks a complex lock in exclusive-write mode.
<b>lock_write_to_read</b>	Downgrades a complex lock from exclusive-write mode to shared-read mode.
<b>lock_set_recursive</b>	Prepares a complex lock for recursive use.
<b>lock_clear_recursive</b>	Prevents a complex lock from being acquired recursively.

By default, complex locks are not recursive (they cannot be acquired in exclusive-write mode multiple times by a single thread). A complex lock can become recursive through the **lock\_set\_recursive** kernel service. A recursive complex lock is not freed until **lock\_done** is called once for each time that the lock was locked.

Complex locks are not spin locks; a kernel thread that attempts to acquire a complex lock may spin briefly (busy-wait: repeatedly execute instructions which do nothing) if the lock is not free. The table shows the

behavior of kernel threads that attempt to acquire a busy complex lock:

Current® Lock Mode	Owner is Running and no Other Thread is Asleep on This Lock	Owner is Sleeping
Exclusive-write	Caller spins initially, but sleeps if the maximum spin threshold is crossed, or if the owner later sleeps.	Caller sleeps immediately.
Shared-read being acquired for exclusive-write	Caller sleeps immediately.	
Shared-read being acquired for shared-read	Lock granted immediately	

**Note:**

1. On uniprocessor systems, the maximum spin threshold is set to one, meaning that a kernel thread will never spin waiting for a lock.
2. The concept of a single owner does not apply to a lock held in shared-read mode.

## Lockl Locks

**Note:** Lockl locks (previously called conventional locks) are only provided to ensure compatibility with existing code. New code should use simple or complex locks.

Lockl locks are exclusive-access and recursive locks. The lockl lock kernel services are:

**lockl**                Locks a conventional lock.  
**unlockl**            Unlocks a conventional lock.

A thread which tries to acquire a busy lockl lock sleeps immediately.

The lock word of a lockl lock must be located in pinned memory if the lockl service is called with interrupts disabled.

## Atomic Operations

Atomic operations are sequences of instructions that guarantee atomic accesses and updates of shared single word variables. This means that atomic operations cannot protect accesses to complex data structures in the way that locks can, but they provide a very efficient way of serializing access to a single word.

The atomic operation kernel services are:

**fetch\_and\_add**                                Increments a single word variable atomically.  
**fetch\_and\_and, fetch\_and\_or**                Manipulates bits in a single word variable atomically.  
**compare\_and\_swap**                            Conditionally updates or returns a single word variable atomically.

Single word variables accessed by atomic operations must be aligned on a full word boundary, and must be located in pinned memory if atomic operation kernel services are called with interrupts disabled.



---

## File Descriptor Management Services

The File Descriptor Management services are supplied by the logical file system for creating, using, and maintaining file descriptors. These services allow for the implementation of system calls that use a file descriptor as a parameter, create a file descriptor, or return file descriptors to calling applications. The following are the File Descriptor Management services:

<b>ufdcreate</b>	Allocates and initializes a file descriptor.
<b>ufdhold</b>	Increments the reference count on a file descriptor.
<b>ufdrele</b>	Decrements the reference count on a file descriptor.
<b>ufdgetf</b>	Gets a file structure pointer from a held file descriptor.
<b>getufdflags</b>	Gets the flags from a file descriptor.
<b>setufdflags</b>	Sets flags in a file descriptor.

---

## Logical File System Kernel Services

The Logical File System services (also known as the **fp\_services**) allow processes running in kernel mode to open and manipulate files in the same way that user-mode processes do. Data access limitations make it unreasonable to accomplish these tasks with system calls, so a subset of the file system calls has been provided with an alternate kernel-only interface.

The Logical File System services are one component of the logical file system, which provides the functions required to map system call requests to virtual file system requests. The logical file system is responsible for resolution of file names and file descriptors. It tracks all open files in the system using the file table. The Logical File System services are lower level entry points into the system call support within the logical file system.

Routines in the kernel that must access data stored in files or that must set up paths to devices are the primary users of these services. This occurs most commonly in device drivers, where a lower level device driver must be accessed or where the device requires microcode to be downloaded. Use of the Logical File System services is not, however, restricted to these cases.

A process can use the Logical File System services to establish access to a file or device by calling:

- The **fp\_open** service with a path name to the file or device it must access.
- The **fp\_opendev** service with the device number of a device it must access.
- The **fp\_getf** service with a file descriptor for the file or device. If the process wants to retain access past the duration of the system call, it must then call the **fp\_hold** service to acquire a private file pointer.

These three services return a file pointer that is needed to call the other Logical File System services. The other services provide the functions that are provided by the corresponding system calls.

## Other Considerations

The Logical File System services are available only in the process environment.

In addition, calling the **fp\_open** service at certain times can cause a deadlock. The lookup on the file name must acquire file system locks. If the process is already holding any lock on a component of the path, the process will be deadlocked. Therefore, do not use the **fp\_open** service when the process is already executing an operation that holds file system locks on the requested path. The operations most likely to cause this condition are those that create files.

## List of Logical File System Kernel Services

These are the Logical File System kernel services:

<b>fp_access</b>	Checks for access permission to an open file.
<b>fp_close</b>	Closes a file.
<b>fp_fstat</b>	Gets the attributes of an open file.
<b>fp_getdevno</b>	Gets the device number or channel number for a device.
<b>fp_getf</b>	Retrieves a pointer to a file structure.
<b>fp_hold</b>	Increments the open count for a specified file pointer.
<b>fp_ioctl</b>	Issues a control command to an open device or file.
<b>fp_lseek</b>	Changes the current offset in an open file.
<b>fp_llseek</b>	Changes the current offset in an open file. Used to access offsets beyond 2GB.
<b>fp_open</b>	Opens special and regular files or directories.
<b>fp_opendev</b>	Opens a device special file.
<b>fp_poll</b>	Checks the I/O status of multiple file pointers, file descriptors, and message queues.
<b>fp_read</b>	Performs a read on an open file with arguments passed.
<b>fp_readv</b>	Performs a read operation on an open file with arguments passed in <b>iovec</b> elements.
<b>fp_rwuio</b>	Performs read or write on an open file with arguments passed in a <b>uio</b> structure.
<b>fp_select</b>	Provides for cascaded, or redirected, support of the select or poll request.
<b>fp_write</b>	Performs a write operation on an open file with arguments passed.
<b>fp_writew</b>	Performs a write operation on an open file with arguments passed in <b>iovec</b> elements.
<b>fp_fsync</b>	Writes changes for a specified range of a file to permanent storage.

---

## Programmed I/O (PIO) Kernel Services

The following is a list of PIO kernel services:

<b>io_map</b>	Attaches to an I/O mapping
<b>io_map_clear</b>	Removes an I/O mapping segment
<b>io_map_init</b>	Creates and initializes an I/O mapping segment
<b>io_unmap</b>	Detaches from an I/O mapping

These kernel services are defined in the **adspace.h** and **ioacc.h** header files.

For a list of PIO macros, see Programmed I/O Services in *Understanding the Diagnostic Subsystem for AIX*.

---

## Memory Kernel Services

The Memory kernel services provide kernel extensions with the ability to:

- Dynamically allocate and free memory
- Pin and unpin code and data
- Access user memory and transfer data between user and kernel memory
- Create, reference, and change virtual memory objects

The following information is provided to assist you in learning more about memory kernel services:

- Memory Management Kernel Services
- Idata Kernel Services
- Memory Pinning Kernel Services
- User Memory Access Kernel Services
- Virtual Memory Management Kernel Services

- Cross-Memory Kernel Services

## Memory Management Kernel Services

The Memory Management services are:

<b>init_heap</b>	Initializes a new heap to be used with kernel memory management services.
<b>xmalloc</b>	Allocates memory.
<b>xmfree</b>	Frees allocated memory.

## ldata Kernel Services

**ldata** stands for "local data". The **ldata** facility supports data localization by allowing kernel subsystems and extensions to create and use **ldata** pools. The element size plus the initial and maximum number of elements to be contained in the pool is specified while creating the pool. The number of elements in the pool can be dynamically increased up to the maximum. Within the sub-pool the elements are contained such that they are cache-aligned and multiples of cache-line size to promote cache friendliness. Also, the elements in each sub-pool are backed by physical memory local to its corresponding SRAD. Allocation of a storage element is satisfied from the per-SRAD sub-pool on which the caller is located and where the storage element is to be predominately accessed. Deallocation of a storage element returns the element to its associated per-SRAD sub-pool.

Services are provided to allow kernel subsystems and extensions to create, destroy, grow the **ldata** pools. There are a couple of advantages of using **ldata** kernel services over raw **xmallocs**:

1. Since the memory allocated by **ldata** kernel services are backed by local node memory, it is faster to read and write the **ldata** region on that node.
2. **ldata** elements can be allocated from the interrupt environment. **xmalloc** kernel service cannot be called from the interrupt environment. Of course, there is an upper limit on a given **ldata** pool -- the maximum number of elements asked at **ldata** creation time.

The **ldata** services are:

<b>ldata_create</b>	Creates a SRAD-aware pinned storage element pool ( <b>ldata</b> pool) and returns its handle.
<b>ldata_destroy</b>	Destroys a <b>ldata</b> pool created by <b>ldata_create</b> .
<b>ldata_grow</b>	Expands the count of available pinned storage elements contained within a <b>ldata</b> pool.
<b>ldata_alloc</b>	Allocates a pinned storage element from a <b>ldata</b> pool.
<b>ldata_free</b>	Frees a pinned storage element to a <b>ldata</b> pool.

## Memory Pinning Kernel Services

The Memory Pinning services are:

<b>ltpin</b>	Pins the address range in the system (kernel) space and frees the page space for the associated pages.
<b>ltunpin</b>	Unpins the address range in system (kernel) address space and reallocates paging space for the specified region.
<b>pin</b>	Pins the address range in the system (kernel) space.
<b>pincode</b>	Pins the code and data associated with a loaded object module.
<b>unpin</b>	Unpins the address range in system (kernel) address space.
<b>unpincode</b>	Unpins the code and data associated with a loaded object module.
<b>xmempin</b>	Pins the specified address range in user or system memory, given a valid cross-memory descriptor.
<b>xmemunpin</b>	Unpins the specified address range in user or system memory, given a valid cross-memory descriptor.

## User-Memory-Access Kernel Services

In a system call or kernel extension running under a user process, data in the user process can be moved in or out of the kernel using the **copyin** and **copyout** services. The **uiomove** service is used for scatter and gather operations. If user data is to be referenced asynchronously, such as from an interrupt handler or a kernel process, the cross memory services must be used.

The User-Memory-Access kernel services are:

<b>copyin, copyin64</b>	Copies data between user and kernel memory.
<b>copyinstr, copyinstr64</b>	Copies a character string (including the terminating null character) from user to kernel space.
<b>copyout, copyout64</b>	Copies data between user and kernel memory.
<b>fubyte, fubyte64</b>	Fetches, or retrieves, a byte of data from user memory.
<b>fuword, fuword64</b>	Fetches, or retrieves, a word of data from user memory.
<b>subyte, subyte64</b>	Stores a byte of data in user memory.
<b>suword, suword64</b>	Stores a word of data in user memory.
<b>uiomove</b>	Moves a block of data between kernel space and a space defined by a <b>uio</b> structure.
<b>ureadc</b>	Writes a character to a buffer described by a <b>uio</b> structure.
<b>uwritec</b>	Retrieves a character from a buffer described by a <b>uio</b> structure.

**Note:** The **copyin64**, **copyout64**, **copyinstr64**, **fubyte64**, **fuword64**, **subyte64**, and **suword64** kernel services are defined as macros when compiling kernel extensions on the 64-bit kernel. The macros invoke the corresponding kernel services without the "64" suffix.

## Virtual Memory Management Kernel Services

These services are described in more detail in "Understanding Virtual Memory Manager Interfaces" on page 68. The Virtual Memory Management services are:

<b>as_att64</b>	Selects, allocates, and maps a specified region in the current user address space.
<b>as_det64</b>	Unmaps and deallocates a region in the specified address space that was mapped with the <b>as_att64</b> kernel service.
<b>as_geth64</b>	Obtains a handle to the virtual memory object for the specified address given in the specified address space. The virtual memory object is protected.
<b>as_getsrval64</b>	Obtains a handle to the virtual memory object for the specified address given in the specified address space.
<b>as_puth64</b>	Indicates that no more references will be made to a virtual memory object that was obtained using the <b>as_geth64</b> kernel service.
<b>as_seth64</b>	Maps a specified region in the specified address space for the specified virtual memory object.
<b>vm_att</b>	Maps a specified virtual memory object to a region in the current address space.
<b>vm_cflush</b>	Flushes the processor's cache for a specified address range.
<b>vm_det</b>	Unmaps and deallocates the region in the current address space that contains a given address.
<b>vm_flushp</b>	Flushes the specified range of pages.
<b>vm_galloc</b>	Allocates a region of global memory in the 64-bit kernel.
<b>vm_gfree</b>	Frees a region of global memory in the kernel previously allocated with the <b>vm_galloc</b> kernel service.
<b>vm_handle</b>	Constructs a virtual memory handle for mapping a virtual memory object with specified access level.
<b>vm_invalidatep</b>	Releases page frames in the specified range for a non-journalled persistent segment or client segment.
<b>vm_ioaccessp</b>	Initiates asynchronous page-in or page-out for the specified range of pages .
<b>vm_makep</b>	Makes a page in client storage.
<b>vm_mount</b>	Adds a file system to the paging device table.
<b>vm_mounte</b>	Adds a file system with a thread-level strategy routine to the paging device table.

<b>vm_move</b>	Moves data between a virtual memory object and a buffer specified in the <b>uio</b> structure.
<b>vm_mvc</b>	Reads or writes partial pages of files.
<b>vm_protectp</b>	Sets the page protection key for a page range.
<b>vm_qmodify</b>	Determines whether a mapped file has been changed.
<b>vm_qpages</b>	Returns the number of in-memory page frames associated with the virtual memory object.
<b>vm_readp</b>	Initiates asynchronous page-in for the range of pages specified.
<b>vm_release</b>	Releases virtual memory resources for the specified address range.
<b>vm_releasep</b>	Releases virtual memory resources for the specified page range.
<b>vm_segmap</b>	Creates the segments associated with a range of bytes in a file and attaches them to the kernel's address space.
<b>vm_setdevid</b>	Modifies the paging-device table entry for a virtual memory object.
<b>vm_thrpgio_pop</b>	Retrieves the latest per-thread context information in a client file system with a thread page-I/O strategy routine.
<b>vm_thrpgio_push</b>	Stores the current per-thread context information in a client file system with a thread page-I/O strategy routine.
<b>vm_uiomove</b>	Moves data between a virtual memory object and a buffer specified in the <b>uio</b> structure.
<b>vm_umount</b>	Removes a file system from the paging device table.
<b>vm_vmid</b>	Converts a virtual memory handle to a virtual memory object (id).
<b>vm_write</b>	Initiates page-out for a page range in the address space.
<b>vm_writep</b>	Initiates page-out for a page range in a virtual memory object.
<b>vms_create</b>	Creates a virtual memory object of the type and size and limits specified.
<b>vms_delete</b>	Deletes a virtual memory object.
<b>vms_iowait</b>	Waits for the completion of all page-out operations for pages in the virtual memory object.

## Cross-Memory Kernel Services

The cross-memory services allow data to be moved between the kernel and an address space other than the current process address space. A data area within one region of an address space is attached by calling the **xmattach** service. As a result, the virtual memory object cannot be deleted while data is being moved in or out of pages belonging to it. A cross-memory descriptor is filled out by the **xmattach** service. The attach operation must be done while under a process. When the data movement is completed, the **xmdetach** service can be called. The detach operation can be done from an interrupt handler.

The **xmemin** service can be used to transfer data from an address space to kernel space. The **xmemout** service can be used to transfer data from kernel space to an address space. These routines may be called from interrupt handler level routines if the referenced buffers are in memory.

Cross-memory services provide the **xmemdma64** service to prepare a page for DMA processing. The **xmemdma64** service can be called from the process or interrupt environments. The **xmemdma64** service returns the real address of the page for use in preparing DMA address lists. When the DMA transfer is completed, the **xmemdma64** service must be called again to unhide the page.

Data movement by DMA or an interrupt handler requires that the pages remain in memory. This is ensured by pinning the data areas using the **xmempin** service. This can only be done under a process, because the memory pinning services page-fault on pages not present in memory.

The **xmemunpin** service unpins pinned pages. This can be done by an interrupt handler if the data area is the global kernel address space. It must be done under the process if the data area is in user process space.

The Cross-Memory services are:

<b>xmattach</b>	Attaches to a user buffer for cross-memory operations.
<b>xmdetach</b>	Detaches from a user buffer used for cross-memory operations.
<b>xmemin</b>	Performs a cross-memory move by copying data from the specified address space to kernel global memory.
<b>xmemout</b>	Performs a cross-memory move by copying data from kernel global memory to a specified address space.
<b>xmemdma64</b>	Prepares a page for DMA I/O or processes a page after DMA I/O is complete. Returns 64-bit real address.
<b>xmemzero</b>	Zeroes a buffer described by a cross-memory descriptor.

---

## Understanding Virtual Memory Manager Interfaces

The virtual memory manager supports functions that allow a wide range of kernel extension data operations.

The following aspects of the virtual memory manager interface are discussed:

- Virtual Memory Objects
- Addressing Data
- Moving Data to or from a Virtual Memory Object
- Data Flushing
- Discarding Data
- Protecting Data
- Executable Data
- Installing Pager Backends
- Referenced Routines

### Virtual Memory Objects

A *virtual memory object* is an abstraction for the contiguous data that can be mapped into a region of an address space. As a data object, it is independent of any address space. The data it represents can be in memory or on an external storage device. The data represented by the virtual memory object can be shared by mapping the virtual memory object into each address space sharing the access, with the access capability of each mapping represented in that address space map.

File systems use virtual memory objects so that the files can be referenced using a mapped file access method. The mapped file access method represents the data through a virtual memory object, and allows the virtual memory manager to handle page faults on the mapped file. When a page fault occurs, the virtual memory manager calls the services supplied by the service provider (such as a virtual file system) to get and put pages. A data provider (such as a file system) maintains any data structures necessary to map between the virtual memory object offset and external storage addressing.

The data provider creates a virtual memory object when it has a request for access to the data. It deletes the virtual memory object when it has no more clients referencing the data in the virtual memory object.

The **vms\_create** service is called to create virtual memory objects. The **vms\_delete** service is called to delete virtual memory objects.

### Addressing Data

Data in a virtual memory object is made addressable in user or kernel processes through the **shmat** subroutine. A kernel extension uses the **vm\_att** kernel service to select and allocate a region in the current (per-process kernel) address space.

The per-process kernel address space initially sees only global kernel memory and the per-process kernel data. The **vm\_att** service allows kernel extensions to allocate additional regions. However, this augmented per-process kernel address space does not persist across system calls. The additional regions must be re-allocated with each entry into the kernel protection domain.

The **vm\_att** service takes as an argument a virtual memory handle representing the virtual memory object and the access capability to be used. The **vm\_handle** service constructs the virtual memory handles.

When the kernel extension has finished processing the data mapped into the current address space, it should call the **vm\_det** service to deallocate the region and remove access.

## Moving Data to or from a Virtual Memory Object

A data provider (such as a file system) can call the **vm\_makep** service to cause a memory page to be instantiated. This permits a page of data to be moved into a virtual memory object page without causing the virtual memory manager to page in the previous data contents from an external source. This is an operation on the virtual memory object, not an address space range.

The **vm\_move** and **vm\_uimove** kernel services move data between a virtual memory object and a buffer specified in a **uio** structure. This allows data providers (such as a file system) to move data to or from a specified buffer to a designated offset in a virtual memory object. This service is similar to **uimove** service, but the trusted buffer is replaced by the virtual memory object, which need not be currently addressable.

## Data Flushing

A kernel extension can initiate the writing of a data area to external storage with the **vm\_write** kernel service, if it has addressability to the data area. The **vm\_writep** kernel service can be used if the virtual memory object is not currently addressable.

If the kernel extension needs to ensure that the data is moved successfully, it can wait on the I/O completion by calling the **vms\_iowait** service, giving the virtual memory object as an argument.

## Discarding Data

The pages specified by a data range can be released from the underlying virtual memory object by calling the **vm\_release** service. The virtual memory manager deallocates any associated paging space slots. A subsequent reference to data in the range results in a page fault.

A virtual memory data provider can release a specified range of pages in a virtual memory object by calling the **vm\_releasep** service. The virtual memory object need not be addressable for this call.

## Protecting Data

The **vm\_protectp** service can change the storage protect keys in a page range in one client storage virtual memory object. This only acts on the resident pages. The pages are referred to through the virtual memory object. They do not need to be addressable in the current address space. A client file system data provider uses this protection to detect stores of in-memory data, so that mapped files can be extended by storing into them beyond their current end of file.

## Executable Data

If the data moved is to become executable, any data remaining in processor cache must be guaranteed to be moved from cache to memory. This is because the retrieval of the instruction does not need to use the data cache. The **vm\_cflush** service performs this operation.

## Installing Pager Backends

The kernel-extension data providers must provide appropriate routines to be called by the virtual memory manager. These routines move a page-sized block of data into or out of a specified page. These services are also referred to as *pager backends*.

For a local device, the device strategy routine is required. A call to the **vm\_mount** service is used to identify the device (through a **dev\_t** value) to the virtual memory manager.

For a remote data provider, the routine required is a strategy routine, which is specified in the **vm\_mount** service. These strategy routines must run as interrupt-level routines. They must not page fault, and they cannot sleep waiting for locks.

When access to a remote data provider or a local device is removed, the **vm\_umount** service must be called to remove the device entry from the virtual memory manager's paging device table.

## Thread Page-I/O Strategy Routine

Pager backends have their strategy routines run at the interrupt level of the Virtual Memory Manager (VMM). Some file systems, however, require that their strategy routines run in the context of the thread that directly causes the page I/O to occur. The VMM accommodates the difference with a blocking, or thread page-I/O strategy routine. If no thread directly causes page I/O to occur, for example, in the case of page replacement, then the strategy routine is invoked in the context of a worker thread from a pool of worker threads that are managed by the VMM.

A file system with a thread page-I/O strategy routine can use all AIX facilities as part of its handling of page I/O or page protection faults. If the file system intends to re-enter the VMM through client segment page faults or to use VMM services that involve client file segments, the file system must first save the per-thread VMM context by calling the **vm\_thrpgio\_push** service. The per-thread context can be restored using the **vm\_thrpgio\_pop** service.

To use a client file system with a thread-level strategy routine, mount the file system with the **vm\_mounte** service with its strategy routine specified and the **D\_THRPGIO** and **D\_ENHANCEDIO** flags set. To remove the file system from the paging device table, call the **vm\_umount** service.

## Referenced Routines

The virtual memory manager exports these routines exported to kernel extensions:

### Services That Manipulate Virtual Memory Objects

<b>vm_att</b>	Selects and allocates a region in the current address space for the specified virtual memory object.
<b>vms_create</b>	Creates virtual memory object of the specified type and size limits.
<b>vms_delete</b>	Deletes a virtual memory object.
<b>vm_det</b>	Unmaps and deallocates the region at a specified address in the current address space.
<b>vm_flushp</b>	Flushes the specified range of pages.
<b>vm_handle</b>	Constructs a virtual memory handle for mapping a virtual memory object with a specified access level.
<b>vm_invalidatep</b>	Releases page frames in the specified range for a non-journaled persistent segment or client segment.
<b>vm_ioaccessp</b>	Initiates asynchronous page-in or page-out for the specified range of pages .
<b>vms_iowait, vms_iowaitf</b>	Waits for the completion of all page-out operations in the virtual memory object.
<b>vm_makep</b>	Makes a page in client storage.



## Services That Manipulate Virtual Memory Objects

<b>vm_move</b>	Moves data between the virtual memory object and buffer specified in the <b>uio</b> structure.
<b>vm_mvc</b>	Reads or writes partial pages of files.
<b>vm_protectp</b>	Sets the page protection key for a page range.
<b>vm_qpages</b>	Returns the number of in-memory page frames associated with the virtual memory object.
<b>vm_readp</b>	Initiates asynchronous page-in for the range of pages specified.
<b>vm_releasep</b>	Releases page frames and paging space slots for pages in the specified range.
<b>vm_segmap</b>	Creates the segments associated with a range of bytes in a file and attaches them to the kernel's address space.
<b>vm_setdevid</b>	Modifies the paging-device table entry for a virtual memory object.
<b>vm_uiomove</b>	Moves data between the virtual memory object and buffer specified in the <b>uio</b> structure.
<b>vm_vmid</b>	Converts a virtual memory handle to a virtual memory object (id).
<b>vm_writep</b>	Initiates page-out for a page range in a virtual memory object.

The following services support address space operations:

<b>vm_cflush</b>	Flushes cache lines for a specified address range.
<b>vm_release</b>	Releases page frames and paging space slots for the specified address range.
<b>vm_write</b>	Initiates page-out for an address range.

The following Memory-Pinning kernel services also support address space operations. They are the **pin** and **unpin** services.

## Services That Support Cross-Memory Operations

Cross Memory Services are listed in "Memory Kernel Services".

## Services that Support the Installation of Pager Backends

<b>vm_mount</b>	Allocates an entry in the paging device table.
<b>vm_umount</b>	Removes a file system from the paging device table.

## Services that Support 64-bit Processes

<b>as_att64</b>	Allocates and maps a specified region in the current user address space.
<b>as_det64</b>	Unmaps and deallocates a region in the current user address space that was mapped with the <b>as_att64</b> kernel service.
<b>as_geth64</b>	Obtains a handle to the virtual memory object for the specified address.
<b>as_puth64</b>	Indicates that no more references will be made to a virtual memory object using the <b>as_geth64</b> kernel service.
<b>as_seth64</b>	Maps a specified region for the specified virtual memory object.
<b>as_getsrval64</b>	Obtains a handle to the virtual memory object for the specified address.
<b>IS64U</b>	Determines if the current user address space is 64-bit or not.

---

## Message Queue Kernel Services

The Message Queue kernel services provide the same message queue functions to a kernel extension as the **msgctl**, **msgget**, **msgsnd**, and **msgxrcv** subroutines make available to a program executing in user mode. Parameters have been added for moving returned information to an explicit parameter to free the return codes for error code usage. Instead of the error information available in the **errno** global variable (as in user mode), the Message Queue services use the service's return code. The error values are the same, except that a memory fault error (**EFAULT**) cannot occur because message buffer pointers in the kernel address space are assumed to be valid.

The Message Queue services can be called only from the process environment because they prevent the caller from specifying kernel buffers. These services can be used as an Interprocess Communication mechanism to other kernel processes or user-mode processes. See Kernel Extension and Device Driver Management Services for more information on the functions that these services provide.

There are four Message Queue services available from the kernel:

<b>kmsgctl</b>	Provides message-queue control operations.
<b>kmsgget</b>	Obtains a message-queue identifier.
<b>kmsgrcv</b>	Reads a message from a message queue.
<b>kmsgsnd</b>	Sends a message using a previously defined message queue.

---

## Network Kernel Services

The Network kernel services are divided into:

- Address Family Domain and Network Interface Device Driver services
- Routing and Interface services
- Loopback services
- Protocol services
- Communications Device Handler Interface services

### Address Family Domain and Network Interface Device Driver Kernel Services

The Address Family Domain and Network Interface Device Driver services enable address family domains (Protocols) and network interface drivers to add and remove themselves from network switch tables.

The **if\_attach** service and **if\_detach** services add and remove network interfaces from the Network Interface List. Protocols search this list to determine an appropriate interface on which to transmit a packet.

Protocols use the **add\_input\_type** and **del\_input\_type** services to notify network interface drivers that the protocol is available to handle packets of a certain type. The Network Interface Driver uses the **find\_input\_type** service to distribute packets to a protocol.

The **add\_netisr** and **del\_netisr** services add and delete network software interrupt handlers. Address families add and delete themselves from the Address Family Domain switch table by using the **add\_domain\_af** and **del\_domain\_af** services. The Address Family Domain switch table is a list of all available protocols that can be used in the **socket** subroutine.

The Address Family Domain and Network Interface Device Driver services are:

<b>add_domain_af</b>	Adds an address family to the Address Family domain switch table.
<b>add_input_type</b>	Adds a new input type to the Network Input table.

<b>add_netisr</b>	Adds a network software interrupt service to the Network Interrupt table.
<b>del_domain_af</b>	Deletes an address family from the Address Family domain switch table.
<b>del_input_type</b>	Deletes an input type from the Network Input table.
<b>del_netisr</b>	Deletes a network software interrupt service routine from the Network Interrupt table.
<b>find_input_type</b>	Finds the given packet type in the Network Input Interface switch table and distributes the input packet according to the table entry for that type.
<b>if_attach</b>	Adds a network interface to the network interface list.
<b>if_detach</b>	Deletes a network interface from the network interface list.
<b>ifunit</b>	Returns a pointer to the <b>ifnet</b> structure of the requested interface.
<b>schednetisr</b>	Schedules or invokes a network software interrupt service routine.

## Routing and Interface Address Kernel Services

The Routing and Interface Address services provide protocols with a means of establishing, accessing, and removing routes to remote hosts or gateways. Routes bind destinations to a particular network interface.

The interface address services accept a destination address or network and return an associated interface address. Protocols use these services to determine if an address is on a directly connected network.

The Routing and Interface Address services are:

<b>ifa_ifwithaddr</b>	Locates an interface based on a complete address.
<b>ifa_ifwithstaddr</b>	Locates the point-to-point interface with a given destination address.
<b>ifa_ifwithnet</b>	Locates an interface on a specific network.
<b>if_down</b>	Marks an interface as down.
<b>if_nostat</b>	Zeroes statistical elements of the interface array in preparation for an attach operation.
<b>rtalloc</b>	Allocates a route.
<b>rtfree</b>	Frees the routing table entry
<b>rtinit</b>	Sets up a routing table entry, typically for a network interface.
<b>rtredirect</b>	Forces a routing table entry with the specified destination to go through the given gateway.
<b>rtrequest</b>	Carries out a request to change the routing table.

## Loopback Kernel Services

The Loopback services enable networking code to be exercised without actually transmitting packets on a network. This is a useful tool for developing new protocols without introducing network variables. Loopback services can also be used to send packets to local addresses without using hardware loopback.

The Loopback services are:

<b>loifp</b>	Returns the address of the software loopback interface structure.
<b>looutput</b>	Sends data through a software loopback interface.

## Protocol Kernel Services

Protocol kernel services provide a means of finding a particular address family as well as a raw protocol handler. The raw protocol handler basically passes raw packets up through sockets so that a protocol can be implemented in user space.

The Protocol kernel services are:

<b>pfctlinput</b>	Starts the <b>ctlinput</b> function for each configured protocol.
-------------------	---

<b>pffindproto</b>	Returns the address of a protocol switch table entry.
<b>raw_input</b>	Builds a <b>raw_header</b> structure for a packet and sends both to the raw protocol handler.
<b>raw_usrreq</b>	Implements user requests for raw protocols.

## Communications Device Handler Interface Kernel Services

The Communications Device Handler Interface services provide a standard interface between network interface drivers and communications device handlers. The **net\_attach** and **net\_detach** services open and close the device handler. Once the device handler has been opened, the **net\_xmit** service can be used to transmit packets. Asynchronous start done notifications are recorded by the **net\_start\_done** service. The **net\_error** service handles error conditions.

The Communications Device Handler Interface services are:

<b>add_netopt</b>	This macro adds a network option structure to the list of network options.
<b>del_netopt</b>	This macro deletes a network option structure from the list of network options.
<b>net_attach</b>	Opens a communications I/O device handler.
<b>net_detach</b>	Closes a communications I/O device handler.
<b>net_error</b>	Handles errors for communication network interface drivers.
<b>net_sleep</b>	Sleeps on the specified wait channel.
<b>net_start</b>	Starts network IDs on a communications I/O device handler.
<b>net_start_done</b>	Starts the done notification handler for communications I/O device handlers.
<b>net_wakeup</b>	Wakes up all sleepers waiting on the specified wait channel.
<b>net_xmit</b>	Transmits data using a communications I/O device handler.
<b>net_xmit_trace</b>	Traces transmit packets. This kernel service was added for those network interfaces that do not use the <b>net_xmit</b> kernel service to trace transmit packets.

---

## Process and Exception Management Kernel Services

The process and exception management kernel services provided by the base kernel provide the capability to:

- Create kernel processes
- Register exception handlers
- Provide process serialization
- Generate and handle signals
- Support event waiting and notification

### Creating Kernel Processes

Kernel extensions use the **creatp** and **initp** kernel services to create and initialize a kernel process. The **setpinit** kernel service allow a kernel process to change its parent process from the one that created it to the **init** process, so that the creating process does not receive the death-of-child process signal upon kernel process termination. “Using Kernel Processes” on page 8 provides additional information concerning use of these services.

### Creating Kernel Threads

Kernel extensions use the **thread\_create** and **kthread\_start** services to create and initialize kernel-only threads. For more information about threads, see “Understanding Kernel Threads” on page 6.

The **thread\_setsched** service is used to control the scheduling parameters, priority and scheduling policy, of a thread.

## Kernel Structures Encapsulation

The **getpid** kernel service is used by a kernel extension in either the process or interrupt environment to determine the current execution environment and obtain the process ID of the current process if in the process environment. The **rusage\_incr** service provides an access to the **rusage** structure.

The thread-specific **uthread** structure is also encapsulated. The **getuerror** and **setuerror** kernel services should be used to access the `ut_error` field. The **thread\_self** kernel service should be used to get the current thread's ID.

## Registering Exception Handlers

The **setjmpx**, **clrjmpx**, and **longjmpx** kernel services allow a kernel extension to register an exception handler by:

- Saving the exception handler's context with the **setjmpx** kernel service
- Removing its saved context with the **clrjmpx** kernel service if no exception occurred
- Starting the next registered exception handler with the **longjmpx** kernel service if it was unable to handle the exception

For more information concerning use of these services, see "Handling Exceptions While in a System Call" on page 31.

## Signal Management

Signals can be posted either to a kernel process or to a kernel thread. The **pidsig** service posts a signal to a specified kernel process; the **kthread\_kill** service posts a signal to a specified kernel thread. A thread uses the **sig\_chk** service to poll for signals delivered to the kernel process or thread in the kernel mode.

For more information about signal management, see "Kernel Process Signal and Exception Handling" on page 11.

## Events Management

The event notification services provide support for two types of interprocess communications:

<b>Primitive</b>	Allows only one process thread waiting on the event.
<b>Shared</b>	Allows multiple processes threads waiting on the event.

The **et\_wait** and **et\_post** kernel services support single waiter event notification by using mutually agreed upon event control bits for the kernel thread being posted. There are a limited number of control bits available for use by kernel extensions. If the **kernel\_lock** is owned by the caller of the **et\_wait** service, it is released and acquired again upon wakeup.

The following kernel services support a shared event notification mechanism that allows for multiple threads to be waiting on the shared event.

<b>e_assert_wait</b>	<b>e_wakeup</b>
<b>e_block_thread</b>	e_wakeup_one
<b>e_clear_wait</b>	e_wakeup_w_result
<b>e_sleep_thread</b>	e_wakeup_w_sig

These services support an unlimited number of shared events (by using caller-supplied event words). The following list indicates methods to wait for an event to occur:

- Calling **e\_assert\_wait** and **e\_block\_thread** successively; the first call puts the thread on the event queue, the second blocks the thread. Between the two calls, the thread can do any job, like releasing several locks. If only one lock, or no lock at all, needs to be released, one of the two other methods should be preferred.
- Calling **e\_sleep\_thread**; this service releases a simple or a complex lock, and blocks the thread. The lock can be automatically reacquired at wakeup.

The **e\_clear\_wait** service can be used by a thread or an interrupt handler to wake up a specified thread, or by a thread that called **e\_assert\_wait** to remove itself from the event queue without blocking when calling **e\_block\_thread**. The other wakeup services are event-based. The **e\_wakeup** and **e\_wakeup\_w\_result** services wake up every thread sleeping on an event queue; whereas the **e\_wakeup\_one** service wakes up only the most favored thread. The **e\_wakeup\_w\_sig** service posts a signal to every thread sleeping on an event queue, waking up all the threads whose sleep is interruptible.

The **e\_sleep** and **e\_sleepl** kernel services are provided for code that was written for previous releases of the operating system. Threads that have called one of these services are woken up by the **e\_wakeup**, **e\_wakeup\_one**, **e\_wakeup\_w\_result**, **e\_wakeup\_w\_sig**, or **e\_clear\_wait** kernel services. If the caller of the **e\_sleep** service owns the **kernel lock**, it is released before waiting and is acquired again upon wakeup. The **e\_sleepl** service provides the same function as the **e\_sleep** service except that a caller-specified lock is released and acquired again instead of the **kernel\_lock**.

## List of Process, Thread, and Exception Management Kernel Services

The Process, Thread, and Exception Management kernel services are listed below.

<b>clrjmpx</b>	Removes a saved context by popping the most recently saved jump buffer from the list of saved contexts.
<b>creatp</b>	Creates a new kernel process.
<b>e_assert_wait</b>	Asserts that the calling kernel thread is going to sleep.
<b>e_block_thread</b>	Blocks the calling kernel thread.
<b>e_clear_wait</b>	Clears the wait condition for a kernel thread.
<b>e_sleep, e_sleep_thread, or e_sleepl</b>	Forces the calling kernel thread to wait for the occurrence of a shared event.
<b>e_sleep_thread</b>	Forces the calling kernel thread to wait the occurrence of a shared event.
<b>e_wakeup, e_wakeup_one, or e_wakeup_w_result</b>	Notifies kernel threads waiting on a shared event of the event's occurrence.
<b>e_wakeup_w_sig</b>	Posts a signal to sleeping kernel threads.
<b>et_post</b>	Notifies a kernel thread of the occurrence of one or more events.
<b>et_wait</b>	Forces the calling kernel thread to wait for the occurrence of an event.
<b>getpid</b>	Gets the process ID of the current process.
<b>getppidx</b>	Gets the parent process ID of the specified process.
<b>initp</b>	Changes the state of a kernel process from idle to ready.
<b>kthread_kill</b>	Posts a signal to a specified kernel-only thread.
<b>kthread_start</b>	Starts a previously created kernel-only thread.
<b>limit_sigs</b>	Changes the signal mask for the calling kernel thread.
<b>longjmpx</b>	Allows exception handling by causing execution to resume at the most recently saved context.
<b>NLuprintf</b>	Submits a request to print an internationalized message to the controlling terminal of a process.
<b>pgsignal</b>	Sends a signal to all of the processes in a process group.
<b>pidsig</b>	Sends a signal to a process.
<b>rusage_incr</b>	Increments a field of the <b>rusage</b> structure.
<b>setjmpx</b>	Allows saving the current execution state or context.

<b>setpinit</b>	Sets the parent of the current kernel process to the <b>init</b> process.
<b>sig_chk</b>	Provides the calling kernel thread with the ability to poll for receipt of signals.
<b>sigsetmask</b>	Changes the signal mask for the calling kernel thread.
<b>sleep</b>	Forces the calling kernel thread to wait on a specified channel.
<b>thread_create</b>	Creates a new kernel-only thread in the calling process.
<b>thread_self</b>	Returns the caller's kernel thread ID.
<b>thread_setsched</b>	Sets kernel thread scheduling parameters.
<b>thread_terminate</b>	Terminates the calling kernel thread.
<b>ue_proc_check</b>	Determines if a process is critical to the system.
<b>uprintf</b>	Submits a request to print a message to the controlling terminal of a process.

---

## RAS Kernel Services

The Reliability, Availability, and Serviceability (RAS) kernel services are used to record the occurrence of hardware or software failures and to capture data about these failures. The recorded information can be examined using the **errpt** or **trcrpt** commands.

The **panic** kernel service is called when a catastrophic failure occurs and the system can no longer operate. The **panic** service performs a system dump. The system dump captures data areas that are registered in the Master Dump Table. The kernel and kernel extensions use the **dmp\_ctl** kernel service to add and delete entries in the Master Dump Table, and record dump routine failures.

The **errsave** and **errlast** kernel service is called to record an entry in the system error log when a hardware or software failure is detected.

The **trcgenk** and **trcgenkt** kernel services are used along with the **trchook** subroutine to record selected system events in the event-tracing facility.

The **ras\_register** and **ras\_unregister** kernel services register and unregister RAS handlers for a specific component. These handlers are called by the kernel when the system needs to communicate various RAS commands to each component.

The **register\_HA\_handler** and **unregister\_HA\_handler** kernel services are used to register high availability event handlers for kernel extensions that need to be aware of events such as processor deallocation.

One of the RAS features is a service that monitors for excessive periods of interrupt disablement on a processor, and logs these events to the error log. The **disablement\_checking\_suspend** and **disablement\_checking\_resume** services exempt a code segment from this detection.

---

## Security Kernel Services

The Security kernel services provide methods for controlling the auditing system and for determining the access rights to objects for the invoking process.

The following services are security kernel services:

<b>suser</b>	Determines the privilege state of a process.
<b>audit_svcstart</b>	Initiates an audit record for a system call.
<b>audit_svcbcopy</b>	Appends event information to the current audit event buffer.
<b>audit_svcfinis</b>	Writes an audit record for a kernel service.

<b>crcopy</b>	Creates a copy of a security credentials structure.
<b>crdup</b>	Creates a copy of the current security credentials structure.
credential macros	Provide a means for accessing the user and group identifier fields within a credentials structure.
<b>crexport</b>	Copies an internal format credentials structure to an external format credentials structure.
<b>crfree</b>	Frees a security credentials structure.
<b>crget</b>	Allocates a new, uninitialized security credentials structure.
<b>crhold</b>	Increments the reference count of a security credentials structure.
<b>crref</b>	Increments the reference count of the current security credentials structure.
<b>crset</b>	Replaces the current security credentials structure.
<b>kcred_genpagvalue</b>	Generates a system-wide unique PAG value for a given PAG name (such as afs).
<b>kcred_getcap</b>	Copies a capability vector from a credentials structure.
<b>kcred_getgroups</b>	Copies the concurrent group set from a credentials structure.
<b>kcred_getpag</b>	Copies a process authentication group (PAG) ID from a credentials structure.
<b>kcred_getpag64</b>	Retrieves 64-bit PAG values from a process's credentials structure.
<b>kcred_getpagid</b>	Returns the process authentication group (PAG) identifier for a PAG name.
<b>kcred_getpagname</b>	Retrieves the name of a process authentication group (PAG).
<b>kcred_getpriv</b>	Copies a privilege vector from a credentials structure.
<b>kcred_setcap</b>	Copies a capabilities set into a credentials structure.
<b>kcred_setgroups</b>	Copies a concurrent group set into a credentials structure.
<b>kcred_setpag</b>	Copies a process authentication group ID into a credentials structure.
<b>kcred_setpag64</b>	Stores 64-bit PAG values in a process's credentials structure.
<b>kcred_setpagname</b>	Copies a process authentication group ID into a credentials structure.
<b>kcred_setpriv</b>	Copies a privilege vector into a credentials structure.

---

## Timer and Time-of-Day Kernel Services

The Timer and Time-of-Day kernel services provide kernel extensions with the ability to be notified when a period of time has passed. The **tstart** service supports a very fine granularity of time. The **timeout** service is built on the **tstart** service and is provided for compatibility with earlier versions of the operating system. The **w\_start** service provides a timer with less granularity, but much cheaper path-length overhead when starting a timer.

The Timer and Time-of-Day kernel services are divided into the following categories:

- Time-of-Day services
- Fine Granularity Timer services
- Timer services for compatibility
- Watchdog Timer services

### Time-Of-Day Kernel Services

The Time-Of-Day kernel services are:

<b>curtime</b>	Reads the current time into a time structure.
<b>kgettickd</b>	Retrieves the current status of the systemwide time-of-day timer-adjustment values.
<b>ksettimer</b>	Sets the systemwide time-of-day timer.
<b>ksettickd</b>	Sets the current status of the systemwide timer-adjustment values.



## Fine Granularity Timer Kernel Services

The Fine Granularity Timer kernel services are:

<b>delay</b>	Suspends the calling process for the specified number of timer ticks.
<b>talloc</b>	Allocates a timer request block before starting a timer request.
<b>tfree</b>	Deallocates a timer request block.
<b>tstart</b>	Submits a timer request.
<b>tstop</b>	Cancels a pending timer request.

For more information about using the Fine Granularity Timer services, see “Using Fine Granularity Timer Services and Structures.”

## Timer Kernel Services for Compatibility

The following Timer kernel services are provided for compatibility:

<b>timeout</b>	Schedules a function to be called after a specified interval.
<b>timeoutcf</b>	Allocates or deallocates callout table entries for use with the <b>timeout</b> kernel service.
<b>untimeout</b>	Cancels a pending timer request.

## Watchdog Timer Kernel Services

The Watchdog timer kernel services are:

<b>w_clear</b>	Removes a watchdog timer from the list of watchdog timers known to the kernel.
<b>w_init</b>	Registers a watchdog timer with the kernel.
<b>w_start</b>	Starts a watchdog timer.
<b>w_stop</b>	Stops a watchdog timer.

---

## Using Fine Granularity Timer Services and Structures

The **talloc**, **tstart**, **tstop**, and **tfree** services provide fine-resolution timing functions. These timer services should be used when the following conditions are required:

- Timing requests for less than one second
- Critical timing
- Absolute timing

A kernel extension uses the **talloc** service to allocate a timer request block (struct **trb**). The kernel extension initializes the fields in the **trb** structure to specify the timer completion function, when the function is to receive control, and how the function is to receive control. The **tstart** service is used to schedule the timer event. The timer completion handler is called when the specified time is reached. The address of the timer request block is passed to the completion handler function. The timer completion handler can initialize and reschedule the timer request block using the **tstart** service. The **tstop** service is used to cancel a pending timer request. The **tfree** service is used to return a previously allocated timer request block to the system. Be sure that the **trb** structure is not active when it is freed by the **tstop** service.

Use the **talloc** service to allocate a timer request block. Do not allocate your own memory to contain timer request blocks.

Do not access or modify any fields of the **trb** structure between the time when a timer request block is started by the **tstart** service and the time when the completion handler runs, or a call to the **tstop** service is completed against the timer request block.

If you use the **tstart** service in the completion handler to reschedule the timer request block that was passed to the completion handler, be sure to initialize the `t->timeout` field and the appropriate **T\_INCTERVAL**, **T\_ABSOLUTE**, and **T\_LOWRES** bits in the `t->flags` field. Use an OR operation to set the `t->flags` bits because the AIX kernel timer service is using other bits in the `t->flags` field when the completion handler is called.

The AIX kernel timer service accesses the timer request block that was passed to the completion handler after the completion handler returned. A race condition might occur if a completion handler makes the timer request block that was passed to the completion handler available for the reuse by another processor. A system outage might occur if the other processor uses the **tstart** service to reschedule the timer request block before the AIX kernel timer service has completed its access to the block. The other processor can avoid the system outage using the **tstop** service to ensure that the timer request block is stopped before using the **tstart** service to reschedule the timer request block.

#### Notes:

1. If a timer request block is rescheduled anywhere other than its completion handler, use the **tstop** service to ensure that the timer request block is not scheduled or running in its completion handler. Failure to do so might result in a system outage.
2. Use the **tstop** service against the timer request block before calling the **tstart** service if a timer request block is restarted using the **tstart** service anywhere other than its completion handler. Failure to do so can result in system stops.
3. Always use the **tstop** service to cancel a timer request before calling the **tfree** service. This is necessary to ensure that the timer request block is not on a system timer queue, running the completion handler, or returning to the system from a completion handler when the timer request block is freed. Failure to do so might result in a system outage.

Do not use the **tfree** service in a timer completion handler to free the timer request block that was passed to the completion handler.

The Watchdog timer services can be used for noncritical times having a one-second resolution. The **timeout** service can be used for noncritical times having a clock-tick resolution.

## Timer Services Data Structures

The **trb** (timer request) structure is found in the `/sys/timer.h` file. The **itimerstruc\_t** structure contains the second/nanosecond structure for time operations and is found in the `sys/time.h` file.

The **trb** structure contains a number of fields. Some of the fields are used for AIX kernel processing. Others are intended for timer service users. A timer service user must set the following fields in the timer request block after calling the **talloc** service and before calling the **tstart** service.

Fields	Description
<code>t-&gt;flags</code>	Specifies the flags. You must set either the <b>T_INCTERVAL</b> flag or the <b>T_ABSOLUTE</b> flag in the field. You can also set the <b>T_LOWRES</b> flag if you want. The <b>T_LOWRES</b> flag causes the system to round the <code>t-&gt;timeout</code> value to the next timer timeout. The advantage of using the <b>T_LOWRES</b> flag is that it prevents an extra timer interrupt from being generated. The timeout is rounded to a larger value. The <b>T_LOWRES</b> flag never causes more than 10 ms to be added to the timeout value because the system maintains timers at an interval of 10 ms.
<code>t-&gt;timeout</code>	Specifies the timeout value. This field is used for both the interval timer and the absolute timer as indicated by the <code>t-&gt;flags</code> field.
<code>t-&gt;func</code>	Specifies the completion-handler function pointer. This function must remain in pinned memory as long as the function might be called. Use the <b>tstop</b> service to cancel any pending timer request blocks that can call the completion handler before unloading the kernel extension containing the function.

Fields	Description
t->ipri	Specifies the interrupt priority at which the completion handler is called. Completion handlers are called in the interrupt environment. The interrupt priority values are defined in the <b>/usr/include/sys/intr.h</b> file. The least favored interrupt priority for a timer completion handler is INTTIMER. The most favored interrupt priority is INTMAX. For more information about interrupt priorities, see “Understanding Interrupts” on page 49 in <i>AIX Version 6.1 Kernel Extensions and Device Support Programming Concepts</i> .
t->t_union	(optional) Passes data to the completion handler.
t->id	(optional) Specifies the ID of the process that schedules the timer request block. Device drivers can set the field to a value of -1.

The t->func, t->t\_union, t->ipri, and t->id fields persist across a call from **tstart** to the completion handler. Only the t->timeout and t->flags fields must be reset before a subsequent call to the **tstart** service.

Other fields in the **trb** structure do not need to be initialized between the **talloc** and **tstart** services.

The **itimerstruc\_t t.it** value substructure is used to store time information for both absolute and incremental timers. The **T\_ABSOLUTE** absolute request flag is defined in the **sys/timer.h** file.

## Coding the Timer Function

The t->func timer function should be declared as follows:

```
void func (t)
struct trb *t;
```

The argument to the **func** completion handler routine is the address of the **trb** structure, not the contents of the t\_union field.

The t->func timer function is called on an interrupt level. Therefore, the code for this function must be in pinned storage and must follow conventions for interrupt handlers.

---

## Using Multiprocessor-Safe Timer Services

On a multiprocessor system, timer request blocks and watchdog timer structures could be accessed simultaneously by several processors. The kernel services shown below potentially alter critical information in these blocks and structures, and therefore check whether it is safe to perform the requested service before proceeding:

<b>tstop</b>	Cancels a pending timer request.
<b>w_clear</b>	Removes a watchdog timer from the list of watchdog timers known to the kernel.
<b>w_init</b>	Registers a watchdog timer with the kernel.

If the requested service cannot be performed, the kernel service returns an error value.

In order to be multiprocessor safe, the caller must check the value returned by these kernel services. If the service was not successful, the caller must take an appropriate action, for example, retrying in a loop. If the caller holds a device driver lock, it should release and then reacquire the lock within this loop in order to avoid deadlock.

Drivers which were written for uniprocessor systems do not check the return values of these kernel services and are not multiprocessor-safe. Such drivers can still run as funnelled device drivers.

---

## Virtual File System (VFS) Kernel Services

The Virtual File System (VFS) kernel services are provided as fundamental building blocks for use when writing a virtual file system. These services present a standard interface for such functions as configuring file systems, creating and freeing v-nodes, and looking up path names.

Most functions involved in the writing of a file system are specific to that file system type. But a limited number of functions must be performed in a consistent manner across the various file system types to enable the logical file system to operate independently of the file system type.

The VFS kernel services are:

<b>common_reclock</b>	Implements a generic interface to the record locking functions.
<b>fidtovp</b>	Maps a file system structure to a file ID.
<b>gfsadd</b>	Adds a file system type to the <b>gfs</b> table.
<b>gfsdel</b>	Removes a file system type from the <b>gfs</b> table.
<b>vfs_hold</b>	Holds a <b>vfs</b> structure and increments the structure's use count.
<b>vfs_unhold</b>	Releases a <b>vfs</b> structure and decrements the structure's use count.
<b>vfsrele</b>	Releases all resources associated with a virtual file system.
<b>vfs_search</b>	Searches the vfs list.
<b>vn_free</b>	Frees a v-node previously allocated by the <b>vn_get</b> kernel service.
<b>vn_get</b>	Allocates a virtual node and associates it with the designated virtual file system.
<b>lookupvp</b>	Retrieves the v-node that corresponds to the named path.

---

## Related Information

Chapter 1, "Kernel Environment," on page 1

"Block I/O Buffer Cache Kernel Services: Overview" on page 48

Understanding the Virtual File System Interface

Communications Physical Device Handler Model Overview

Understanding File Descriptors in *AIX Version 6.1 General Programming Concepts: Writing and Debugging Programs*.

## Subroutine References

The **msgctl** subroutine, **msgget** subroutine, **msgsnd** subroutine, **msgxrcv** subroutine in *AIX Version 6.1 Technical Reference: Base Operating System and Extensions Volume 1*.

The **trchook** subroutine in *AIX Version 6.1 Technical Reference: Base Operating System and Extensions Volume 2*.

## Commands References

The **iostat** command in *AIX Version 6.1 Commands Reference, Volume 3*.

The **vmstat** command in *AIX Version 6.1 Commands Reference, Volume 6*.

## Technical References

The **talloc** kernel service, **tfree** kernel service, **tstart** kernel service, **tstop** kernel service in *AIX Version 6.1 Technical Reference: Kernel and Subsystems Volume 1*.

---

## Chapter 5. Asynchronous I/O Subsystem

*Synchronous input/output (I/O)* occurs while you wait. Applications processing cannot continue until the I/O operation is complete.

In contrast, *asynchronous I/O (AIO)* operations run in the background and do not block user applications. This improves performance, because I/O operations and application processing can run simultaneously.

Many applications, such as databases and file servers, take advantage of the ability to overlap processing and I/O. These AIO operations use various kinds of devices and files. Additionally, multiple AIO operations can run at the same time on one or more devices or files. Using AIO usually improves I/O throughput for these types of applications. The actual performance, however, depends partly on the number of concurrent I/O requests that the application can issue at one time. When the AIO fast path is not used, the performance also depends on how many AIO server processes that handle the I/O requests are running. For more information about the fast path, see “How Many AIO Servers Am I Currently Using?” on page 84.

Each AIO request has a corresponding control block in the application’s address space. When an AIO request is made, a handle is established in the control block. This handle is used to retrieve the status and the return values of the request.

Applications use the **aio\_read** and **aio\_write** subroutines to perform the I/O. Control returns to the application from the subroutine, as soon as the request has been queued. The application can then continue processing while the disk operation is being performed.

A kernel process (kproc), called an AIO server (AIOS), is in charge of each request from the time it is taken off the queue until it completes. The number of servers limits the number of disk I/O operations that can be in progress in the system simultaneously.

The default value of the *minservers* tunable is 3, and that of the *maxservers* tunable is 30. In systems that seldom run applications that use AIO, this is usually adequate. For environments with many disk drives and key applications that use AIO, the defaults might be too low. The result of a deficiency of servers is that disk I/O seems much slower than it should be. Not only do requests spend inordinate lengths of time in the queue, but the low ratio of servers to disk drives means that the seek-optimization algorithms have too few requests to work with for each drive.

**Note:** AIO does not work if the control block or buffer is created using mmap (mapping segments).

There are two AIO subsystems. The original AIX AIO, now called LEGACY AIO, has the same function names as the Portable Operating System Interface (POSIX) compliant POSIX AIO. The major differences between the two involve different parameter passing. Both subsystems are defined in the **/usr/include/sys/aio.h** file. The **\_AIO\_AIX\_SOURCE** macro is used to distinguish between the two versions.

**Note:** The **\_AIO\_AIX\_SOURCE** macro used in the **/usr/include/sys/aio.h** file must be defined when using this file to compile an AIO application with the LEGACY AIO function definitions. The default compile using the **aio.h** file is for an application with the new POSIX AIO definitions. To use the LEGACY AIO function definitions do the following in the source file:

```
#define _AIO_AIX_SOURCE
#include <sys/aio.h>
```

or when compiling on the command line, type the following:

```
xlc ... -D_AIO_AIX_SOURCE ... classic_aio_program.c
```

For each AIO function there is a legacy and a POSIX definition. LEGACY AIO has an additional  `aio_nwait`  function, which although not a part of POSIX definitions, has been included in POSIX AIO to help those who want to port from LEGACY to POSIX definitions. POSIX AIO has an additional  `aio_fsync`  function, which is not included in LEGACY AIO. For a list of these functions, see “Asynchronous I/O Subroutines” on page 87.

---

## How Do I Know if I Need to Use AIO?

Using the  `vmstat`  command with an interval and count value, you can determine if the processor is idle waiting for disk I/O. The  `wa`  column details the percentage of time the processor was idle with pending local disk I/O.

If there is at least one outstanding I/O to a local disk when the wait process is running, the time is classified as waiting for I/O. Unless AIO is being used by the process, an I/O request to disk causes the calling process to block (or sleep) until the request is complete. After a process’s I/O request completes, it is placed on the run queue.

A  `wa`  value consistently over 25 percent might indicate that the disk subsystem is not balanced properly, or it might be the result of a disk-intensive workload.

**Note:** AIO does not relieve an overly busy disk drive. Using the  `iostat`  command with an interval and count value, you can determine if any disks are overly busy. Monitor the  `%tm_act`  column for each disk drive on the system. On some systems, a  `%tm_act`  of 35.0 or higher for one disk can cause noticeably slower performance. The relief for this case could be to move data from more busy to less busy disks, but simply having AIO does not relieve an overly busy disk problem.

## SMP Systems

For SMP systems, the  `us` ,  `sy` ,  `id`  and  `wa`  columns are only averages over all processors. But remember that the I/O wait statistic per processor is not really a processor-specific statistic; it is a global statistic. An I/O wait is distinguished from idle time only by the state of a pending I/O. If there is any pending disk I/O, and the processor is not busy, then it is an I/O wait time. Disk I/O is not tracked by processors, so when there is any I/O wait, all processors get charged (assuming they are all equally idle).

## How Many AIO Servers Am I Currently Using?

To determine how many POSIX AIO Servers are currently running, type the following information on the command line:

```
pstat -a | grep posix_aio_server | wc -l
```

**Requirement:** You must run this command as the root user.

To determine how many LEGACY AIO Servers are currently running, type the following information on the command line:

```
pstat -a | egrep ' aio_server' | wc -l
```

**Requirement:** You must run this command as the root user.

A timeout value, specified by the  `server_inactivity`  tunable, causes a server to exit if the server is idle without serving an AIO request. If, when the server exits, the total number of servers falls below the value of the  `minservers`  tunable multiplied by the number of processors, the server becomes idle, waiting for an AIO request to serve. By reducing the number of idle processes that are not being used to serve AIO requests, overall system performance is enhanced.

If the data being accessed asynchronously is located in a Journaled File System (JFS), all I/O is routed through the AIOS  `kprocs` .

If the data being accessed asynchronously is located on a raw logical volume or the Enhanced Journaled File System (JFS2) in conjunction with concurrent I/O (CIO), the I/O is routed using a fast path and does not go through the AIOS kprocs. In that case the number of servers that are running is not relevant.

However, if you want to confirm that an application that uses raw logic volumes or JFS2 with CIO is taking advantage of AIO, you can disable the fast path option using the **ioo** command. When this option is disabled, even raw logical volume and JFS2 I/O are forced through the AIOS kprocs. At that point, the **pstat** command listed in the preceding discussion works. Do not run the system with this option disabled for any length of time.

You can use the AIO fast path for files that are opened with CIO on the JFS2 file system.

## How Many AIO Servers Do I Need?

You can set the value for the maximum number of servers in the following ways:

- Limit the maximum number of servers to a number equal to ten times the number of disks that are to be used concurrently.
- Set the maximum number of servers to a high number (for example, 100), depending on the system and the number of processors.
  - Change the *server\_inactivity* tunable value to 600 (the default value is 300).
  - Monitor the number of servers every 10 minutes (600 seconds) throughout the course of a normal workload. If the number of servers is constantly at the aggregate value for the *maxservers* tunable (the *maxservers* tunable value times the number of processors) and there is unused processor and IO bandwidth, consider increasing the *maxservers* tunable value. If you find that at times the system is not performing as expected, consider lowering the *maxservers* tunable value.

The goal is to find the performance balance between processor usage and IO bandwidth.

- Take statistics using **vmstat -s** before any high I/O activity begins, and again at the end. Check the field *iodone*. From this you can determine how many physical I/Os are being handled in a given wall clock period. Then increase the maximum number of servers and see if you can get more activity or event completions (*iodones*) in the same time period.

In general, consider changing the *minservers* tunable value only when an application will be issuing a high number of I/Os that is beyond the *server\_inactivity* timeout value (thus causing a high number of AIOS kprocs to be created in bursts). To help smooth this condition, either increase the *minservers* value to keep the required number of AIOS kprocs active, or increase the *server\_inactivity* tunable value so that the system will naturally keep them active if the application keeps issuing requests within the timeout window.

### Requirement:

- You must set the *minservers* tunable value at a level so that optimal performance can be obtained across an average workload. You do not need to restart the system to effect a change to the *minservers* or *maxservers* tunable.
- The value of the *minservers* tunable cannot exceed that of the *maxservers* tunable.

## Tunable Values for Asynchronous I/O

The *aio\_minreqs*, *aio\_minservers*, *aio\_maxservers*, and *aio\_server\_inactivity* nonrestricted tunables are for the LEGACY AIO subsystem; and the *posix\_aio\_minreqs*, *posix\_aio\_minservers*, *posix\_aio\_maxservers*, *posix\_aio\_server\_inactivity* nonrestricted tunables are for the POSIX AIO subsystem. For more information about each tunable, see the **ioo** command.

---

## Functions of Asynchronous I/O

Functions provided by the AIO facilities are:

- Large File-Enabled Asynchronous I/O
- Nonblocking I/O
- Notification of I/O completion
- Cancellation of I/O requests

## Large File-Enabled Asynchronous I/O

The fundamental data structure associated with all AIO operations is **struct aiocb**. Within this structure is the `aio_offset` field, which is used to specify the offset for an I/O operation.

Due to the signed 32-bit definition of `aio_offset`, the default AIO interfaces are limited to an offset of 2Gb minus 1. To overcome this limitation, a new AIO control block with a signed 64-bit offset field and a new set of AIO interfaces has been defined. These 64-bit definitions end with "64."

The large offset-enabled AIO interfaces are available under the `_LARGE_FILES` compilation environment and under the `_LARGE_FILE_API` programming environment. For further information, see *Writing Programs That Access Large Files in AIX Version 6.1 General Programming Concepts: Writing and Debugging Programs*.

Under the `_LARGE_FILES` compilation environment, AIO applications written to the default interfaces interpret the following redefinitions:

Item	Redefined To Be	Header File
<code>struct aiocb</code>	<code>struct aiocb64</code>	<code>sys/aio.h</code>
<code>aio_read()</code>	<code>aio_read64()</code>	<code>sys/aio.h</code>
<code>aio_write()</code>	<code>aio_write64()</code>	<code>sys/aio.h</code>
<code>aio_cancel()</code>	<code>aio_cancel64()</code>	<code>sys/aio.h</code>
<code>aio_suspend()</code>	<code>aio_suspend64()</code>	<code>sys/aio.h</code>
<code>aio_listio()</code>	<code>aio_listio64()</code>	<code>sys/aio.h</code>
<code>aio_return()</code>	<code>aio_return64()</code>	<code>sys/aio.h</code>
<code>aio_error()</code>	<code>aio_error64()</code>	<code>sys/aio.h</code>

For information on using the `_LARGE_FILES` environment, see *Porting Applications to the Large File Environment in AIX Version 6.1 General Programming Concepts: Writing and Debugging Programs*.

In the `_LARGE_FILE_API` environment, the 64-bit application programming interfaces (APIs) are visible. This environment requires recoding of applications to the new 64-bit API name. For further information on using the `_LARGE_FILE_API` environment, see *Using the 64-Bit File System Subroutines in AIX Version 6.1 General Programming Concepts: Writing and Debugging Programs*.

## Nonblocking I/O

After issuing an I/O request, the user application can proceed without being blocked while the I/O operation is in progress. The I/O operation occurs while the application is running. Specifically, when the application issues an I/O request, the request is queued. The application can then resume running before the I/O operation starts.

To manage AIO, each AIO request has a corresponding control block in the application's address space. This control block contains the control and status information for the request. It can be used again when the I/O operation finishes.



## Notification of I/O Completion

After issuing an AIO request, the user application can determine when and how the I/O operation finishes. This information is provided in three ways:

- The application can poll the status of the I/O operation.
- The system can asynchronously notify the application when the I/O operation is done.
- The application can block until the I/O operation is complete.

## Polling the Status of the I/O Operation

The application can periodically poll the status of the I/O operation. The status of each I/O operation is provided in the application's address space in the control block associated with each request. Portable applications can retrieve the status by using the `aio_error` subroutine. The `aio_suspend` subroutine suspends the calling process until one or more AIO requests finish.

## Asynchronously Notifying the Application When the I/O Operation Completes

Asynchronously notifying the I/O completion is done by signals. Specifically, an application can request that a `SIGIO` signal be delivered when the I/O operation is complete. To do this, the application sets a flag in the control block at the time it issues the I/O request. If several requests are issued, the application can poll the status of the requests to determine which ones completed.

## Blocking the Application until the I/O Operation Is Complete

The third way to determine whether an I/O operation is complete is to let the calling process become blocked and wait until at least one of the I/O requests it is waiting for is complete. This method is similar to synchronous style I/O. It is useful for applications that, after performing some processing, need to wait for I/O completion before proceeding.

## Cancellation of I/O Requests

Some I/O requests can be canceled. Cancellation is not guaranteed and might succeed or not depending upon the state of the individual request. If a request is in the queue and the I/O operations have not yet started, the request can be canceled. Typically, a request can no longer be canceled after the actual I/O operation begins.

---

## Asynchronous I/O Subroutines

The following 64-bit subroutines are provided for performing AIO:

Subroutine	Purpose
<code>aio_cancel</code> or <code>aio_cancel64</code>	Cancels one or more outstanding AIO requests.
<code>aio_error</code> or <code>aio_error64</code>	Retrieves the error status of an AIO request.
<code>aio_fsync</code>	Synchronizes asynchronous files.
<code>lio_listio</code> or <code>lio_listio64</code>	Initiates a list of AIO requests with a single call.
<code>aio_nwait</code>	Suspends the calling process until <i>n</i> AIO requests are completed.
<code>aio_read</code> or <code>aio_read64</code>	Reads asynchronously from a file.
<code>aio_return</code> or <code>aio_return64</code>	Retrieves the return status of an AIO request.
<code>aio_suspend</code> or <code>aio_suspend64</code>	Suspends the calling process until one or more AIO requests finishes.
<code>aio_write</code> or <code>aio_write64</code>	Writes asynchronously to a file.

## Order and Priority of Asynchronous I/O Calls

An application can issue several AIO requests on the same file or device. However, because the I/O operations are performed asynchronously, the order in which they are handled might not be the order in which the I/O calls are made. The application must enforce ordering of its own I/O requests if ordering is required.

**Note:** Priority among the I/O requests is available only for POSIX AIO.

For files that support **seek** operations, seeking can be done as part of the asynchronous read or write operations. The whence and offset fields are provided in the control block of the request to set the *seek* parameters. The *seek* pointer is updated when the asynchronous read or write call returns.

---

## Subroutines Affected by Asynchronous I/O

The following existing subroutines are affected by AIO:

- The **close** subroutine
- The **exit** subroutine
- The **exec** subroutine
- The **fork** subroutine

If the application closes a file, or calls the **\_exit** or **exec** subroutines while it has some outstanding I/O requests, the requests are canceled. If they cannot be canceled, the application is blocked until the requests finish. When a process calls the **fork** subroutine, its AIO is not inherited by the child process.

One fundamental limitation in AIO is page hiding. When an unbuffered (raw) AIO is issued, the page that contains the user buffer is hidden during the actual I/O operation. This ensures cache consistency. However, the application can access the memory locations that fall within the same page as the user buffer. This can cause the application to block as a result of a page fault. To alleviate this, allocate page aligned buffers and do not touch the buffers until the I/O request using it finishes.

---

## 64-bit Enhancements

AIO has been enhanced to support 64-bit enabled applications. On 64-bit platforms, both 32-bit and 64-bit AIO can occur simultaneously.

The **aio\_cb** structure, the fundamental data structure associated with all AIO operations, has changed. The element of this struct, **aio\_return**, is now defined as `ssize_t`. Previously, it was defined as an `int`. AIO supports large files by default. An application compiled in 64-bit mode can do AIO to a large file without any additional `#define` specifications or special opening of those files.

---

## LEGACY AIO Extended Functionality

The extended functionality supported by LEGACY AIO includes extended asynchronous I/O control block (AIOCB), I/O priorities and cache hints, and I/O completion ports.

### Extended AIOCB

LEGACY AIO supports functionality that is not available in POSIX AIO. To extend the LEGACY AIOCB, the **aio\_reqprio** and **aio\_fp** fields are deprecated, and the following new fields are introduced:

Field	Version
<b>aio_version</b>	All versions.
<b>aio_priority</b>	AIOCBX_VERS1
<b>aio_cache_hint</b>	AIOCBX_VERS1
<b>aio_iocpfd</b>	AIOCBX_VERS2

A new flag, **AIO\_EXTENDED**, has also been added to the **aio\_flags** field. If the **AIO\_EXTENDED** flag is not set, LEGACY AIO completely ignores any new extended fields. If the **AIO\_EXTENDED** flag is set within **aio\_flags**, and the **aio\_version** field contains a value greater than 0 and less than or equal to **AIOCBX\_VERSION**, all extended fields with a version indicated in the preceding table that are less than or equal to the version number specified in the AIOCB are in force. Future extensions to the legacy AIOCB structure will use new version values and introduce new extended fields beyond what is currently defined within the AIOCB structure.

Except for the **aio\_version** field, all extended fields are required to ignore a value of 0 (zero). A user of any extended field must ensure that all other unused extended fields are initialized to zero. Use either the **bzero** or **memset** function on the entire AIOCB structure prior to setting any field in the structure.

## I/O Priorities and Cache Hints

To use I/O priorities and cache hints with AIO, set the **AIO\_EXTENDED** flag in the **aio\_flags** field and the **aio\_version** field to a value of **AIOCBX\_VERS1** or greater. All other extended fields that are defined must be set to 0 if they are not used. The following fields are used with this extended functionality:

Field	Version
<b>aio_priority</b>	AIOCBX_VERS1
<b>aio_cache_hint</b>	AIOCBX_VERS1

The **aio\_priority** and **aio\_cache\_hint** values take effect only on a 64-bit kernel under the following conditions:

- The file descriptor being operated on by AIO belongs to the raw character interface of an LVM logical volume.
- The LVM logical volume resides on a device that supports I/O priorities and cache hints.

The **aio\_read**, **aio\_read64**, **aio\_write**, **aio\_write64**, **lio\_listio**, and **lio\_listio64** interfaces are all compatible with an extended AIOCB. Other interfaces (such as **aio\_cancel**) ignore the extended fields.

The valid values for **aio\_priority** and **aio\_cache\_hint** are described in the **sys/extendio.h** file. The **aio\_priority** must be either **IOPRIORITY\_UNSET** (0) or a value from 1 to 15. Lower I/O priority values are considered to be more important than higher values. For example, a value of 1 is considered highest priority and a value of 15 is considered lowest priority. The **aio\_cache\_hint** must be either **CH\_AGE\_OUT\_FAST** or **CH\_PAGE\_WRITE**. These cache hint values are mutually exclusive. If **CH\_AGE\_OUT\_FAST** is set, the I/O buffer can be aged out quickly from the storage device buffer cache. This is useful in situations where the application is already caching the I/O buffer and redundant caching within the storage layer can be avoided. If **CH\_PAGE\_WRITE** is set, the I/O buffer is written only to the storage device cache and not to the disk.

## Using I/O Completion Ports with AIO Requests

To use I/O completion ports (IOCP) with AIO requests, set the **AIO\_EXTENDED** flag in the **aio\_flags** field and the **aio\_version** field to a value of **AIOCBX\_VERS2** or higher. All other extended fields defined must be set to 0 if they are not used. The following fields are used with this extended functionality:

Field	Version
<b>aio_iocpfd</b>	AIOCBX_VERS2

A limitation of the AIO interface that is used in a threaded environment is that **aio\_nwait()** collects completed I/O requests for ALL threads in the same process. In other words, one thread collects completed I/O requests that are submitted by another thread. Another problem is that multiple threads cannot invoke the collection routines (such as **aio\_nwait()**) at the same time. If one thread issues **aio\_nwait()** while another thread is calling it, the second **aio\_nwait()** returns **EBUSY**. This limitation can affect I/O performance when many I/Os must run at the same time and a single thread cannot run fast enough to collect all the completed I/Os.

Using I/O completion ports with AIO requests provides the capability for an application to capture results of various AIO operations on a per-thread basis in a multithreaded environment. This functionality provides threads with a method of receiving completion status for only the AIO requests initiated by the thread.

The IOCP subsystem only provides completion status by generating completion packets for AIO requests. The I/O cannot be submitted for regular files through IOCP.

The current behavior of AIO remains unchanged. An application is free to use any existing AIO interfaces in combination with I/O completion ports. The application is responsible for "harvesting" completion packets for any noncanceled AIO requests that it has associated with a completion port.

The application must associate a file with a completion port using the **CreateloCompletionPort** IOCP routine. The file can be associated with multiple completion ports, and a completion port can have multiple files associated with it. When making the association, the application must use an application-defined *CompletionKey* to differentiate between AIO completion packets and socket completion packets. The application can use different *CompletionKeys* to differentiate among individual files (or in any other manner) as necessary.

The application must also associate AIO requests with the same completion port as the corresponding file. It does this by initializing the **aio\_iocpfd** of the AIOCB with the file descriptor of the completion port. An AIOCB can be associated with only one completion port, but a completion port can have multiple AIOCBs associated with it. The association between a completion port and an AIOCB must be done *before* the request is made. This is accomplished using an AIO routine, such as **aio\_write**, **aio\_read**, or **lio\_listio**. If the value in the **aio\_iocpfd** field is not a valid completion port file descriptor, the attempt to start the request fails and no I/O is performed.

An association must be made directly between a completion port and an AIOCB. For example, if you want to call **lio\_listio()**, each AIOCB in the **lio\_listio** chain must be associated individually prior to the call. It is not necessary to have all AIOCBs in the chain associated with a completion port.

After an association is made, it remains until the application explicitly clears it by using a value of 0 for the **aio\_iocpfd** field, or the AIOCB is destroyed. A completion packet is created only when I/O completes for an AIOCB that has been associated with a completion port.

A summary of the steps that an application takes to use I/O completion ports with AIO requests is as follows:

1. Opens a regular file for I/O.
2. Calls the **CreateloCompletionPort** routine to create an I/O completion port (IOCP), using the file descriptor for the regular file and an application-defined *CompletionKey*, which is used to differentiate AIO requests from socket I/O. The **CreateloCompletionPort** function returns an IOCP file descriptor that corresponds to the newly created IOCP.
3. Allocates and clears (using the **bzero** function) an AIO control block. Indicates that I/O completion ports are to be used with AIO requests by setting the **AIO\_EXTENDED** flag of the AIOCB's **aio\_flags** field. Also sets the **aio\_version** field to a value of **AIOCBX\_VERS2** or higher.
4. Associates the AIO request with the IOCP by initializing the **aio\_iocpfd** field in the AIOCB to contain the IOCP file descriptor returned by the **CreateloCompletionPort** routine.
5. Starts the AIO request using existing AIO interfaces. Multiple requests can be started using the **lio\_listio** interface.
6. Calls the **GetQueuedCompletionStatus** function with the IOCP file descriptor to collect the results of the completed AIO requests on a particular IOCP. The application provides the address of a pointer in the **LPOVERLAPPED** argument to **GetQueuedCompletionStatus**, so the corresponding AIOCB pointer can be returned. Details of the AIO request can be determined by examining the returned AIOCB.
7. After all I/O is complete, the application is responsible for closing all file descriptors.

---

## Related Information

### Subroutine References

The **aio\_cancel** or **aio\_cancel64** subroutine, **aio\_error** or **aio\_error64** subroutine, **aio\_read** or **aio\_read64** subroutine, **aio\_return** or **aio\_return64** subroutine, **aio\_suspend** or **aio\_suspend64** subroutine, **aio\_write** or **aio\_write64** subroutine, **lio\_listio** or **lio\_listio64** subroutine in *AIX Version 6.1 Technical Reference: Base Operating System and Extensions Volume 1*.

### Commands References

The **chdev** command in *AIX Version 6.1 Commands Reference, Volume 1*.

The **mkdev** command in *AIX Version 6.1 Commands Reference, Volume 3*.

The **rmdev** command in *AIX Version 6.1 Commands Reference, Volume 4*.



---

## Chapter 6. Device Configuration Subsystem

Devices are usually pieces of equipment that attach to a computer. Devices include printers, adapters, and disk drives. Additionally, devices are special files that can handle device-related tasks.

System users cannot operate devices until device configuration occurs. To configure devices, the Device Configuration Subsystem is available.

Read about general configuration characteristics and procedures in:

- “Scope of Device Configuration Support”
- “Device Configuration Subsystem Overview”
- “General Structure of the Device Configuration Subsystem” on page 94

---

### Scope of Device Configuration Support

The term *device* has a wider range of meaning in this operating system than in traditional operating systems. Traditionally, *devices* refers to hardware components such as disk drives, tape drives, printers, and keyboards. Pseudo-devices, such as the console, **error** special file, and **null** special file, are also included in this category. However, in this operating system, all of these devices are referred to as *kernel devices*, which have device drivers and are known to the system by major and minor numbers.

Also, in this operating system, hardware components such as buses, adapters, and enclosures (including racks, drawers, and expansion boxes) are considered devices.

---

### Device Configuration Subsystem Overview

Devices are organized hierarchically within the system. This organization requires lower-level device dependence on upper-level devices in child-parent relationships. The system device (sys0) is the highest-level device in the system node, which consists of all physical devices in the system.

Each device is classified into functional classes, functional subclasses and device types (for example, printer *class*, parallel *subclass*, 4201 Proprinter *type*). These classifications are maintained in the device configuration databases with all other device information.

The Device Configuration Subsystem consists of:

#### High-level Commands

Maintain (add, delete, view, change) configured devices within the system. These commands manage all of the configuration functions and are performed by invoking the appropriate device methods for the device being configured. These commands call device methods and low-level commands.

#### Device Methods

The system uses the high-level **Configuration Manager (cfgmgr)** command used to invoke automatic device configurations through system boot phases and the user can invoke the command during system run time. *Configuration rules* govern the **cfgmgr** command.

#### Database

Define, configure, change, unconfigure, and undefine devices. The device methods are used to identify or change the device *states* (operational modes). Maintains data through the *ODM* (Object Data Manager) by object classes. Predefined Device Objects contain configuration data for all devices that can possibly be used by the system. Customized Device Objects contain data for *device instances* that are actually in use by the system.

---

## General Structure of the Device Configuration Subsystem

The Device Configuration Subsystem can be viewed from the following different levels:

- High-level perspective
- Device method level
- Low-level perspective

Data that is used by the three levels is maintained in the *Configuration database*. The database is managed as object classes by the Object Data Manager (ODM). All information relevant to support the device configuration process is stored in the configuration database.

The system cannot use any device unless it is configured.

The database has two components: the Predefined database and the Customized database. The *Predefined database* contains configuration data for all devices that could possibly be supported by the system. The *Customized database* contains configuration data for the devices actually defined or configured in that particular system.

The *Configuration manager* (**cfgmgr** command) performs the configuration of a system's devices automatically when the system is booted. This high-level program can also be invoked through the system keyboard to perform automatic device configuration. The configuration manager command configures devices as specified by *Configuration rules*.

### High-Level Perspective

From a high-level, user-oriented perspective, device configuration comprises the following basic tasks:

- Adding a device to the system
- Deleting a device from the system
- Changing the attributes of a device
- Showing information about a device

From a high-level, system-oriented perspective, device configuration provides the basic task of automatic device configuration: running the configuration manager program.

A set of high-level commands accomplish all of these tasks during run time: **chdev**, **mkdev**, **lsattr**, **lsconn**, **lsdev**, **lsparent**, **rmdev**, and **cfgmgr**. High-level commands can invoke device methods and low-level commands.

### Device Method Level

Beneath the high-level commands (including the **cfgmgr** Configuration Manager program) is a set of functions called *device methods*. These methods perform well-defined configuration steps, including these five functions:

- Defining a device in the configuration database
- Configuring a device to make it available
- Changing a device to make a change in its characteristics
- Unconfiguring a device to make it unavailable
- Undefined a device from the configuration database

“Understanding Device States” on page 99 discusses possible device states and how the various methods affect device state changes.



The high-level device commands (including **cfgmgr**) can use the device methods. These methods insulate high-level configuration programs from kernel-specific, hardware-specific, and device-specific configuration steps. Device methods can invoke low-level commands.

## Low-Level Perspective

Beneath the device methods is a set of low-level library routines that can be directly called by device methods as well as by high-level configuration programs.

---

## Device Configuration Database Overview

The Configuration database is an object-oriented database. The Object Data Manager (ODM) provides facilities for accessing and manipulating it through object classes.

The following databases are used in the configuration process:

<b>Predefined database</b>	Contains information about all possible types of devices that can be defined for the system.
<b>Customized database</b>	Describes all devices currently defined for use in the system. Items are referred to as <i>device instances</i> .

ODM Device Configuration Object Classes in *AIX Version 6.1 Technical Reference: Kernel and Subsystems Volume 2* provides access to the object classes that make up the Predefined and Customized databases.

Devices must be defined in the database for the system to make use of them. For a device to be in the Defined state, the Configuration database must contain a complete description of it. This information includes items such as the device driver name, the device major and minor numbers, the device method names, the device attributes, connection information, and location information.

---

## Basic Device Configuration Procedures Overview

At system boot time, **cfgmgr** is automatically invoked to configure all devices detected as well as any device whose device information is stored in the Configuration database. At run time, you can configure a specific device by directly invoking (or indirectly invoking through a usability interface layer) high-level device commands.

High-level device commands invoke methods and allow the user to add, delete, show, and change devices and their associated attributes.

When a specific device is defined through its define method, the information from the Predefined database for that type of device is used to create the information describing the specific device instance. This specific device instance information is then stored in the Customized database. For more information on define methods, see *Writing a Define Method in AIX Version 6.1 Technical Reference: Kernel and Subsystems Volume 2*.

The process of configuring a device is often highly device-specific. The configure method for a kernel device must:

- Load the device's driver into the kernel.
- Pass the device dependent structure (DDS) describing the device instance to the driver. For more information on DDS, see "Device Dependent Structure (DDS) Overview" on page 103.
- Create a special file for the device in the **/dev** directory. For more information, see *Special Files in AIX Version 6.1 Files Reference*.

For more information on configure methods, see *Writing a Configure Method in AIX Version 6.1 Technical Reference: Kernel and Subsystems Volume 2*.

Of course, many devices do not have device drivers. For this type of device the configured state is not as meaningful. However, it still has a Configure method that simply marks the device as configured or performs more complex operations to determine if there are any devices attached to it.

The configuration process requires that a device be defined or configured before a device attached to it can be defined or configured. At system boot time, the Configuration Manager first configures the system device. The remaining devices are configured by traversing down the parent-child connections layer by layer. The Configuration Manager then configures any pseudo-devices that need to be configured.

---

## Device Configuration Manager Overview

The Configuration Manager is a rule-driven program that automatically configures devices in the system during system boot and run time. When the Configuration Manager is invoked, it reads rules from the Configuration Rules object class and performs the indicated actions. For more information on Configuration Rules, see *Configuration Rules (Config\_Rules) Object Class in AIX Version 6.1 Technical Reference: Kernel and Subsystems Volume 2*.

Devices in the system are organized in clusters of tree structures known as *nodes*. Each tree is a logical subsystem by itself. For example, the system node consists of all the physical devices in the system. The top of the node is the system device. Below the bus and connected to it are the adapters. The bottom of the hierarchy contains devices to which no other devices are connected. Most pseudo-devices, including low -function terminal (LFT) and pseudo-terminal (pty) devices, are organized as separate tree structures or nodes.

## Devices Graph

See “Understanding Device Dependencies and Child Devices” on page 101 for more information.

## Configuration Rules

Each rule in the Configuration Rules (Config\_Rules) object class specifies a program name that the Configuration Manager must execute. These programs are typically the configuration programs for the devices at the top of the nodes. When these programs are invoked, the names of the next lower-level devices that need to be configured are returned.

The Configuration Manager configures the next lower-level devices by invoking the configuration methods for those devices. In turn, those configuration methods return a list of to-be-configured device names. The process is repeated until no more device names are returned. As a result, all devices in the same node are configured in transverse order. The following are different types of rules:

- Phase 1
- Phase 2
- Service

The system boot process is divided into two phases. In each phase, the Configuration Manager is invoked. During phase 1, the Configuration Manager is called with a **-f** flag, which specifies that *phase = 1* rules are to be executed. This results in the configuration of base devices into the system, so that the root file system can be used. During phase 2, the Configuration Manager is called with a **-s** flag, which specifies that *phase = 2* rules are to be executed. This results in the configuration of the rest of the devices into the system.

For more information on the booting process, see *Understanding System Boot Processing in Operating system and device management*.

The Configuration Manager invokes the programs in the order specified by the sequence value in the rule. In general, the lower the sequence number within a given phase, the higher the priority. Thus, a rule with a 2 sequence number is executed before a rule with a sequence number of 5. An exception is made for 0 sequence numbers, which indicate a don't-care condition. Any rule with a sequence number of 0 is executed last. The Configuration Rules (Config\_Rules) object class provides an example of this process.

If device names are returned from the program invoked, the Configuration Manager finishes traversing the node tree before it invokes the next program. Note that some program names might not be associated with any devices, but they must be included to configure the system.

## Invoking the Configuration Manager

During system boot time, the Configuration Manager is run in two phases. In phase 1, it configures the base devices needed to successfully start the system. These devices include the root volume group, which permits the configuration database to be read in from the root file system.

In phase 2, the Configuration Manager configures the remaining devices using the configuration database from the root file system. During this phase, different rules are used, depending on whether the system was booted in normal mode or in service mode. If the system is booted in service mode, the rules for service mode are used. Otherwise, the phase 2 rules are used.

The Configuration Manager can also be invoked during run time to configure all the detectable devices that might have been turned off at system boot or added after the system boot. In this case, the Configuration Manager uses the phase 2 rules.

---

## Device Classes, Subclasses, and Types Overview

To manage the wide variety of devices it supports more easily, the operating system classifies them hierarchically. One advantage of this arrangement is that device methods and high-level commands can operate against a whole set of similar devices.

Devices are categorized into the following main groups:

- Functional classes
- Functional subclasses
- Device types

Devices are organized into a set of *functional classes* at the highest level. From a user's point of view, all devices belonging to the same class perform the same functions. For example, all printer devices basically perform the same function of generating printed output.

However, devices within a class can have different interfaces. A class can therefore be partitioned into a set of *functional subclasses* in which devices belonging to the same subclass have similar interfaces. For example, serial printers and parallel printers form two subclasses of printer devices.

Finally, a device subclass is a collection of *device types*. All devices belonging to the same device type share the same manufacturer's model name and number. For example, 3812-2 (model 2 Pageprinter) and 4201 (Proprinter II) printers represent two types of printers.

Devices of the same device type can be managed by different drivers if the type belongs to more than one subclass. For example, the 4201 printer belongs to both the serial interface and parallel interface subclasses of the printer class, although there are different drivers for the two interfaces. However, a device of a particular class, subclass, and type can be managed by only one device driver.

Devices in the system are organized in clusters of tree structures known as *nodes*. For example, the system node consists of all the physical devices in the system. At the top of the node is the system

device. Below the bus and connected to it are the adapters. The bottom of the hierarchy contains the devices to which no other devices are connected. Most pseudo-devices, including LFT and PTY, are organized as separate nodes.

---

## Writing a Device Method

*Device methods* are programs associated with a device that perform basic device configuration operations. These operations consist of defining, undefining, configuring, unconfiguring, and reconfiguring a device. Some devices also use optional start and stop operations.

The following are the basic device methods:

<b>Define</b>	Creates a device instance in the Customized database.
<b>Configure</b>	Configures a device instance already represented in the Customized database. This method is responsible for making a device available for use in the system.
<b>Change</b>	Reconfigures a device by allowing device characteristics or attributes to be changed.
<b>Unconfigure</b>	Makes a configured device unavailable for use in the system. The device instance remains in the Customized database but must be reconfigured before it can be used.
<b>Undefine</b>	Deletes a device instance from the Customized database.

## Invoking Methods

One device method can invoke another device method. For instance, a Configure method for a device may need to invoke the Define method for child devices. The Change method can invoke the Unconfigure and Configure methods. To ensure proper operation, a method that invokes another method must always use the `odm_run_method` subroutine.

## Example Methods

See the `/usr/samples` directory for example device method source code. These source code excerpts are provided for example purposes only. The examples do not function as written.

---

## Understanding Device Methods Interfaces

Device methods are not executed directly from the command line. They are only invoked by the Configuration Manager at boot time or by the `cfgmgr`, `mkdev`, `chdev`, and `rmdev` configuration commands at run time. As a result, any device method you write should meet well-defined interfaces.

The parameters that are passed into the methods as well as the exit codes returned must both satisfy the requirements for each type of method. Additionally, some methods must write information to the `stdout` and `stderr` files.

These interfaces are defined for each of the device methods in the individual articles on writing each method.

To better understand how these interfaces work, one needs to understand, at least superficially, the flow of operations through the Configuration Manager and the run-time configuration commands.

## Configuration Manager

The Configuration Manager begins by invoking a Node Configuration program listed in one of the rules in the Configuration Rules (Config\_Rules) object class. A node is a group of devices organized into a tree structure representing the various interconnections of the devices. The Node Configuration program is responsible for starting the configuration process for a node. It does this by querying the Customized database to see if the device at the top of the node is represented in the database. If so, the program writes the logical name of the device to the `stdout` file and then returns to the Configuration Manager.

The Configuration Manager intercepts the Node Configuration program's **stdout** file to obtain the name of the device that was written. It then invokes the Configure method for that device. The device's Configure method performs the steps necessary to make the device available. If the device is not an intermediate one, the Configure method simply returns to the Configuration Manager. However, if the device is an intermediate device that has child devices, the Configure method must determine whether any of the child devices need to be configured. If so, the Configure method writes the names of all the child devices to be configured to the **stdout** file and then returns to the Configuration Manager.

The Configuration Manager intercepts the Configure method's **stdout** file to retrieve the names of the children. It then invokes, one at a time, the Configure methods for each child device. Each of these Configure methods operates as described for the parent device. For example, it might simply exit when complete, or write to its **stdout** file a list of additional device names to be configured and then exit. The Configuration Manager will continue to intercept the device names written to the **stdout** file and to invoke the Configure methods for those devices until the Configure methods for all the devices have been run and no more names are written to the **stdout** file.

## Run-Time Configuration Commands

User configuration commands invoke device methods during run time.

**mkdev** The **mkdev** command is invoked to define or configure, or define and configure, devices at run time. If just defining a device, the **mkdev** command invokes the Define method for the device. The Define method creates the customized device instance in the Customized Devices (CuDv) object class and writes the name assigned to the device to the **stdout** file. The **mkdev** command intercepts the device name written to the **stdout** file by the Define method to learn the name of the device. If user-specified attributes are supplied with the **-a** flag, the **mkdev** command then invokes the Change method for the device.

If defining and configuring a device, the **mkdev** command invokes the Define method, gets the name written to the **stdout** file with the Define method, invokes the Change method for the device if user-specified attributes were supplied, and finally invokes the device's Configure method.

If only configuring a device, the device must already exist in the CuDv object class and its name must be specified to the **mkdev** command. In this case, the **mkdev** command simply invokes the Configure method for the device.

**chdev** The **chdev** command is used to change the characteristics, or attributes, of a device. The device must already exist in the CuDv object class, and the name of the device must be supplied to the **chdev** command. The **chdev** command simply invokes the Change method for the device.

**rmdev** The **rmdev** command can be used to undefine or unconfigure, or unconfigure and undefine, a device. In all cases, the device must already exist in the CuDv object class and the name of the device must be supplied to the **rmdev** command. The **rmdev** command then invokes the Undefine method, the Unconfigure method, or the Unconfigure method followed by the Undefine method, depending on the function requested by the user.

**cfgmgr** The **cfgmgr** command can be used to configure all detectable devices that did not get configured at boot time. This might occur if the devices had been powered off at boot time. The **cfgmgr** command is the Configuration Manager and operates in the same way at run time as it does at boot time. The boot time operation is described in Device Configuration Manager Overview .

---

## Understanding Device States

Device methods are responsible for changing the state of a device in the system. A device can be in one of four states as represented by the Device Status Flag descriptor in the device's object in the Customized Devices (CuDv) object class.

<b>Defined</b>	Represented in the Customized database, but neither configured nor available for use in the system.
<b>Available</b>	Configured and available for use.
<b>Undefined</b>	Not represented in the Customized database.

**Stopped** Configured, but not available for use by applications. (Optional state)  
**Note:** Start and stop methods are only supported on the **inet0** device.

The Define method is responsible for creating a device instance in the Customized database and setting the state to Defined. The Configure method performs all operations necessary to make the device usable and then sets the state to Available.

The Change method usually does not change the state of the device. If the device is in the Defined state, the Change method applies all changes to the database and leaves the device defined. If the device is in the Available state, the Change method attempts to apply the changes to both the database and the actual device, while leaving the device available. However, if an error occurs when applying the changes to the actual device, the Change method might need to unconfigure the device, thus changing the state to Defined.

Any Unconfigure method you write must perform the operations necessary to make a device unusable. Basically, this method undoes the operations performed by the Configure method and sets the device state to Defined. Finally, the Undefine method actually deletes all information for a device instance from the Customized database, thus reverting the instance to the Undefined state.

The Stopped state is an optional state that some devices require. A device that supports this state needs Start and Stop methods. The Stop method changes the state from Available to Stopped. The Start method changes it from Stopped back to Available.

**Note:** Start and stop methods are only supported on the **inet0** device.

---

## Adding an Unsupported Device to the System

The operating system provides support for a wide variety of devices. However, some devices are not currently supported. You can add a currently unsupported device only if you also add the necessary software to support it.

To add a currently unsupported device to your system, you might need to:

- Modify the Predefined database
- Add appropriate device methods
- Add a device driver
- Use **installp** procedures

## Modifying the Predefined Database

To add a currently unsupported device to your system, you must modify the Predefined database. To do this, you must add information about your device to three predefined object classes:

- Predefined Devices (PdDv) object class
- Predefined Attribute (PdAt) object class
- Predefined Connection (PdCn) object class

To describe the device, you must add one object to the PdDv object class to indicate the class, subclass, and device type. You must also add one object to the PdAt object class for each device attribute, such as interrupt level or block size. Finally, you must add objects to the PdCn object class if the device is an intermediate device. If the device is an intermediate device, you must add an object for each different connection location on the intermediate device.

You can use the **odmadd** Object Data Manager (ODM) command from the command line or in a shell script to populate the necessary Predefined object classes from stanza files.

The Predefined database is populated with devices that are present at the time of installation. For some supported devices, such as serial and parallel printers and SCSI disks, the database also contains generic device objects. These generic device objects can be used to configure other similar devices that are not explicitly supported in the Predefined database. If new devices are added after installation, additional device support might need to be installed.

For example, if you have a serial printer that closely resembles a printer supported by the system, and the system's device driver for serial printers works on your printer, you can add the device driver as a printer of type **osp** (other serial printer). If these generic devices successfully add your device, you do not need to provide additional system software.

## Adding Device Methods

You must add device methods when adding system support for a new device. Primary methods needed to support a device are:

- Define
- Configure
- Change
- Undefine
- Unconfigure

When adding a device that closely resembles devices already supported, you might be able to use one of the methods of the already supported device. For example, if you are adding a new type of SCSI disk whose interfaces are identical to supported SCSI disks, the existing methods for SCSI disks may work. If so, all you need to do is populate the Predefined database with information describing the new SCSI disk, which will be similar to information describing a supported SCSI disk.

If you need instructions on how to write a device method, see [Writing a Device Method](#) .

## Adding a Device Driver

If you add a new device, you will probably need to add a device driver. However, if you are adding a new device that closely resembles an already supported device, you might be able to use the existing device driver. For example, when you are adding a new type of SCSI disk whose interfaces are identical to supported SCSI disks, the existing SCSI disk device driver might work.

## Using installp Procedures

The **installp** procedures provide a method for adding the software and Predefined information needed to support your new device. You might need to write shell scripts to perform tasks such as populating the Predefined database.

---

## Understanding Device Dependencies and Child Devices

The dependencies that one device has on another can be represented in the configuration database in two ways. One way usually represents physical connections such as a keyboard device connected to a particular keyboard adapter. The keyboard device has a dependency on the keyboard adapter in that it cannot be configured until after the adapter is configured. This relationship is usually referred to as a parent-child relationship, with the adapter as parent and the keyboard device as child. These relationships are represented with the Parent Device Logical Name and Location Where Device Is Connected descriptors in the Customized Devices (CuDv) object class.

The second method represents a logical connection. A device method can add an object identifying both a dependent device and the device upon which it depends to the Customized Dependency (CuDep) object class. The dependent device is considered to *have* a dependency, and the depended-upon device is

considered to *be* a dependency. CuDep objects are usually added to the database to represent a situation in which one device requires access to another device. For example, the lft0 device depends upon a particular keyboard or display device.

These two types of dependencies differ significantly. The configuration process uses parent-child dependencies at boot time to configure all devices that make up a node. The CuDep dependency is usually only used by a device's Configure method to record the names of the devices on which it depends.

For device methods, the parent-child relationship is the more important. Parent-child relationships affect device-method activities in these ways:

- A parent device cannot be unconfigured if it has a configured child.
- A parent device cannot be undefined if it has a defined or configured child.
- A child device cannot be defined if the parent is not defined or configured.
- A child device cannot be configured if the parent is not configured.
- A parent device's configuration cannot be changed if it has a configured child. This guarantees that the information about the parent that the child's device driver might be using remains valid.

However, when a device is listed as a dependency of another device in the CuDep object class, the only effect is to prevent the depended-upon device from being undefined. The name of the dependency is important to the dependent device. If the depended-upon device were allowed to be undefined, a third device could be defined and assigned the same name.

Writers of Unconfigure and Change methods for a depended-upon device should not worry about whether the device is listed as a dependency. If the depended-upon device is actually open by the other device, the Unconfigure and Change operations will fail because their device is busy. But if the depended-upon device is *not* currently open, the Unconfigure or Change operations can be performed without affecting the dependent device.

The possible parent-child connections are defined in the Predefined Connection (PdCn) object class. Each predefined device type that can be a parent device is represented in this object class. There is an object for each connection location (such as slots or ports) describing the subclass of devices that can be connected at that location. The subclass is used to identify each device because it indicates the devices' connection type (for example, SCSI or rs232).

There is no corresponding predefined object class describing the possible CuDep dependencies. A device method can be written so that it already knows what the dependencies are. If predefined data is required, it can be added as predefined attributes for the dependent device in the Predefined Attribute (PdAt) object class.

---

## Accessing Device Attributes

The predefined device attributes for each type of predefined device are stored in the Predefined Attribute (PdAt) object class. The objects in the PdAt object class identify the default values as well as other possible values for each attribute. The Customized Attribute (CuAt) object class contains only attributes for customized device instances that have been changed from their default values.

When a customized device instance is created by a Define method, its attributes assume the default values. As a result, no objects are added to the CuAt object class for the device. If an attribute for the device is changed from the default value by the Change method, an object to describe the attribute's current value is added to the CuAt object class for the attribute. If the attribute is subsequently changed back to the default value, the Change method deletes the CuAt object for the attribute.

Any device methods that need the current attribute values for a device must access both the PdAt and CuAt object classes. If an attribute appears in the CuAt object class, then the associated object identifies the current value. Otherwise, the default value from the PdAt attribute object identifies the current value.



## Modifying an Attribute Value

When modifying an attribute value, methods you write must obtain the objects for that attribute from both the PdAt and CuAt object classes.

Any method you write must be able to handle the following four scenarios:

- If the new value differs from the default value and no object currently exists in the CuAt object class, any method you write must add an object into the CuAt object class to identify the new value.
- If the new value differs from the default value and an object already exists in the CuAt object class, any method you write must update the CuAt object with the new value.
- If the new value is the same as the default value and an object exists in the CuAt object class, any method you write must delete the CuAt object for the attribute.
- If the new value is the same as the default value and no object exists in the CuAt object class, any method you write does not need to do anything.

Your methods can use the **getattr** and **putattr** subroutines to get and modify attributes. The **getattr** subroutine checks both the PdAt and CuAt object classes before returning an attribute to you. It always returns the information in the form of a CuAt object even if returning the default value from the PdAt object class.

Use the **putattr** subroutine to modify these attributes.

---

## Device Dependent Structure (DDS) Overview

A *device dependent structure* (DDS) contains information that describes a device instance to the device driver. It typically contains information about device-dependent attributes as well as other information the driver needs to communicate with the device. In many cases, information about a device's parent is included. (For instance, a driver needs information about the adapter and the bus the adapter is plugged into to communicate with a device connected to an adapter.)

A device's DDS is built each time the device is configured. The Configure method can fill in the DDS with fixed values, computed values, and information from the Configuration database. Most of the information from the Configuration database usually comes from the attributes for the device in the Customized Attribute (CuAt) object class, but can come from any of the object classes. Information from the database for the device's parent device or parent's parent device can also be included. The DDS is passed to the device driver with the **SYS\_CFGDD** flag of the **sysconfig** subroutine, which calls the device driver's **ddconfig** subroutine with the **CFG\_INIT** command.

## How the Change Method Updates the DDS

The Change method is invoked when changing the configuration of a device. The Change method must ensure consistency between the Configuration database and the view that any device driver might have of the device. This is accomplished by:

1. Not allowing the configuration to be changed if the device has configured children; that is, children in either the Available or Stopped states. This ensures that a DDS built using information in the database about a parent device remains valid because the parent cannot be changed.
2. If a device has a device driver and the device is in either the Available or Stopped state, the Change method must communicate to the device driver any changes that would affect the DDS. This can be accomplished with **ioctl** operations, if the device driver provides the support to do so. It can also be accomplished by taking the following steps:
  - a. Terminating the device instance by calling the **sysconfig** subroutine with the **SYS\_CFGDD** operation. This operation calls the device driver's **ddconfig** subroutine with the **CFG\_TERM** command.
  - b. Rebuilding the DDS using the changed information.

- c. Passing the new DDS to the device driver by calling the **sysconfig** subroutine **SYS\_CFGDD** operation. This operation then calls the **ddconfig** subroutine with the **CFG\_INIT** command.

Many Change methods simply invoke the device's Unconfigure method, apply changes to the database, and then invoke the device's Configure method. This process ensures the two stipulated conditions since the Unconfigure method, and thus the change, will fail, if the device has Available or Stopped children. Also, if the device has a device driver, its Unconfigure method terminates the device instance. Its Configure method also rebuilds the DDS and passes it to the driver.

## Guidelines for DDS Structure

There is no single defined DDS format. Writers of device drivers and device methods must agree upon a particular device's DDS format. When obtaining information about a parent device, you might want to group that information together in the DDS.

When building a DDS for a device connected to an adapter card, you will typically need the following adapter information:

<b>slot number</b>	Obtained from the <b>connwhere</b> descriptor of the adapter's Customized Device (CuDv) object.
<b>bus resources</b>	Obtained from attributes for the adapter in the Customized Attribute (CuAt) or Predefined Attribute (PdAt) object classes. These include attributes for bus interrupt levels, interrupt priorities, bus memory addresses, bus I/O addresses, and DMA arbitration levels.

The following attribute must be obtained for the adapter's parent bus device:

<b>bus_id</b>	Identifies the I/O bus. This field is needed by the device driver to access the I/O bus.
---------------	--

**Note:** The **getattr** device configuration subroutine should be used whenever attributes are obtained from the Configuration database. This subroutine returns the Customized attribute value if the attribute is represented in the Customized Attribute object class. Otherwise, it returns the default value from the Predefined Attribute object class.

Finally, a DDS generally includes the device's logical name. This is used by the device driver to identify the device when logging an error for the device.

## Example of DDS

The following example provides a guide for using DDS format.

```
/* Device DDS */
struct device_dds {
    /* Bus information */
    ulong bus_id;          /* I/O bus id          */
    ushort us_type;       /* Bus type, i.e. BUS_MICRO_CHANNEL*/
    /* Adapter information */
    int slot_num;         /* Slot number         */
    ulong io_addr_base;   /* Base bus i/o address */
    int bus_intr_lvl;     /* bus interrupt level  */
    int intr_priority;    /* System interrupt priority */
    int dma_lvl;         /* DMA arbitration level */
    /* Device specific information */
    int block_size;       /* Size of block in bytes */
    int abc_attr;         /* The abc attribute    */
    int xyz_attr;         /* The xyz attribute    */
    char resource_name[16]; /* Device logical name  */
};
```

---

## List of Device Configuration Commands

The high-level device configuration commands are:

<b>chdev</b>	Changes a device's characteristics.
<b>lsdev</b>	Displays devices in the system and their characteristics.
<b>mkdev</b>	Adds a device to the system.
<b>rmdev</b>	Removes a device from the system.
<b>lsattr</b>	Displays attribute characteristics and possible values of attributes for devices in the system.
<b>lscnn</b>	Displays the connections a given device, or kind of device, can accept.
<b>lsparent</b>	Displays the possible parent devices that accept a specified connection type or device.
<b>cfgmgr</b>	Configures devices by running the programs specified in the Configuration Rules (Config_Rules) object class.

Associated commands are:

<b>bootlist</b>	Alters the list of boot devices seen by ROS when the machine boots.
<b>lscfg</b>	Displays diagnostic information about a device.
<b>restbase</b>	Reads the base customized information from the boot image and restores it into the Device Configuration database used during system boot phase 1.
<b>savebase</b>	Saves information about base customized devices in the Device Configuration Database onto the boot device.

---

## List of Device Configuration Subroutines

Following are the preexisting conditions for using the device configuration library subroutines:

- The caller has initialized the Object Data Manager (ODM) before invoking any of these library subroutines. This is done using the **initialize\_odm** subroutine. Similarly, the caller must terminate the ODM (using the **terminate\_odm** subroutine) after these library subroutines have completed.
- Because all of these library subroutines (except the **attrval**, **getattr**, and **putattr** subroutines) access the Customized Device Driver (CuDvDr) object class, this class must be exclusively locked and unlocked at the proper times. The application does this by using the **odm\_lock** and **odm\_unlock** subroutines. In addition, those library subroutines that access the CuDvDr object class exclusively lock this class with their own internal locks.

Following are the device configuration library subroutines:

<b>attrval</b>	Verifies that attributes are within range.
<b>genmajor</b>	Generates the next available major number for a device driver instance.
<b>genminor</b>	Generates the smallest unused minor number, a requested minor number for a device if it is available, or a set of unused minor numbers.
<b>genseq</b>	Generates a unique sequence number for creating a device's logical name.
<b>getattr</b>	Returns attribute objects from either the Predefined Attribute (PdAt) or Customized Attribute (CuAt) object class, or both.
<b>getminor</b>	Gets from the CuDvDr object class the minor numbers for a given major number.
<b>loadext</b>	Loads or unloads and binds or unbinds device drivers to or from the kernel.
<b>putattr</b>	Updates attribute information in the CuAt object class or creates a new object for the attribute information.
<b>reldevno</b>	Releases the minor number or major number, or both, for a device instance.
<b>relmajor</b>	Releases the major number associated with a specific device driver instance.

---

## Related Information

Understanding System Boot Processing in *Operating system and device management*

Special Files in *AIX Version 6.1 Files Reference*

Initial Printer Configuration in *Printers and printing*

Machine Device Driver, Loading a Device Driver in *AIX Version 6.1 Technical Reference: Kernel and Subsystems Volume 2*.

Writing a Define Method, Writing a Configure Method, Writing a Change Method, Writing an Unconfigure Method, Writing an Undefine Method, Writing Optional Start and Stop Methods, How Device Methods Return Errors, Device Methods for Adapter Cards: Guidelines in *AIX Version 6.1 Technical Reference: Kernel and Subsystems Volume 2*

Configuration Rules (Config\_Rules) Object Class, Customized Dependency (CuDep) Object Class, Customized Devices (CuDv) Object Class, Predefined Attribute (PdAt) Object Class, Predefined Connection (PdCn) Object Class, Adapter-Specific Considerations For the Predefined Devices (PdDv) Object Class, Adapter-Specific Considerations For the Predefined Attributes (PdAt) Object Class, Predefined Devices Object Class, ODM Device Configuration Object Classes in *AIX Version 6.1 Technical Reference: Kernel and Subsystems Volume 2*.

## Subroutine References

The **getattr** subroutine, **ioctl** subroutine, **odm\_run\_method** subroutine, **putattr** subroutine in *AIX Version 6.1 Technical Reference: Base Operating System and Extensions Volume 1*.

The **sysconfig** subroutine in *AIX Version 6.1 Technical Reference: Base Operating System and Extensions Volume 2*.

## Commands References

The **cfgmgr** command, **chdev** command in *AIX Version 6.1 Commands Reference, Volume 1*.

The **mkdev** command in *AIX Version 6.1 Commands Reference, Volume 3*.

The **rmdev** command in *AIX Version 6.1 Commands Reference, Volume 4*.

## Technical References

The SYS\_CFGDD **sysconfig** operation in *AIX Version 6.1 Technical Reference: Base Operating System and Extensions Volume 1*.

The **ddconfig** device driver entry point in *AIX Version 6.1 Technical Reference: Kernel and Subsystems Volume 1*.

---

## Chapter 7. Communications I/O Subsystem

The Communication I/O Subsystem design introduces a more efficient, streamlined approach to attaching data link control (DLC) processes to communication and LAN adapters.

The Communication I/O Subsystem consists of one or more physical device handlers (PDHs) that control various communication adapters. The interface to the physical device handlers can support any number of processes, the limit being device-dependent.

**Note:** A PDH, as used for the Communications I/O, provides both the device head role for interfacing to users, and the device handler role for performing I/O to the device.

A communications PDH is a special type of multiplexed character device driver. Information common to all communications device handlers is discussed here. Additionally, individual communications PDHs have their own adapter-specific sets of information. Refer to the following to learn more about the adapter types:

- Serial Optical Link Device Handler Overview

Each adapter type requires a device driver. Each PDH can support one or more adapters of the same type.

There are two interfaces a user can use to access a PDH. One is from a user-mode process (application space), and the other is from a kernel-mode process (within the kernel).

---

### User-Mode Interface to a Communications PDH

The user-mode process uses system calls (**open**, **close**, **select**, **poll**, **ioctl**, **read**, **write**) to interface to the PDH to send or receive data. The **poll** or **select** subroutine notifies a user-mode process of available receive data, available transmit, and status and exception conditions.

---

### Kernel-Mode Interface to a Communications PDH

The kernel-mode interface to a communications PDH differs from the interface supported for a user-mode process in the following ways:

- Kernel services are used instead of system calls. This means that, for example, the **fp\_open** kernel service is used instead of the **open** subroutine. The same holds true for the **fp\_close**, **fp\_ioctl**, and **fp\_write** kernel services.
- The **ddread** entry point, **ddselect** entry point, and **CIO\_GET\_STAT** (Get Status) **ddioctl** operation are not supported in kernel mode. Instead, kernel-mode processes specify at open time the addresses of their own procedures for handling receive data available, transmit available and status or exception conditions. The PDH directly calls the appropriate procedure, whenever that condition arises. These kernel procedures must execute and return quickly since they are executing within the priority of the PDH.
- The **ddwrite** operation for a kernel-mode process differs from a user-mode process in that there are two ways to issue a **ddwrite** operation to transmit data:
  - Transmit each buffer of data with the **fp\_write** kernel service.
  - Use the fast write operation, which allows the user to directly call the **ddwrite** operation (no context switching) for each buffer of data to be transmitted. This operation helps increase the performance of transmitted data. A **fp\_ioctl** (**CIO\_GET\_FASTWRT**) kernel service call obtains the functional address of the write function. This address is used on all subsequent write function calls. Support of the fast write operation is optional for each device.

---

## CDLI Device Drivers

Some device drivers have a different design and use the services known as Common Data Link Interface (CDLI). The following device drivers use CDLI:

- Forum-Compliant ATM LAN Emulation Device Driver
- Fiber Distributed Data Interface (FDDI) Device Driver
- High-Performance (8fc8) Token-Ring Device Driver
- High-Performance (8fa2) Token-Ring Device Driver
- Ethernet Device Drivers

---

## Communications Physical Device Handler Model Overview

A physical device handler (PDH) must provide eight common entry points. An individual PDH names its entry points by placing a unique identifier in front of the supported command type. The following are the required eight communications PDH entry points:

<b>ddconfig</b>	Performs configuration functions for a device handler. Supported the same way that the common <b>ddconfig</b> entry point is.
<b>ddmpx</b>	Allocates or deallocates a channel for a multiplexed device handler. Supported the same way as the common <b>ddmpx</b> device handler entry point.
<b>ddopen</b>	Performs data structure allocation and initialization for a communications PDH. Supported the same way as the common <b>ddopen</b> entry point. Time-consuming tasks, such as port initialization and connection establishment, are deferred until the ( <b>CIO_START</b> ) <b>ddioctl</b> call is issued. A PDH can support multiple users of a single port.
<b>ddclose</b>	Frees up system resources used by the specified communications device until they are needed again. Supported the same way as the common <b>ddclose</b> entry point.
<b>ddwrite</b>	Queues a message for transmission or blocks until the message can be queued. The <b>ddwrite</b> entry point can attempt to queue a transmit request (nonblocking) or wait for it to be queued (blocking), depending on the setting of the <b>DNDELAY</b> flag. The caller has the additional option of requesting an asynchronous acknowledgment when the transmission actually completes.
<b>ddread</b>	Returns a message of data to a user-mode process. Supports blocking or nonblocking reads depending on the setting of the <b>DNDELAY</b> flag. A blocking read request does not return to the caller until data is available. A nonblocking read returns with a message of data if it is immediately available. Otherwise, it returns a length of 0 (zero).
<b>ddselect</b>	Checks to see if a specified event or events has occurred on the device for a user-mode process. Supported the same way as the common <b>ddselect</b> entry point.
<b>ddioctl</b>	Performs the special I/O operations requested in an <b>ioctl</b> subroutine. Supported the same way as the common <b>ddioctl</b> entry point. In addition, a communications PDH must support the following four options: <ul style="list-style-type: none"><li>• <b>CIO_START</b></li><li>• <b>CIO_HALT</b></li><li>• <b>CIO_QUERY</b></li><li>• <b>CIO_GET_STAT</b></li></ul>

Individual PDHs can add additional commands. Hardware initialization and other time-consuming activities, such as call establishment, are performed during the **CIO\_START** operation.

## Use of mbuf Structures in the Communications PDH

PDHs use **mbuf** structures to buffer send and receive data. These structures allow the PDH to gather data when transmitting frames and scatter for receive operations. The **mbuf** structures are internal to the kernel and are used only by kernel-mode processes and PDHs.

PDHs and kernel-mode processes require a set of utilities for obtaining and returning **mbuf** structures from a buffer pool.

Kernel-mode processes use the Berkeley **mbuf** scheme for transmit and receive buffers. The structure for an **mbuf** is defined in the `/usr/include/sys/mbuf.h` file.

## Common Communications Status and Exception Codes

In general, communication device handlers return codes from a group of common exception codes. However, device handlers for specific communication devices can return device-specific exception codes. Common exception codes are defined in the `/usr/include/sys/comio.h` file and include the following:

<b>CIO_OK</b>	Indicates that the operation was successful.
<b>CIO_BUF_OVFLW</b>	Indicates that the data was lost due to buffer overflow.
<b>CIO_HARD_FAIL</b>	Indicates that a hardware failure was detected.
<b>CIO_NOMBUF</b>	Indicates that the operation was unable to allocate <b>mbuf</b> structures.
<b>CIO_TIMEOUT</b>	Indicates that a time-out error occurred.
<b>CIO_TX_FULL</b>	Indicates that the transmit queue is full.
<b>CIO_NET_RCVRY_ENTER</b>	Enters network recovery.
<b>CIO_NET_RCVRY_EXIT</b>	Indicates the device handler is exiting network recovery.
<b>CIO_NET_RCVRY_MODE</b>	Indicates the device handler is in Recovery mode.
<b>CIO_INV_CMD</b>	Indicates that an invalid command was issued.
<b>CIO_BAD_MICROCODE</b>	Indicates that the microcode download failed.
<b>CIO_NOT_DIAG_MODE</b>	Indicates that the command could not be accepted because the adapter is not open in Diagnostic mode.
<b>CIO_BAD_RANGE</b>	Indicates that the parameter values have failed a range check.
<b>CIO_NOT_STARTED</b>	Indicates that the command could not be accepted because the device has not yet been started by the first call to <b>CIO_START</b> operation.
<b>CIO_LOST_DATA</b>	Indicates that the receive packet was lost.
<b>CIO_LOST_STATUS</b>	Indicates that a status block was lost.
<b>CIO_NETID_INV</b>	Indicates that the network ID was not valid.
<b>CIO_NETID_DUP</b>	Indicates that the network ID was a duplicate of an existing ID already in use on the network.
<b>CIO_NETID_FULL</b>	Indicates that the network ID table is full.

---

## Status Blocks for Communications Device Handlers Overview

Status blocks are used to communicate status and exception information.

User-mode processes receive a status block whenever they request a **CIO\_GET\_STAT** operation. A user-mode process can wait for the next available status block by issuing a **ddselect** entry point with the specified **POLLPRI** event.

A kernel-mode process receives a status block through the **stat\_fn** procedure. This procedure is specified when the device is opened with the **ddopen** entry point.

Status blocks contain a code field and possible options. The code field indicates the type of status block code (for example, **CIO\_START\_DONE**). A status block's options depend on the block code. The C structure of a status block is defined in the `/usr/include/sys/comio.h` file.

The following are the common status codes:

- **CIO\_START\_DONE**
- **CIO\_HALT\_DONE**
- **CIO\_TX\_DONE**
- **CIO\_NULL\_BLK**
- **CIO\_LOST\_STATUS**
- **CIO\_ASYNC\_STATUS**

Additional device-dependent status block codes may be defined.

## CIO\_START\_DONE

This block is provided by the device handler when the **CIO\_START** operation completes:

option[0]	The <b>CIO_OK</b> or <b>CIO_HARD_FAIL</b> status/exception code from the common or device-dependent list.
option[1]	The low-order two bytes are filled in with the <code>netid</code> field. This field is passed when the <b>CIO_START</b> operation is invoked.
option[2]	Device-dependent.
option[3]	Device-dependent.

## CIO\_HALT\_DONE

This block is provided by the device handler when the **CIO\_HALT** operation completes:

option[0]	The <b>CIO_OK</b> status/exception code from the common or device-dependent list.
option[1]	The low-order two bytes are filled in with the <code>netid</code> field. This field is passed when the <b>CIO_START</b> operation is invoked.
option[2]	Device-dependent.
option[3]	Device-dependent.

## CIO\_TX\_DONE

The following block is provided when the physical device handler (PDH) is finished with a transmit request for which acknowledgment was requested:

option[0]	The <b>CIO_OK</b> or <b>CIO_TIMEOUT</b> status/exception code from the common or device-dependent list.
option[1]	The <code>write_id</code> field specified in the <b>write_extension</b> structure passed in the <code>ext</code> parameter to the <b>ddwrite</b> entry point.
option[2]	For a kernel-mode process, indicates the <b>mbuf</b> pointer for the transmitted frame.
option[3]	Device-dependent.

## CIO\_NULL\_BLK

This block is returned whenever a status block is requested but there are none available:

option[0]	Not used
option[1]	Not used
option[2]	Not used
option[3]	Not used

## CIO\_LOST\_STATUS

This block is returned once after one or more status blocks is lost due to status queue overflow. The **CIO\_LOST\_STATUS** block provides the following:

option[0]	Not used
option[1]	Not used
option[2]	Not used
option[3]	Not used



## CIO\_ASYNC\_STATUS

This status block is used to return status and exception codes that occur unexpectedly:

option[0]	The <b>CIO_HARD_FAIL</b> or <b>CIO_LOST_DATA</b> status/exception code from the common or device-dependent list
option[1]	Device-dependent
option[2]	Device-dependent
option[3]	Device-dependent

---

## MPQP Device Handler Interface Overview for the ARTIC960Hx PCI Adapter

The ARTIC960Hx PCI Adapter (PCI MPQP) device handler is a component of the communication I/O subsystem. The PCI MPQP device handler interface is made up of the following eight entry points:

<b>tsclose</b>	Resets the PCI MPQP device to a known state and returns system resources back to the system on the last close for that adapter. The port no longer transmits or receives data.
<b>tsconfig</b>	Provides functions for initializing and terminating the PCI MPQP device handler and adapter.
<b>tsioctl</b>	Provides the following functions for controlling the PCI MPQP device: <b>CIO_START</b> Initiates a session with the PCI MPQP device handler. <b>CIO_HALT</b> Ends a session with the PCI MPQP device handler. <b>CIO_QUERY</b> Reads the counter values accumulated by the PCI MPQP device handler. <b>CIO_GET_STAT</b> Gets the status of the current PCI MPQP adapter and device handler. <b>MP_CHG_PARMS</b> Permits the data link control (DLC) to change certain profile parameters after the PCI MPQP device has been started.
<b>tsopen</b>	Opens a channel on the PCI MPQP device for transmitting and receiving data.
<b>tsmpx</b>	Provides allocation and deallocation of a channel.
<b>tsread</b>	Provides the means for receiving data to the PCI MPQP device.
<b>tsselect</b>	Provides the means for determining which specified events have occurred on the PCI MPQP device.
<b>tswrite</b>	Provides the means for transmitting data to the PCI MPQP device.

## Binary Synchronous Communication (BSC) with the PCI MPQP Adapter

The PCI MPQP adapter software performs low-level BSC frame-type determination to facilitate character parsing at the kernel-mode process level. Frames received without errors are parsed. A message type is returned in the status field of the extension block along with a pointer to the receive buffer. The message type indicates the type of frame that was received.

For control frames that only contain control characters, the message type is returned and no data is transferred from the board. For example, if an ACK0 was received, the message type MP\_ACK0 is returned in the status field of the extension block. In addition, a NULL pointer for the receive buffer is returned. If an error occurs, the error status is logged by the device driver. Unlogged buffer overrun errors are an exception.

**Note:** In BSC communications, the caller receives either a message type or an error status.

Read operations must be performed using the **readx** subroutine because the **read\_extension** structure is needed to return BSC function results.

### BSC Message Types Detected by the PCI MPQP Adapter

BSC message types are defined in the `/usr/include/sys/mpqp.h` file. The PCI MPQP adapter can detect the following message types:

MP_ACK0	MP_DISC	MP_STX_ETX
MP_ACK1	MP_SOH_ITB	MP_STX_ENQ
MP_WACK	MP_SOH_ETB	MP_DATA_ACK0
MP_NAK	MP_SOH_ETX	MP_DATA_ACK1
MP_ENQ	MP_SOH_ENQ	MP_DATA_NAK
MP_EOT	MP_STX_ITB	MP_DATA_ENQ
MP_RVI	MP_STX_ETB	

### Receive Errors Logged by the PCI MPQP Adapter

The PCI MPQP adapter detects many types of receive errors. As errors occur they are logged and the appropriate statistical counter is incremented. The kernel-mode process is not notified of the error. The following are the possible BSC receive errors logged by the PCI MPQP adapter:

- Receive overrun
- A cyclical redundancy check (CRC) or longitudinal redundancy check (LRC) framing error
- Parity error
- Clear to Send (CTS) timeout
- Data synchronization lost
- ID field greater than 15 bytes (BSC)
- Invalid pad at end of frame (BSC)
- Unexpected or invalid data (BSC)

If status and data information are available, but no extension block is provided, the **read** operation returns the data, but not the status information.

**Note:** Errors, such as buffer overflow errors, can occur during the read data operation. In these cases, the return value is the byte count. Therefore, status should be checked even if no **errno** global value is returned.

### Description of the PCI MPQP Card

The PCI MPQP card is a 4-port multiprotocol adapter that supports BSC and SDLC on the EIA232-D, X.21, and V.35 physical interfaces. When using the X.21 physical interface, X.21 centralized multipoint operation on a leased-circuit public data network is not supported.

---

## Serial Optical Link Device Handler Overview

The serial optical link (SOL) device handler is a component of the communication I/O subsystem. The device handler can support one to four serial optical ports. An optical port consists of two separate pieces. The serial link adapter is on the system planar and is packaged with two to four adapters in a single chip. The serial optical channel converter plugs into a slot on the system planar and provides two separate optical ports.

## Special Files

There are two separate interfaces to the serial optical link device handler. The special file **/dev/ops0** provides access to the optical port subsystem. An application that opens this special file has access to all the ports, but it does not need to be aware of the number of ports available. Each write operation includes a destination processor ID. The device handler sends the data out the correct port to reach that processor. In case of a link failure, the device handler uses any link that is available.

The **/dev/op0**, **/dev/op1**, ..., **/dev/opn** special files provide a diagnostic interface to the serial link adapters and the serial optical channel converters. Each special file corresponds to a single optical port that can only be opened in Diagnostic mode. A diagnostic open allows the diagnostic ioctls to be used, but normal reads and writes are not allowed. A port that is open in this manner cannot be opened with the **/dev/ops0** special file. In addition, if the port has already been opened with the **/dev/ops0** special file, attempting to open a **/dev/opx** special file will fail unless a forced diagnostic open is used.

## Entry Points

The SOL device handler interface consists of the following entry points:

<b>sol_close</b>	Resets the device to a known state and frees system resources.
<b>sol_config</b>	Provides functions to initialize and terminate the device handler, and query the vital product data (VPD).
<b>sol_fastwrt</b>	Provides the means for kernel-mode users to transmit data to the SOL device driver.
<b>sol_ioctl</b>	Provides various functions for controlling the device. The valid <b>sol_ioctl</b> operations are: <b>CIO_GET_FASTWRT</b> Gets attributes needed for the <b>sol_fastwrt</b> entry point. <b>CIO_GET_STAT</b> Gets the device status. <b>CIO_HALT</b> Halts the device. <b>CIO_QUERY</b> Queries device statistics. <b>CIO_START</b> Starts the device. <b>IOCINFO</b> Provides I/O character information. <b>SOL_CHECK_PRID</b> Checks whether a processor ID is connected. <b>SOL_GET_PRIDS</b> Gets connected processor IDs.
<b>sol_mpx</b>	Provides allocation and deallocation of a channel.
<b>sol_open</b>	Initializes the device handler and allocates the required system resources.
<b>sol_read</b>	Provides the means for receiving data.
<b>sol_select</b>	Determines if a specified event has occurred on the device.
<b>sol_write</b>	Provides the means for transmitting data.

---

## Configuring the Serial Optical Link Device Driver

When configuring the serial optical link (SOL) device driver, consider the physical and logical devices, and changeable attributes of the SOL subsystem.

## Physical and Logical Devices

The SOL subsystem consists of several physical and logical devices in the ODM configuration database:

Device	Description
<b>slc</b> (serial link chip)	There are two serial link adapters in each COMBO chip. The <b>slc</b> device is automatically detected and configured by the system.
<b>otp</b> (optic two-port card)	Also known as the serial optical channel converter (SOCC). There is one SOCC possible for each <b>slc</b> . The <b>otp</b> device is automatically detected and configured by the system.
<b>op</b> (optic port)	There are two optic ports per <b>otp</b> . The <b>op</b> device is automatically detected and configured by the system.
<b>ops</b> (optic port subsystem)	This is a logical device. There is only one created at any time. The <b>ops</b> device requires some additional configuration initially, and is then automatically configured from that point on. The <b>/dev/ops0</b> special file is created when the <b>ops</b> device is configured. The <b>ops</b> device cannot be configured when the processor ID is set to -1.

## Changeable Attributes of the Serial Optical Link Subsystem

The system administrator can change the following attributes of the serial optical link subsystem:

**Note:** If your system uses serial optical link to make a direct, point-to-point connection to another system or systems, special conditions apply. You must start interfaces on two systems at approximately the same time, or a method error occurs. If you wish to connect to at least one machine on which the interface has already been started, this is not necessary.

<b>Processor ID</b>	This is the address by which other machines connected by means of the optical link address this machine. The processor ID can be any value in the range of 1 to 254. To avoid a conflict on the network, this value is initially set to -1, which is not valid, and the <b>ops</b> device cannot be configured. <b>Note:</b> If you are using TCP/IP over the serial optical link, the processor ID must be the same as the low-order octet of the IP address. It is not possible to successfully configure TCP/IP if the processor ID does not match.
<b>Receive Queue Size</b>	This is the maximum number of packets that is queued for a user-mode caller. The default value is 30 packets. Any integer in the range from 30 to 150 is valid.
<b>Status Queue Size</b>	This is the maximum number of status blocks that will be queued for a user-mode caller. The default value is 10. Any integer in the range from 3 to 20 is valid.

The standard SMIT interface is available for setting these attributes, listing the serial optical channel converters, handling the initial configuration of the **ops** device, generating a trace report, generating an error report, and configuring TCP/IP.

---

## Forum-Compliant ATM LAN Emulation Device Driver

The **Forum-Compliant ATM LAN Emulation (LANE)** device driver allows communications applications and access methods that would normally operate over local area network (LAN) attachments to operate over high-speed ATM networks. This **ATM LANE** function supports LAN Emulation Client (LEC) as specified in *The ATM Forum Technical Committee LAN Emulation Over ATM Version 1.0*, as well as MPOA Client (MPC) through a subset of *ATM Forum LAN Emulation Over ATM Version 2 - LUNI Specification*, and *ATM Forum Multi-Protocol Over ATM Version 1.0*.

The **ATM LANE** device driver emulates the operation of Standard Ethernet, IEEE 802.3 Ethernet, and IEEE 802.5 Token Ring LANs. It encapsulates each LAN packet and transfers its LAN data over an ATM network at up to OC12 speeds (622 megabits per second). This data can also be bridged transparently to a traditional LAN with ATM/LAN bridges such as the IBM 2216.

Each LEC participates in an emulated LAN containing additional functions such as:

- A LAN Emulation Configuration Server (LECS) that provides automated configuration of the LEC's operational attributes.
- A LAN Emulation Server (LES) that provides address resolution
- A Broadcast and Unknown Server (BUS) that distributes packets sent to a broadcast address or packets sent without knowing the ATM address of the remote station (for example, whenever an ARP response has not been received yet).

There is always at least one ATM switch and a possibility of additional switches, bridges, or concentrators.

ATM supports UNI3.0, UNI3.1, and UNI4.0 signalling.

In support of Ethernet jumbo frames, LE Clients can be configured with maximum frame size values greater than 1516 bytes. Supported forum values are: 1516, 4544, 9234, and 18190.

Incoming Add Party requests are supported for the Control Distribute and Multicast Forward Virtual Circuits (VCs). This allows multiple LE clients to be open concurrently on the same ELAN without additional hardware.

LANE and MPOA are both enabled for IPV4 TCP checksum offload. Transmit offload is automatically enabled when it is supported by the adapter. Receive offload is configured by using the `rx_checksum` attribute. The `NDD_CHECKSUM_OFFLOAD` device driver flag is set to indicate general offload capability and also indicates that transmit offload is operational.

Transmit offload of IP-fragmented TCP packets is not supported. Transmit packets that MPOA needs to fragment are offloaded in the MPOA software, instead of in the adapter. UDP offloading is also not supported.

The **ATM LANE** device driver is a dynamically loadable device driver. Each LE Client or MPOA Client is configurable by the operator, and the LANE driver is loaded into the system as part of that configuration process. If an LE Client or MPOA Client has already been configured, the LANE driver is automatically reloaded at reboot time as part of the system configuration process.

The interface to the **ATM LANE** device driver is through kernel services known as Network Services.

Interfacing to the **ATM LANE** device driver is achieved by calling the device driver's entry points for opening the device, closing the device, transmitting data, and issuing device control commands, just as you would interface to any of the Common Data Link Interface (CDLI) LAN device drivers.

The **ATM LANE** device driver interfaces with all hardware-level ATM device drivers that support CDLI, ATM Call Management, and ATM Signaling.

## Adding ATM LANE Clients

At least one ATM LAN Emulation client must be added to the system to communicate over an ATM network using the **ATM Forum LANE** protocol. A user with root authority can add Ethernet or Token-Ring clients using the `smit atmle_panel` fast path.

Entries are required for the Local LE Client's **LAN MAC Address** field and possibly the **LES ATM Address** or **LECS ATM Address** fields, depending on the support provided at the server. If the server accepts the well-known ATM address for LECS, the value of the **Automatic Configuration via LECS** field can be set to **Yes**, and the **LES** and **LECS ATM Address** fields can be left blank. If the server does not support the well-known ATM address for LECS, an ATM address must be entered for either LES (manual configuration) or LECS (automatic configuration). All other configuration attribute values are optional. If used, you can accept the defaults for ease-of-use.

Configuration help text is also available within the SMIT LE Client add and change menus.

## Configuration Parameters for the ATM LANE Device Driver

The **ATM LANE** device driver supports the following configuration parameters for each LE Client:

<b>addl_drvr</b>	Specifies the CDLI demultiplexer being used by the LE Client. The value set by the <b>ATM LANE</b> device driver is <b>/usr/lib/methods/cfgdmxtok</b> for Token Ring emulation and <b>/usr/lib/methods/cfgdmxeth</b> for Ethernet. This is not an operator-configurable attribute.
<b>addl_stat</b>	Specifies the routine being used by the LE client to generate device-specific statistics for the <b>entstat</b> and <b>tokstat</b> commands. The values set by the <b>ATM LANE</b> device driver are: <ul style="list-style-type: none"><li>• <b>/usr/sbin/atmle_ent_stat</b></li><li>• <b>/usr/sbin/atmle_tok_stat</b></li></ul> The <b>addl_stat</b> attribute is not operator-configurable.
<b>arp_aging_time</b>	Specifies the maximum timeout period (in seconds) that the LE Client maintains an LE_ARP cache entry without verification (ATM Forum LE Client parameter <i>C17</i> ). The default value is 300 seconds.
<b>arp_cache_size</b>	Specifies the maximum number of LE_ARP cache entries that are held by the LE Client before removing the least recently used entry. The default value is 128 entries.
<b>arp_response_timeout</b>	Specifies the maximum timeout period (in seconds) for LE_ARP request/response exchanges (ATM Forum LE Client parameter <i>C20</i> ). The default value is 1 second.
<b>atm_device</b>	Specifies the logical name of the physical <b>ATM</b> device driver that this LE Client is to operate with, as specified in the CuDv database (for example, <b>atm0</b> , <b>atm1</b> , <b>atm2</b> , ...). The default is <b>atm0</b> .
<b>auto_cfg</b>	Specifies whether the LE Client is to be automatically configured. Select <b>Yes</b> if the LAN Emulation Configuration Server (LECS) is to be used by the LE Client to obtain the ATM address of the LE ARP Server, as well as any additional configuration parameters provided by the LECS. The default value is No (manual configuration). The attribute values are: <b>Yes</b> auto configuration <b>No</b> manual configuration  <b>Note:</b> Configuration parameters provided by LECS override configuration values provided by the operator.
<b>debug_trace</b>	Specifies whether this LE Client should keep a real time debug log within the kernel and allow full system trace capability. Select <b>Yes</b> to enable full tracing capability for this LE Client. Select <b>No</b> for optimal performance when minimal tracing is desired. The default is <b>Yes</b> (full tracing capability).

<b>elan_name</b>	<p>Specifies the name of the Emulated LAN this LE Client wishes to join (ATM Forum LE Client parameter <i>C5</i>). This is an SNMPv2 DisplayString of 1-32 characters, or it might be left blank (unused). See RFC1213 for a definition of an SNMPv2 DisplayString.</p> <p><b>Note:</b></p> <ol style="list-style-type: none"> <li>Any operator configured <b>elan_name</b> should match exactly what is expected at the LECS/LES server when attempting to join an ELAN. Some servers can alias the ELAN name and allow the operator to specify a logical name that correlates to the actual name. Other servers might require the exact name to be specified.</li> </ol> <p>Previous versions of LANE would accept any <b>elan_name</b> from the server, even when configured differently by the operator. However, with multiple LECS/LES now possible, it is desirable that only the ELAN identified by the network administrator is joined. Use the <b>force_elan_name</b> attribute below to ensure that the name you have specified is the only ELAN joined.</p> <p>If no <b>elan_name</b> attribute is configured at the LEC, or the <b>force_elan_name</b> attribute is disabled, the server can stipulate whatever <b>elan_name</b> is available.</p> <p>Failure to use an ELAN name that is identical to the server's when specifying the <b>elan_name</b> and <b>force_elan_name</b> attributes causes the LEC to fail the join process, with <b>entstat/tokstat</b> status indicating Driver Flag <b>Limbo</b>.</p> <ol style="list-style-type: none"> <li>Blanks can be inserted within an <b>elan_name</b> by typing a tilde (~) character whenever a blank character is desired. This allows a network administrator to specify an ELAN name with imbedded blanks as in the default of some servers.</li> </ol> <p>Any tilde (~) character that occupies the first character position of the <b>elan_name</b> remains unchanged (that is, the resulting name can start with a tilde (~) but all remaining tilde characters are converted to blanks).</p>
<b>failsafe_time</b>	<p>Specifies the maximum timeout period (in seconds) that the LE Client attempts to recover from a network outage. A value of zero indicates that you should continue recovery attempts unless a nonrecoverable error is encountered. The default value is 0 (unlimited).</p>
<b>flush_timeout</b>	<p>Specifies the maximum timeout period (in seconds) for FLUSH request/response exchanges (ATM Forum LE Client parameter <i>C21</i>). The default value is 4 seconds.</p>
<b>force_elan_name</b>	<p>Specifies that the Emulated LAN Name returned from the LECS or LES servers must exactly match the name entered in the <b>elan_name</b> attribute above. Select <b>Yes</b> if the <b>elan_name</b> field must match the server configuration and join parameters. This allows a specific ELAN to be joined when multiple LECS and LES servers are available on the network. The default value is No, which allows the server to specify the ELAN Name.</p>
<b>fwd_delay_time</b>	<p>Specifies the maximum timeout period (in seconds) that the LE Client maintains an entry for a non-local MAC address in its LE_ARP cache without verification, when the <b>Topology Change</b> flag is True (ATM Forum LE Client parameter <i>C18</i>). The default value is 15 seconds.</p>
<b>fwd_dsc_timeout</b>	<p>Specifies the timeout period (in seconds) that can elapse without an active Multicast Forward VCC from the BUS. (ATM Forum LE Client parameter <i>C33</i>). If the timer expires without an active Multicast Forward VCC, the LE Client attempts recovery by re-establishing its Multicast Send VCC to the BUS. The default value is 60 seconds.</p>
<b>init_ctl_time</b>	<p>Specifies the initial control timeout period (in seconds) for most request/response control frame interactions (ATM Forum LE Client parameter <i>C7i</i>). This timeout is increased by its initial value after each timeout expiration without a response, but does not exceed the value specified by the Maximum Control Timeout attribute (<b>max_ctl_time</b>). The default value is 5 seconds.</p>
<b>lan_type</b>	<p>Identifies the type of local area network being emulated (ATM Forum LE Client parameter <i>C2</i>). Both Ethernet/IEEE 802.3 and Token Ring LANs can be emulated using ATM Forum LANE. The attribute values are:</p> <ul style="list-style-type: none"> <li>• Ethernet/IEEE802.3</li> <li>• TokenRing</li> </ul>

<b>lecs_atm_addr</b>	<p>If you are doing auto configuration using the <b>LE Configuration Server (LECS)</b>, this field specifies the ATM address of LECS. It can remain blank if the address of LECS is not known and the LECS is connected by way of PVC (VPI=0, VCI=17) or the well-known address, or is registered by way of ILMI. If the 20-byte address of the LECS is known, it must be entered as hexadecimal numbers using a period (.) as the delimiter between bytes. Leading zeros of each byte can be omitted, for example:</p> <p>47.0.79.0.0.0.0.0.0.0.0.0.0.0.a0.3.0.0.1</p> <p>(the LECS well-known address)</p>
<b>les_atm_addr</b>	<p>If you are doing manual configuration (without the aid of an <b>LECS</b>), this field specifies the ATM address of the LE ARP Server (LES) (ATM Forum LE Client parameter <i>C9</i>). This 20-byte address must be entered as hexadecimal numbers using a period (.) as the delimiter between bytes. Leading zeros of each byte can be omitted, for example:</p> <p>39.11.ff.22.99.99.99.0.0.0.1.49.10.0.5a.68.0.a.1</p>
<b>local_lan_addr</b>	<p>Specifies the local unicast LAN MAC address that is represented by this LE Client and registered with the LE Server (ATM Forum LE Client parameter <i>C6</i>). This 6-byte address must be entered as hexadecimal numbers using a period (.) as the delimiter between bytes. Leading zeros of each byte can be omitted.</p> <p>Ethernet Example: 2.60.8C.2C.D2.DC Token Ring Example: 10.0.5A.4F.4B.C4</p>
<b>max_arp_retries</b>	<p>Specifies the maximum number of times an LE_ARP request can be retried (ATM Forum LE Client parameter <i>C13</i>). The default value is 1.</p>
<b>max_config_retries</b>	<p>Specifies the number of times a configuration control frame such as LE_JOIN_REQUEST should be retried. Duration (in seconds) between retries is derived from the <b>init_ctl_time</b> and <b>max_ctl_time</b> attributes. The default is 1.</p>
<b>max_ctl_time</b>	<p>Specifies the maximum timeout period (in seconds) for most request and response control frame interactions (ATM Forum LE Client parameter <i>C7</i>). The default value is 30 seconds.</p>
<b>max_frame_size</b>	<p>Specifies the maximum AAL-5 send data-unit size of data frames for this LE Client. In general, this value should coincide with the LAN type and speed as follows:</p> <p><b>Unspecified</b> for auto LECS configuration</p> <p><b>1516 bytes</b> for Ethernet and IEEE 802.3 networks</p> <p><b>4544 bytes</b> for 4 Mbps Token Rings or Ethernet jumbo frames</p> <p><b>9234 bytes</b> for 16 Mbps Token Rings or Ethernet jumbo frames</p> <p><b>18190 bytes</b> for 16 Mbps Token Rings or Ethernet jumbo frames</p>
<b>max_queued_frames</b>	<p>Specifies the maximum number of outbound packets that are held for transmission per LE_ARP cache entry. This queueing occurs when the Maximum Unknown Frame Count (<b>max_unknown_fct</b>) has been reached, or when flushing previously transmitted packets while switching to a new virtual channel. The default value is 60 packets.</p>
<b>max_rdy_retries</b>	<p>Specifies the maximum number of READY_QUERY packets sent in response to an incoming call that has not yet received data or a READY_IND packet. The default value is 2 retries.</p>
<b>max_unknown_fct</b>	<p>Specifies the maximum number of frames for a given unicast LAN MAC address that can be sent to the Broadcast and Unknown Server (BUS) within time period Maximum Unknown Frame Time (<b>max_unknown_ftm</b>) (ATM Forum LE Client parameter <i>C10</i>). The default value is 1.</p>



<b>max_unknown_ftm</b>	Specifies the maximum timeout period (in seconds) that a given unicast LAN address can be sent to the Broadcast and Unknown Server (BUS). The LE Client sends no more than Maximum Unknown Frame Count ( <b>max_unknown_fct</b> ) packets to a given unicast LAN destination within this timeout period (ATM Forum LE Client parameter <i>C11</i> ). The default value is 1 second.
<b>mpoa_enabled</b>	Specifies whether Forum MPOA and LANE-2 functions should be enabled for this LE Client. Select <b>Yes</b> if MPOA will be operational on the LE Client. Select <b>No</b> when traditional LANE-1 functionality is required. The default is No (LANE-1).
<b>mpoa_primary</b>	Specifies whether this LE Client is to be the primary configurator for MPOA using LAN Emulation Configuration Server (LECS). Select <b>Yes</b> if this LE Client will be obtaining configuration information from the LECS for the MPOA Client. This attribute is only meaningful if running auto config with an LECS, and indicates that the MPOA configuration TLVs from this LEC is available to the MPC. Only one LE Client can be active as the MPOA primary configurator. The default is No.
<b>path_sw_delay</b>	Specifies the maximum timeout period (in seconds) that frames sent on any path in the network take to be delivered (ATM Forum LE Client parameter <i>C22</i> ). The default value is 6 seconds.
<b>peak_rate</b>	Specifies the forward and backward peak bit rate in K-bits per second that are used by this LE Client to set up virtual channels. Specify a value that is compatible with the lowest speed remote device with which you expect this LE Client to be communicating. Higher values might cause congestion in the network. A value of zero allows the LE Client to adjust its <b>peak_rate</b> to the actual speed of the adapter. If the adapter does not provide its maximum peak rate value, the LE Client defaults its <b>peak_rate</b> to 25600. Any non-zero value specified is accepted and used by the LE Client up to the maximum value allowed by the adapter. The default value is 0, which uses the adapter's maximum peak rate.
<b>ready_timeout</b>	Specifies the maximum timeout period (in seconds) in which data or a <b>READY_IND</b> message is expected from a calling party (ATM Forum LE Client parameter <i>C28</i> ). The default value is 4 seconds.
<b>ring_speed</b>	Specifies the Token Ring speed as viewed by the ifnet layer. The value set by the <b>ATM LANE</b> device driver is 16 Mbps for Token Ring emulation and ignored for Ethernet. This is not an operator-configurable attribute.
<b>rx_checksum</b>	Specifies whether this LE Client should offload TCP receive checksums to the ATM hardware. Select <b>Yes</b> if TCP checksums should be handled in hardware. Select <b>No</b> if TCP checksums should be handled in software. The default is <b>Yes</b> (enable hardware receive checksum).
<b>soft_restart</b>	<b>Note:</b> The ATM adapter must also have receive checksum enabled to be functional. Specifies whether active data virtual circuits (VCs) are to be maintained during connection loss of ELAN services such as the LE ARP Server (LES) or Broadcast and Unknown Server (BUS). Normal ATM Forum operation forces a disconnect of data VCs when LES/BUS connections are lost. This option to maintain active data VCs might be advantageous when server backup capabilities are available. The default value is No.
<b>vcc_activity_timeout</b>	Specifies the maximum timeout period (in seconds) for inactive Data Direct Virtual Channel Connections (VCCs). Any switched Data Direct VCC that does not transmit or receive data frames in this timeout period is terminated (ATM Forum LE Client parameter <i>C12</i> ). The default value is 1200 seconds (20 minutes).

## Device Driver Configuration and Unconfiguration

The **atmle\_config** entry point performs configuration functions for the **ATM LANE** device driver.

### Device Driver Open

The **atmle\_open** function is called to open the specified network device.

The **LANE** device driver does an asynchronous open. It starts the process of attaching the device to the network, sets the **NDD\_UP** flag in the **ndd\_flags** field, and returns 0. The network attachment continues in the background where it is driven by network activity and system timers.

**Note:** The Network Services **ns\_alloc** routine that calls this open routine causes the open to be synchronous. It waits until the **NDD\_RUNNING** or the **NDD\_LIMBO** flag is set in the **ndd\_flags** field or 15 seconds have passed.

If the connection is successful, the **NDD\_RUNNING** flag is set in the **ndd\_flags** field, and an **NDD\_CONNECTED** status block is sent. The **ns\_alloc** routine returns at this time.

If the device connection fails, the **NDD\_LIMBO** flag is set in the **ndd\_flags** field, and an **NDD\_LIMBO\_ENTRY** status block is sent.

If the device is eventually connected, the **NDD\_LIMBO** flag is disabled, and the **NDD\_RUNNING** flag is set in the **ndd\_flags** field. Both **NDD\_CONNECTED** and **NDD\_LIMBO\_EXIT** status blocks are sent.

## Device Driver Close

The **atmle\_close** function is called by the Network Services **ns\_free** routine to close the specified network device. This function resets the device to a known state and frees system resources associated with the device.

The device will not be detached from the network until the device's transmit queue is allowed to drain.

## Data Transmission

The **atmle\_output** function transmits data using the network device.

If the destination address in the packet is a broadcast address, the **M\_BCAST** flag in the **p\_mbuf->m\_flags** field should be set prior to entering this routine. A broadcast address is defined as FF.FF.FF.FF.FF.FF (hex) for both Ethernet and Token Ring and C0.00.FF.FF.FF.FF (hex) for Token Ring.

If the destination address in the packet is a multicast or group address, the **M\_MCAST** flag in the **p\_mbuf->m\_flags** field should be set prior to entering this routine. A multicast or group address is defined as any nonindividual address other than a broadcast address.

The device driver keeps statistics based on the **M\_BCAST** and **M\_MCAST** flags.

Token Ring LANE emulates a duplex device. If a Token Ring packet is transmitted with a destination address that matches the LAN MAC address of the local LE Client, the packet is received. This is also True for Token Ring packets transmitted to a broadcast address, enabled functional address, or an enabled group address. Ethernet LANE, on the other hand, emulates a simplex device and does not receive its own broadcast or multicast transmit packets.

## Data Reception

When the **LANE** device driver receives a valid packet from a network ATM device driver, the **LANE** device driver calls the **nd\_receive** function that is specified in the **ndd\_t** structure of the network device. The **nd\_receive** function is part of a CDLI network demuxer. The packet is passed to the **nd\_receive** function in mbufs.

The **LANE** device driver passes one packet to the **nd\_receive** function at a time.

The device driver sets the **M\_BCAST** flag in the **p\_mbuf->m\_flags** field when a packet is received that has an all-stations broadcast destination address. This address value is defined as FF.FF.FF.FF.FF.FF (hex) for both Token Ring and Ethernet and is defined as C0.00.FF.FF.FF.FF (hex) for Token Ring.

The device driver sets the **M\_MCAST** flag in the **p\_mbuf->m\_flags** field when a packet is received that has a nonindividual address that is different than an all-stations broadcast address.

Any packets received from the network are discarded if they do not fit the currently emulated **LAN** protocol and frame format are discarded.

## Asynchronous Status

When a status event occurs on the device, the **LANE** device driver builds the appropriate status block and calls the **nd\_status** function that is specified in the **ndd\_t** structure of the network device. The **nd\_status** function is part of a CDLI network demuxer.

The following status blocks are defined for the **LANE** device driver:

### Hard Failure

When an error occurs within the internal operation of the **ATM LANE** device driver, it is considered unrecoverable. If the device was operational at the time of the error, the **NDD\_LIMBO** and **NDD\_RUNNING** flags are disabled, and the **NDD\_DEAD** flag is set in the **ndd\_flags** field, and a hard failure status block is generated.

<b>code</b>	Set to <b>NDD_HARD_FAIL</b>
<b>option[0]</b>	Set to <b>NDD_UCODE_FAIL</b>

### Enter Network Recovery Mode

When the device driver detects an error that requires initiating recovery logic to make the device temporarily unavailable, the following status block is returned by the device driver:

<b>code</b>	Set to <b>NDD_LIMBO_ENTER</b>
<b>option[0]</b>	Set to <b>NDD_UCODE_FAIL</b>

**Note:** While the device driver is in this recovery logic, the network connections might not be fully functional. The device driver notifies users when the device is fully functional by way of an **NDD\_LIMBO\_EXIT** asynchronous status block.

When a general error occurs during operation of the device, this status block is generated.

### Exit Network Recovery Mode

When the device driver has successfully completed recovery logic from the error that made the device temporarily unavailable, the following status block is returned by the device driver. This status block means the device is now fully functional.

<b>code</b>	Set to <b>NDD_LIMBO_EXIT</b>
<b>option[0]</b>	The <b>option</b> field is not used.

## Device Control Operations

The **atmle\_ctl** function is used to provide device control functions.

### ATMLE\_MIB\_GET

This control requests the **LANE** device driver's current ATM LAN Emulation MIB statistics.

The user should pass in the address of an **atmle\_mibs\_t** structure as defined in **usr/include/sys/atmle\_mibs.h**. The driver returns **EINVAL** if the buffer area is smaller than the required structure.

The **ndd\_flags** field can be checked to determine the current state of the **LANE** device.

### ATMLE\_MIB\_QUERY

This control requests the **LANE** device driver's ATM LAN Emulation MIB support structure.

The user should pass in the address of an **atmle\_mibs\_t** structure as defined in **usr/include/sys/atmle\_mibs.h**. The driver returns EINVAL if the buffer area is smaller than the required structure.

The device driver does *not* support any variables for read\_write or write only. If the syntax of a member of the structure is some integer type, the level of support flag is stored in the whole field, regardless of the size of the field. For those fields defined as character arrays, the value is returned only in the first byte in the field.

## **NDD\_CLEAR\_STATS**

This control requests all the statistics counters kept by the **LANE** device driver to be zeroed.

## **NDD\_DISABLE\_ADDRESS**

This command disables the receipt of packets destined for a multicast/group address; and for Token Ring, it disables the receipt of packets destined for a functional address. For Token Ring, the functional address indicator (bit 0, the most significant bit of byte 2) indicates whether the address is a functional address (the bit is a 0) or a group address (the bit is a 1).

In all cases, the **length** field value is required to be 6. Any other value causes the **LANE** device driver to return EINVAL.

**Functional Address:** The reference counts are decremented for those bits in the functional address that are enabled (set to 1). If the reference count for a bit goes to zero, the bit is disabled in the functional address mask for this LE Client.

If no functional addresses are active after receipt of this command, the **TOK\_RECEIVE\_FUNC** flag in the **ndd\_flags** field is reset. If no functional or multicast/group addresses are active after receipt of this command, the **NDD\_ALTADDRS** flag in the **ndd\_flags** field is reset.

**Multicast/Group Address:** If a multicast/group address that is currently enabled is specified, receipt of packets destined for that group address is disabled. If an address is specified that is not currently enabled, EINVAL is returned.

If no functional or multicast/group addresses are active after receipt of this command, the **NDD\_ALTADDRS** flag in the **ndd\_flags** field is reset. Additionally for Token Ring, if no multicast/group address is active after receipt of this command, the **TOK\_RECEIVE\_GROUP** flag in the **ndd\_flags** field is reset.

## **NDD\_DISABLE\_MULTICAST**

The **NDD\_DISABLE\_MULTICAST** command disables the receipt of *all* packets with unregistered multicast addresses, and only receives those packets whose multicast addresses were registered using the **NDD\_ENABLE\_ADDRESS** command. The *arg* and *length* parameters are not used. The **NDD\_MULTICAST** flag in the **ndd\_flags** field is reset only after the reference count for multicast addresses has reached zero.

## **NDD\_ENABLE\_ADDRESS**

The **NDD\_ENABLE\_ADDRESS** command enables the receipt of packets destined for a multicast/group address; and additionally for Token Ring, it enables the receipt of packets destined for a functional address. For Ethernet, the address is entered in canonical format, which is left-to-right byte order with the I/G (Individual/Group) indicator as the least significant bit of the first byte. For Token Ring, the address format is entered in noncanonical format, which is left-to-right bit and byte order and has a functional address indicator. The functional address indicator (the most significant bit of byte 2) indicates whether the address is a functional address (the bit value is 0) or a group address (the bit value is 1).

In all cases, the **length** field value is required to be 6. Any other length value causes the **LANE** device driver to return EINVAL.

**Functional Address:** The Token-Ring network architecture provides bit-specific functional addresses for widely used functions, such as Ring Parameter Server or Configuration Report Server. Ring stations use functional address masks to identify these functions. The specified address is OR'ED with the currently specified functional addresses, and the resultant address is set as the functional address for the device. Functional addresses are encoded in a bit-significant format, thereby allowing multiple individual groups to be designated by a single address.

For example, if function G is assigned a functional address of C0.00.00.08.00.00 (hex), and function M is assigned a functional address of C0.00.00.00.00.40 (hex), then ring station Y, whose node contains function G and M, would have a mask of C0.00.00.08.00.40 (hex). Ring station Y would receive packets addressed to either function G or M or to an address like C0.00.00.08.00.48 (hex) because that address contains bits specified in the mask.

**Note:** The **LANE** device driver forces the first 2 bytes of the functional address to be C0.00 (hex). In addition, bits 6 and 7 of byte 5 of the functional address are forced to 0.

The **NDD\_ALTADDRS** and **TOK\_RECEIVE\_FUNC** flags in the **ndd\_flags** field are set.

Because functional addresses are encoded in a bit-significant format, reference counts are kept on each of the 31 least significant bits of the address. Reference counts are not kept on the 17 most significant bits (the C0.00 (hex) of the functional address and the functional address indicator bit).

**Multicast/Group Address:** A multicast/group address table is used by the **LANE** device driver to store address filters for incoming multicast/group packets. If the **LANE** device driver is unable to allocate kernel memory when attempting to add a multicast/group address to the table, the address is not added and ENOMEM is returned.

If the **LANE** device driver is successful in adding a multicast/group address, the **NDD\_ALTADDRS** flag in the **ndd\_flags** field is set. Additionally for Token Ring, the **TOK\_RECEIVE\_GROUP** flag is set, and the first 2 bytes of the group address are forced to be C0.00 (hex).

## **NDD\_ENABLE\_MULTICAST**

The **NDD\_ENABLE\_MULTICAST** command enables the receipt of packets with any multicast (or group) address. The *arg* and *length* parameters are not used. The **NDD\_MULTICAST** flag in the **ndd\_flags** field is set.

## **NDD\_DEBUG\_TRACE**

This control requests a LANE or MPOA driver to toggle the current state of its **debug\_trace** configuration flag.

This control is available to the operator through the LANE Ethernet **entstat -t** or LANE Token Ring **tokstat -t** commands, or through the MPOA **mpcstat -t** command. The current state of the **debug\_trace** configuration flag is displayed in the output of each command as follows:

- For the **entstat** and **tokstat** commands, **NDD\_DEBUG\_TRACE** is enabled only if you see Driver Flags: Debug.
- For the **mpcstat** command, you see Debug Trace: Enabled.

## **NDD\_GET\_ALL\_STATS**

This control requests all current LANE statistics, based on both the generic LAN statistics and the **ATM LANE** protocol in progress.

For Ethernet, pass in the address of an **ent\_ndd\_stats\_t** structure as defined in the file **/usr/include/sys/cdli\_entuser.h**.

For Token Ring, pass in the address of a **tok\_ndd\_stats\_t** structure as defined in the file **/usr/include/sys/cdli\_tokuser.h**.

The driver returns EINVAL if the buffer area is smaller than the required structure.

The **ndd\_flags** field can be checked to determine the current state of the LANE device.

### **NDD\_GET\_STATS**

This control requests the current generic LAN statistics based on the **LAN** protocol being emulated.

For Ethernet, pass in the address of an **ent\_ndd\_stats\_t** structure as defined in the file **/usr/include/sys/cdli\_entuser.h**.

For Token Ring, pass in the address of a **tok\_ndd\_stats\_t** structure as defined in file **/usr/include/sys/cdli\_tokuser.h**.

The **ndd\_flags** field can be checked to determine the current state of the LANE device.

### **NDD\_MIB\_ADDR**

This control requests the current receive addresses that are enabled on the **LANE** device driver. The following address types are returned, up to the amount of memory specified to accept the address list:

- Local LAN MAC Address
- Broadcast Address FF.FF.FF.FF.FF.FF (hex)
- Broadcast Address C0.00.FF.FF.FF.FF (hex)
- (returned for Token Ring only)
- Functional Address Mask
- (returned for Token Ring only, and only if at least one functional address has been enabled)
- Multicast/Group Address 1 through n
- (returned only if at least one multicast/group address has been enabled)

Each address is 6-bytes in length.

### **NDD\_MIB\_GET**

This control requests the current MIB statistics based on whether the LAN being emulated is Ethernet or Token Ring.

If Ethernet, pass in the address of an **ethernet\_all\_mib\_t** structure as defined in the file **/usr/include/sys/ethernet\_mibs.h**.

If Token Ring, pass in the address of a **token\_ring\_all\_mib\_t** structure as defined in the file **/usr/include/sys/tokenring\_mibs.h**.

The driver returns EINVAL if the buffer area is smaller than the required structure.

The **ndd\_flags** field can be checked to determine the current state of the LANE device.

### **NDD\_MIB\_QUERY**

This control requests **LANE** device driver's MIB support structure based on whether the LAN being emulated is Ethernet or Token Ring.

If Ethernet, pass in the address of an **ethernet\_all\_mib\_t** structure as defined in the file **/usr/include/sys/ethernet\_mibs.h**.

If Token Ring, pass in the address of a **token\_ring\_all\_mib\_t** structure as defined in the file **/usr/include/sys/tokenring\_mibs.h**.

The driver returns EINVAL if the buffer area is smaller than the required structure.

The device driver does *not* support any variables for read\_write or write only. If the syntax of a member of the structure is some integer type, the level of support flag is stored in the whole field, regardless of the size of the field. For those fields which are defined as character arrays, the value is returned only in the first byte in the field.

## Tracing and Error Logging in the ATM LANE Device Driver

The **LANE** device driver has two trace points:

- 3A1 - Normal Code Paths
- 3A2 - Error Conditions

Tracing can be enabled through SMIT or with the **trace** command.

```
trace -a -j 3a1,3a2
```

Tracing can be disabled through SMIT or with the **trcstop** command. After trace is stopped, the results can be formatted into readable text with the **trcrpt** command.

```
trcrpt > /tmp/trc.out
```

### LANE error log templates:

**ERRID\_ATMLE\_MEM\_ERR**

An error occurred while attempting to allocate memory or pin the code. This error log entry accompanies return code ENOMEM on an open or control operation.

**ERRID\_ATMLE\_LOST\_SW**

The **LANE** device driver lost contact with the ATM switch. The device driver enters Network Recovery Mode in an attempt to recover from the error and is temporarily unavailable during the recovery procedure. This generally occurs when the cable is unplugged from the switch or ATM adapter.

**ERRID\_ATMLE\_REGAIN\_SW**

Contact with the ATM switch has been re-established (for example, the cable has been plugged back in).

**ERRID\_ATMLE\_NET\_FAIL**

The device driver has gone into Network Recovery Mode in an attempt to recover from a network error and is temporarily unavailable during the recovery procedure. User intervention is not required for this error unless the problem persists.

**ERRID\_ATMLE\_RCVRY\_CMPLTE**

The network error that caused the **LANE** device driver to go into error recovery mode has been corrected.

## Adding an ATM MPOA Client

A Multi-Protocol Over ATM (MPOA) Client (MPC) can be added to the system to allow ATM LANE packets that would normally be routed through various LANE IP Subnets or Logical IP Subnets (LISs) within an ATM network, to be sent and received over shortcut paths that do not contain routers. MPOA can provide significant savings on end-to-end throughput performance for large data transfers, and can free up resources in routers that might otherwise be used up handling packets that could have bypassed routers altogether.

Only one MPOA Client is established per node. This MPC can support multiple ATM ports, containing LE Clients/Servers and MPOA Servers. The key requirement being, that for this MPC to create shortcut paths, each remote target node must also support MPOA Client, and must be directly accessible using the matrix of switches representing the ATM network.

A user with root authority can add this MPOA Client using the **smit mpoa\_panel** fast path, or click **Devices** → **Communication** → **ATM Adapter** → **Services** → **Multi-Protocol Over ATM (MPOA)**.

No configuration entries are required for the MPOA Client. Ease-of-use default values are provided for each of the attributes derived from ATM Forum recommendations.

Configuration help text is also available within MPOA Client SMIT to aid in making any modifications to attribute default values.

## Configuration Parameters for ATM MPOA Client

The **ATM LANE** device driver supports the following configuration parameters for the MPOA Client:

<b><i>auto_cfg</i></b>	Auto Configuration with LEC/LECS. Specifies whether the MPOA Client is to be automatically configured using LANE Configuration Server (LECS). Select <b>Yes</b> if a primary LE Client is used to obtain the MPOA configuration attributes, which overrides any manual or default values. The default value is No (manual configuration). The attribute values are: <b>Yes</b> - auto configuration <b>No</b> - manual configuration
<b><i>debug_trace</i></b>	Specifies whether this MPOA Client should keep a real time debug log within the kernel and allow full system trace capability. Select <b>Yes</b> to enable full tracing capabilities for this MPOA Client. Select <b>No</b> for optimal performance when minimal tracing is desired. The default is <b>Yes</b> (full tracing capability).
<b><i>fragment</i></b>	Enables MPOA fragmentation and specifies whether fragmentation should be performed on packets that exceed the maximum transmission unit (MTU) returned in the MPOA Resolution Reply. Select <b>Yes</b> to have outgoing packets fragmented as needed. Select <b>No</b> to avoid having outgoing packets fragmented. Selecting No causes outgoing packets to be sent down the LANE path when fragmentation must be performed. Incoming packets are always fragmented as needed even if No has been selected. The default value is <b>Yes</b> .
<b><i>hold_down_time</i></b>	Failed resolution request retry Hold Down Time (in seconds). Specifies the length of time to wait before reinitiating a failed address resolution attempt. This value is normally set to a value greater than <i>retry_time_max</i> . This attribute correlates to ATM Forum MPC Configuration parameter <i>MPC-p6</i> . The default value is 160 seconds.
<b><i>init_retry_time</i></b>	Initial Request Retry Time (in seconds). Specifies the length of time to wait before sending the first retry of a request that does not receive a response. This attribute correlates to ATM Forum MPC Configuration parameter <i>MPC-p4</i> . The default value is 5 seconds.
<b><i>retry_time_max</i></b>	Maximum Request Retry Time (in seconds). Specifies the maximum length of time to wait when retrying requests that have not received a response. Each retry duration after the initial retry are doubled (2x) until the retry duration reaches this Maximum Request Retry Time. All subsequent retries wait this maximum value. This attribute correlates to ATM Forum MPC Configuration parameter <i>MPC-p5</i> . The default value is 40 seconds.
<b><i>sc_setup_count</i></b>	Shortcut Setup Frame Count. This attribute is used in conjunction with <i>sc_setup_time</i> to determine when to establish a shortcut path. After the MPC has forwarded at least <i>sc_setup_count</i> packets to the same target within a period of <i>sc_setup_time</i> , the MPC attempts to create a shortcut VCC. This attribute correlates to ATM Forum MPC Configuration parameter <i>MPC-p1</i> . The default value is 10 packets.
<b><i>sc_setup_time</i></b>	Shortcut Setup Frame Time (in seconds). This attribute is used in conjunction with <i>sc_setup_count</i> above to determine when to establish a shortcut path. After the MPC has forwarded at least <i>sc_setup_count</i> packets to the same target within a period of <i>sc_setup_time</i> , the MPC attempts to create a shortcut VCC. This attribute correlates to ATM Forum MPC Configuration parameter <i>MPC-p2</i> . The default value is 1 second.
<b><i>vcc_inact_time</i></b>	VCC Inactivity Timeout value (in minutes). Specifies the maximum length of time to keep a shortcut VCC enabled when there is no send or receive activity on that VCC. The default value is 20 minutes.



## Tracing and Error Logging in the ATM MPOA Client

The ATM MPOA Client has two trace points:

- 3A3 - Normal Code Paths
- 3A4 - Error Conditions

Tracing can be enabled through SMIT or with the **trace** command.

```
trace -a -j 3a3,3a4
```

Tracing can be disabled through SMIT or with the **trcstop** command. After trace is stopped, the results can be formatted into readable text with the **trcrpt** command.

```
trcrpt > /tmp/trc.out
```

### MPOA Client error log templates

Each of the MPOA Client error log templates are prefixed with **ERRID\_MPOA**. An example of an MPOA error entry is as follows:

#### **ERRID\_MPOA\_MEM\_ERR**

An error occurred while attempting to allocate kernel memory.

## Getting Client Status

Three commands are available to obtain status information related to ATM **LANE** clients.

- The **entstat** command and **tokstat** command are used to obtain general ethernet or tokenring device status.
- The **lecstat** command is used to obtain more specific information about a **LANE** client.
- The **mpcstat** command is used to obtain MPOA client status information.

For more information see, **entstat Command**, **lecstat Command**, **mpcstat Command**, and **tokstat Command** in *AIX Version 6.1 Commands Reference*.

---

## Fiber Distributed Data Interface (FDDI) Device Driver

**Note:** The information in this section is specific to AIX 5.1 and earlier.

The FDDI device driver is a dynamically loadable device driver. The device driver is automatically loaded into the system at device configuration time as part of the configuration process.

The interface to the device is through the kernel services known as Network Services.

Interfacing to the device driver is achieved by calling the device driver's entry points for opening the device, closing the device, transmitting data, doing a remote dump, and issuing device control commands.

The FDDI device driver supports the SMT 7.2 standard.

## Configuration Parameters for FDDI Device Driver

### Software Transmit Queue

The driver provides a software transmit queue to supplement the hardware queue. The queue is configurable and contains between 3 and 250 mbufs. The default is 30 mbufs.

### Alternate Address

The driver supports specifying a configurable alternate address to be used instead of the address burned in on the card. This address must have the local bit set. Addresses between 0x400000000000 and 0x7FFFFFFFFFFFFF are supported. The default is 0x400000000000.

### **Enable Alternate Address**

The driver supports enabling the alternate address set with the Alternate Address parameter. Values are YES and NO, with NO as the default.

### **PMF Password**

The driver provides the ability to configure a PMF password. The password default is 0, meaning no password.

### **Max T-Req**

The driver enables the user to configure the card's maximum T-Req.

### **TVX Lower Bound**

The driver enables the user to configure the card's TVX Lower Bound.

### **User Data**

The driver enables the user to set the user data field on the adapter. This data can be any string up to 32 bytes of data. The default is a zero length string.

## **FDDI Device Driver Configuration and Unconfiguration**

The **fddi\_config** entry point performs configuration functions for the FDDI device driver.

### **Device Driver Open**

The **fddi\_open** function is called to open the specified network device.

The device is initialized. When the resources have been successfully allocated, the device is attached to the network.

If the station is not connected to another running station, the device driver opens, but is unable to transmit Logical Link Control (LLC) packets. When in this mode, the device driver sets the CFDDI\_NDD\_LLC\_DOWN flag (defined in **/usr/include/sys/cdli\_fddiuser.h**). When the adapter is able to make a connection with at least one other station this flag is cleared and LLC packets can be transmitted.

### **Device Driver Close**

The **fddi\_close** function is called to close the specified network device. This function resets the device to a known state and frees system resources used by the device.

The device is not detached from the network until the device's transmit queue is allowed to drain.

### **Data Transmission**

The **fddi\_output** function transmits data using the network device.

The FDDI device driver supports up to three mbuf's for each packet. It cannot gather from more than three locations to a packet.

The FDDI device driver does *not* accept user-memory mbufs. It uses **bcopy** on small frames which does not work on user memory.

The driver supports up to the entire mtu in a single mbuf.

The driver requires that the entire mac header be in a single mbuf.

The driver will not accept chained frames of different types. The user should not send Logical Link Control (LLC) and station management (SMT) frames in the same call to output.

The user needs to fill the frame out completely before calling the output routine. The mac header for a FDDI packet is defined by the **cfddi\_hdr\_t** structure defined in **/usr/include/sys/cdli\_fddiuser.h**. The first

byte of a packet is used as a flag for routing the packet on the adapter. For most driver users the value of the packet should be set to `FDDI_TX_NORM`. The possible flags are:

#### **CFDDI\_TX\_NORM**

Transmits the frame onto the ring. This is the normal flag value.

#### **CFDDI\_TX\_LOOPBACK**

Moves the frame from the adapter's transmit queue to its receive queue as if it were received from the media. The frame is not transmitted onto the media.

#### **CFDDI\_TX\_PROC\_ONLY**

Processes the status information frame (SIF) or parameter management frame (PMF) request frame and sends a SIF or PMF response to the host. The frame is not transmitted onto the media. This flag is *not* valid for LLC packets.

#### **CFDDI\_TX\_PROC\_XMIT**

Processes the SIF or PMF request frames and sends a SIF or PMF response to the host. The frame is also transmitted onto the media. This flag is *not* valid for LLC packets.

## **Data Reception**

When the FDDI device driver receives a valid packet from the network device, the FDDI device driver calls the `nd_receive` function that is specified in the `ndd_t` structure of the network device. The `nd_receive` function is part of a CDLI network demuxer. The packet is passed to the `nd_receive` function in mbufs.

## **Reliability, Availability, and Serviceability for FDDI Device Driver**

The FDDI device driver has three trace points. The IDs are defined in the `/usr/include/sys/cdli_fddiuser.h` file.

For FDDI the type of data in an error log is the same for every error log. Only the specifics and the title of the error log change. Information that follows includes an example of an error log and a list of error log entries.

### **Example FDDI Error Log**

Detail Data

FILE NAME

line: 332 file: fddiintr\_b.c

POS REGISTERS

F48E D317 3CC7 0008

SOURCE ADDRESS

4000 0000 0000

ATTACHMENT CLASS

0000 0001

MICRO CHANNEL AND PIO EXCEPTION CODES

0000 0000 0000 0000 0000 0000

FDDI LINK STATISTICS

0080 0000 04A0 0000 0000 0000 0001 0000 0000 0000

0001 0008 0008 0005 0005 0012 0003 0002 0000 0000

0000 0000 0000 0000 0000 0000 0000 0000

SELF TESTS

0000 0000 0000 0000 0000 0000 0000 0000 0000 0000

0000 0000 0000

DEVICE DRIVER INTERNAL STATE

0fdd 0fdd 0000 0000 0000 0000 0000 0000

## Error Log Entries

The FDDI device driver returns the following are the error log entries:

### **ERRID\_CFDDI\_RMV\_ADAP**

This error indicates that the adapter has received a disconnect command from a remote station. The FDDI device driver will initiate shutdown of the device. The device is no longer functional due to this error. User intervention is required to bring the device back online.

If there is no local LAN administrator, user action is required to make the device available. For the device to be brought back online, the device needs to be reset. This can be accomplished by having all users of the FDDI device driver close the device. When all users have closed the device and the device is reset, the device can be brought back online.

### **ERRID\_CFDDI\_ADAP\_CHECK**

This error indicates that an FDDI adapter check has occurred. If the device was connected to the network when this error occurred, the FDDI device goes into Network Recovery Mode in an attempt to recover from the error. The device is temporarily unavailable during the recovery procedure. User intervention is not required to bring the device back online.

### **ERRID\_CFDDI\_DWNLD**

Indicates that the microcode download to the FDDI adapter has failed. If this error occurs during the configuration of the device, the configuration of the device fails. User intervention is required to make the device available.

### **ERRID\_CFDDI\_RCVRY\_ENTER**

Indicates that the FDDI device driver has entered Network Recovery Mode in an attempt to recover from an error. The error which caused the device to enter this mode, is error logged before this error log entry. The device is not fully functional until the device has left this mode. User intervention is not required to bring the device back online.

### **ERRID\_CFDDI\_RCVRY\_EXIT**

Indicates that the FDDI device driver has successfully recovered from the error which caused the device to go into Network Recovery Mode. The device is now fully functional.

### **ERRID\_CFDDI\_RCVRY\_TERM**

Indicates that the FDDI device driver was unable to recover from the error which caused the device to go into Network Recovery Mode and has terminated recovery logic. The termination of recovery logic might be due to an irrecoverable error being detected or the device being closed. If termination is due to an irrecoverable error, that error will be error logged before this error log entry. User intervention is required to bring the device back online.

### **ERRID\_CFDDI\_MC\_ERR**

Indicates that the FDDI device driver has detected a Micro Channel error. The device driver initiates recovery logic in an attempt to recover from the error. User intervention is not required for this error unless the problem persists.

### **ERRID\_CFDDI\_TX\_ERR**

Indicates that the FDDI device driver has detected a transmission error. User intervention is not required unless the problem persists.

### **ERRID\_CFDDI\_PIO**

Indicates the FDDI device driver has detected a program IO error. The device driver initiates recovery logic in an attempt to recover from the error. User intervention is not required for this error unless the problem persists.

### **ERRID\_CFDDI\_DOWN**

Indicates that the FDDI device has been shutdown due to an irrecoverable error. The FDDI device is no longer functional due to the error. The irrecoverable error which caused the device to be shutdown is error logged before this error log entry. User intervention is required to bring the device back online.

**ERRID\_CFDDI\_SELF\_TEST**

Indicates that the FDDI adapter has received a run self-test command from a remote station. The device is unavailable while the adapter's self-tests are being run. If the tests are successful, the FDDI device driver initiates logic to reconnect the device to the network. Otherwise, the device will be shutdown.

**ERRID\_CFDDI\_SELF\_ERR**

Indicates that an error occurred during the FDDI self-tests. User intervention is required to bring the device back online.

**ERRID\_CFDDI\_PATH\_ERR**

Indicates that an error occurred during the FDDI adapter's path tests. The FDDI device driver will initiate recovery logic in an attempt to recover from the error. The FDDI device will temporarily be unavailable during the recovery procedure. User intervention is not required to bring the device back online.

**ERRID\_CFDDI\_PORT**

Indicates that a port on the FDDI device is in a stuck condition. User intervention is not required for this error. This error typically occurs when a cable is not correctly connected.

**ERRID\_CFDDI\_BYPASS**

Indicates that the optical bypass switch is in a stuck condition. User intervention is not required for this error.

**ERRID\_CFDDI\_CMD\_FAIL**

Indicates that a command to the adapter has failed.

---

## High-Performance (8fc8) Token-Ring Device Driver

**Note:** The information in this section is specific to AIX 5.1 and earlier.

The 8fc8 Token-Ring device driver is a dynamically loadable device driver. The device driver automatically loads into the system at device configuration time as part of the configuration process.

The interface to the device is through the kernel services known as Network Services.

Interfacing to the device driver is achieved by calling the device driver's entry points for opening the device, closing the device, transmitting data, doing a remote dump, and issuing device control commands.

The Token-Ring device driver interfaces with the Token-Ring High-Performance Network Adapter (8fc8). It provides a Micro Channel-based connection to a Token-Ring network. The adapter is IEEE 802.5 compatible and supports both 4 and 16 megabit per second networks. The adapter supports only a Shielded Twisted-Pair (STP) Token-Ring connection.

## Configuration Parameters for Token-Ring Device Driver

### Ring Speed

The device driver will support a user configurable parameter that indicates if the Token-Ring is to be run at 4 or 16 megabits per second.

### Software Transmit Queue

The device driver will support a user configurable transmit queue, that can be set to store between 32 and 160 transmit request pointers. Each transmit request pointer corresponds to a transmit request, which might be for several buffers of data.

### Attention MAC frames

The device driver will support a user configurable parameter that indicates if attention MAC frames should be received.

### Beacon MAC frames

The device driver will support a user configurable parameter that indicates if beacon MAC frames should be received.

### Network Address

The driver supports the use of the device's hardware address as the network address or an alternate network address configured through software. When an alternate address is used, any valid individual address can be used. The most significant bit of the address must be set to zero (definition of an individual address).

## Device Driver Configuration and Unconfiguration

The **tok\_config** entry point performs configuration functions Token-Ring device driver.

### Device Driver Open

The **tok\_open** function is called to open the specified network device.

The Token Ring device driver does an asynchronous open. It starts the process of attaching the device to the network, sets the NDD\_UP flag in the ndd\_flags field, and returns 0. The network attachment will continue in the background where it is driven by device activity and system timers.

**Note:** The Network Services **ns\_alloc** routine that calls this open routine causes the open to be synchronous. It waits until the NDD\_RUNNING flag is set in the ndd\_flags field or 60 seconds have passed.

If the connection is successful, the NDD\_RUNNING flag will be set in the ndd\_flags field and a NDD\_CONNECTED status block will be sent. The **ns\_alloc** routine will return at this time.

If the device connection fails, the NDD\_LIMBO flag will be set in the ndd\_flags field and a NDD\_LIMBO\_ENTRY status block will be sent.

If the device is eventually connected, the NDD\_LIMBO flag will be turned off and the NDD\_RUNNING flag will be set in the ndd\_flags field. Both NDD\_CONNECTED and NDD\_LIMBO\_EXIT status blocks will be set.

### Device Driver Close

The **tok\_close** function is called to close the specified network device. This function resets the device to a known state and frees system resources associated with the device.

The device will not be detached from the network until the device's transmit queue is allowed to drain.

## Data Transmission

The **tok\_output** function transmits data using the network device.

The device driver does *not* support mbufs from user memory (which have the M\_EXT flag set).

If the destination address in the packet is a broadcast address, the M\_BCAST flag in the p\_mbuf->m\_flags field should be set prior to entering this routine. A broadcast address is defined as 0xFFFF FFFF FFFF or 0xC000 FFFF FFFF. If the destination address in the packet is a multicast address the M\_MCAST flag in the p\_mbuf->m\_flags field should be set prior to entering this routine. A multicast address is defined as a non-individual address other than a broadcast address. The device driver will keep statistics based upon the M\_BCAST and M\_MCAST flags.

If a packet is transmitted with a destination address that matches the adapter's address, the packet will be received. This is true for the adapter's physical address, broadcast addresses (0xC000 FFFF FFFF or 0xFFFF FFFF FFFF), enabled functional addresses, or an enabled group address.

## Data Reception

When the Token-Ring device driver receives a valid packet from the network device, the Token-Ring device driver calls the **nd\_receive** function that is specified in the `ndd_t` structure of the network device. The **nd\_receive** function is part of a CDLI network demuxer. The packet is passed to the **nd\_receive** function in mbufs.

The Token-Ring device driver passes one packet to the **nd\_receive** function at a time.

The device driver sets the `M_BCAST` flag in the `p_mbuf->m_flags` field when a packet is received that has an all-stations broadcast address. This address is defined as `0xFFFF FFFF FFFF` or `0xC000 FFFF FFFF`.

The device driver sets the `M_MCAST` flag in the `p_mbuf->m_flags` field when a packet is received that has a non-individual address that is different than the all-stations broadcast address.

The adapter does not pass invalid packets to the device driver.

## Asynchronous Status

When a status event occurs on the device, the Token-Ring device driver builds the appropriate status block and calls the **nd\_status** function that is specified in the `ndd_t` structure of the network device. The **nd\_status** function is part of a CDLI network demuxer.

The following status blocks are defined for the Token-Ring device driver.

### Hard Failure

When a hard failure has occurred on the Token-Ring device, the following status blocks can be returned by the Token-Ring device driver. One of these status blocks indicates that a fatal error occurred.

**NDD\_PIO\_FAIL:** When a PIO error occurs, it is retried 3 times. If the error still occurs, it is considered unrecoverable and this status block is generated.

<b>code</b>	Set to <code>NDD_HARD_FAIL</code>
<b>option[0]</b>	Set to <code>NDD_PIO_FAIL</code>
<b>option[]</b>	The remainder of the status block may be used to return additional status information.

**TOK\_RECOVERY\_THRESH:** When most network errors occur, they are retried. Some errors are retried with no limit and others have a recovery threshold. Errors that have a recovery threshold and fail all the retries specified by the recovery threshold are considered unrecoverable and generate the following status block:

<b>code</b>	Set to <code>NDD_HARD_FAIL</code>
<b>option[0]</b>	Set to <code>TOK_RECOVERY_THRESH</code>
<b>option[1]</b>	The specific error that occurred. Possible values are: <ul style="list-style-type: none"><li>• <code>TOK_DUP_ADDR</code> - duplicate node address</li><li>• <code>TOK_PERM_HW_ERR</code> - the device has an unrecoverable hardware error</li><li>• <code>TOK_RING_SPEED</code> - ring beaconing on physical insertion to the ring</li><li>• <code>TOK_RMV_ADAP</code> - remove ring station MAC frame received</li></ul>

## Enter Network Recovery Mode

When the device driver has detected an error that requires initiating recovery logic that will make the device temporarily unavailable, the following status block is returned by the device driver:

**Note:** While the device driver is in this recovery logic, the device might not be fully functional. The device driver will notify users when the device is fully functional by way of an NDD\_LIMBO\_EXIT asynchronous status block.

**NDD\_ADAP\_CHECK:** When an adapter check has occurred, this status block is generated.

<b>code</b>	Set to NDD_LIMBO_ENTER
<b>option[0]</b>	Set to NDD_ADAP_CHECK
<b>option[1]</b>	The adapter check interrupt information is stored in the 2 high-order bytes. The adapter also returns three two-byte parameters. Parameter 0 is stored in the 2 low-order bytes.
<b>option[2]</b>	Parameter 1 is stored in the 2 high-order bytes. Parameter 2 is stored in the 2 low-order bytes.

**NDD\_AUTO\_RMV:** When an internal hardware error following the beacon automatic-removal process has been detected, this status block is generated.

<b>code</b>	Set to NDD_LIMBO_ENTER
<b>option[0]</b>	Set to NDD_AUTO_RMV

**NDD\_BUS\_ERR:** The device has detected a I/O channel error.

<b>code</b>	Set to NDD_LIMBO_ENTER
<b>option[0]</b>	Set to NDD_BUS_ERR
<b>option[1]</b>	Set to error information from the device.

**NDD\_CMD\_FAIL:** The device has detected an error in a command the device driver issued to it.

<b>code</b>	Set to NDD_LIMBO_ENTER
<b>option[0]</b>	Set to NDD_CMD_FAIL
<b>option[1]</b>	Set to error information from the device.

**NDD\_TX\_ERROR:** The device has detected an error in a packet given to the device.

<b>code</b>	Set to NDD_LIMBO_ENTER
<b>option[0]</b>	Set to NDD_TX_ERROR
<b>option[1]</b>	Set to error information from the device.

**NDD\_TX\_TIMEOUT:** The device has detected an error in a packet given to the device.

<b>code</b>	Set to NDD_LIMBO_ENTER
<b>option[0]</b>	Set to NDD_TX_TIMEOUT

**TOK\_ADAP\_INIT:** When the initialization of the device fails, this status block is generated.

<b>code</b>	Set to NDD_LIMBO_ENTER
<b>option[0]</b>	Set to TOK_ADAP_INIT
<b>option[1]</b>	Set to error information from the device.



**TOK\_ADAP\_OPEN:** When a general error occurs during open of the device, this status block is generated.

**code** Set to NDD\_LIMBO\_ENTER  
**option[0]** Set to TOK\_ADAP\_OPEN  
**option[1]** Set to the device open error code from the device.

**TOK\_DMA\_FAIL:** A d\_complete has failed.

**code** Set to NDD\_LIMBO\_ENTER  
**option[0]** Set to TOK\_DMA\_FAIL

**TOK\_RING\_SPEED:** When an error code of 0x27 (physical insertion, ring beaconing) occurs during open of the device, this status block is generated.

**code** Set to NDD\_LIMBO\_ENTER  
**option[0]** Set to TOK\_RING\_SPEED

**TOK\_RMV\_ADAP:** The device has received a remove ring station MAC frame indicating that a network management function had directed this device to get off the ring.

**code** Set to NDD\_LIMBO\_ENTER  
**option[0]** Set to TOK\_RMV\_ADAP

**TOK\_WIRE\_FAULT:** When an error code of 0x11 (lobe media test, function failure) occurs during open of the device, this status block is generated.

**code** Set to NDD\_LIMBO\_ENTER  
**option[0]** Set to TOK\_WIRE\_FAULT

## Exit Network Recovery Mode

When the device driver has successfully completed recovery logic from the error that made the device temporarily unavailable, the following status block is returned by the device driver. This status block means the device is now fully functional.

**code** Set to NDD\_LIMBO\_EXIT  
**option[]** The option fields are not used.

## Network Device Driver Status

When the device driver has status or event information to report, the following status block is returned by the device driver:

**Ring Beaconing:** When the Token-Ring device has detected a beaconing condition (or the ring has recovered from one), the following status block is generated by the Token-Ring device driver:

**code** Set to NDD\_STATUS  
**option[0]** Set to TOK\_BEACONING  
**option[1]** Set to the ring status received from the device.

## Device Connected

When the device is successfully connected to the network the following status block is returned by the device driver:

**code** Set to NDD\_CONNECTED  
**option[]** The option fields are not used.

## Device Control Operations

The **tok\_ctl** function is used to provide device control functions.

### NDD\_GET\_STATS

The user should pass in the **tok\_ndd\_stats\_t** structure as defined in **usr/include/sys/cdli\_tokuser.h**. The driver will fail a call with a buffer smaller than the structure.

The statistics that are returned contain statistics obtained from the device. If the device is inoperable, the statistics that are returned will not contain the current device statistics. The copy of the **ndd\_flags** field can be checked to determine the state of the device.

### NDD\_MIB\_QUERY

The *arg* parameter specifies the address of the **token\_ring\_all\_mib\_t** structure. This structure is defined in the **/usr/include/sys/tokenring\_mibs.h** file.

The device driver does *not* support any variables for read\_write or write only. If the syntax of a member of the structure is some integer type, the level of support flag will be stored in the whole field, regardless of the size of the field. For those fields defined as character arrays, the value will be returned only in the first byte in the field.

### NDD\_MIB\_GET

The *arg* parameter specifies the address of the **token\_ring\_all\_mib\_t** structure. This structure is defined in the **/usr/include/sys/tokenring\_mibs.h** file.

If the device is inoperable, the upstream field of the **Dot5Entry\_t** structure will be zero instead of containing the nearest active upstream neighbor (NAUN). Also the statistics that are returned contain statistics obtained from the device. If the device is inoperable, the statistics that are returned will not contain the current device statistics. The copy of the **ndd\_flags** field can be checked to determine the state of the device.

### NDD\_ENABLE\_ADDRESS

This command enables the receipt of packets with a functional or a group address. The functional address indicator (bit 0 "the MSB" of byte 2) indicates whether the address is a functional address (the bit is a 0) or a group address (the bit is a 1). The length field is not used because the address must be 6 bytes in length.

**Functional Address:** The specified address is ORed with the currently specified functional addresses and the resultant address is set as the functional address for the device. Functional addresses are encoded in a bit-significant format, thereby allowing multiple individual groups to be designated by a single address.

The Token-Ring network architecture provides bit-specific functional addresses for widely-used functions, such as configuration report server. Ring stations use functional address masks to identify these functions. For example, if function G is assigned a functional address of 0xC000 0008 0000, and function M is assigned a function address of 0xC000 0000 0040, then ring station Y, whose node contains function G and M, would have a mask of 0xC000 0008 0040. Ring station Y would receive packets addressed to either function G or M or to an address like 0xC000 0008 0048 because that address contains bits specified in the mask.

**Note:** The device forces the first 2 bytes of the functional address to be 0xC000. In addition, bits 6 and 7 of byte 5 of the functional address are forced to a 0 by the device.

The `NDD_ALTADDRS` and `TOK_RECEIVE_FUNC` flags in the `ndd_flags` field are set.

Because functional addresses are encoded in a bit-significant format, reference counts are kept on each of the 31 least significant bits of the address. Reference counts are not kept on the 17 most significant bits (the `0xC000` of the functional address and the functional address indicator bit).

**Group Address:** If no group address is currently enabled, the specified address is set as the group address for the device. The group address will not be set and `EINVAL` will be returned if a group address is currently enabled.

The device forces the first 2 bytes of the group address to be `0xC000`.

The `NDD_ALTADDRS` and `TOK_RECEIVE_GROUP` flags in the `ndd_flags` field are set.

## **NDD\_DISABLE\_ADDRESS**

This command disables the receipt of packets with a functional or a group address. The functional address indicator (bit 0 "the MSB" of byte 2) indicates whether the address is a functional address (the bit is a 0) or a group address (the bit is a 1). The length field is not used because the address must be 6 bytes in length.

**Functional Address:** The reference counts are decremented for those bits in the functional address that are a one (on). If the reference count for a bit goes to zero, the bit will be "turned off" in the functional address for the device.

If no functional addresses are active after receipt of this command, the `TOK_RECEIVE_FUNC` flag in the `ndd_flags` field is reset. If no functional or group addresses are active after receipt of this command, the `NDD_ALTADDRS` flag in the `ndd_flags` field is reset.

**Group Address:** If the group address that is currently enabled is specified, receipt of packets with a group address is disabled. If a different address is specified, `EINVAL` will be returned.

If no group address is active after receipt of this command, the `TOK_RECEIVE_GROUP` flag in the `ndd_flags` field is reset. If no functional or group addresses are active after receipt of this command, the `NDD_ALTADDRS` flag in the `ndd_flags` field is reset.

## **NDD\_MIB\_ADDR**

The following addresses are returned:

- Device Physical Address (or alternate address specified by user)
- Broadcast Address `0xFFFF FFFF FFFF`
- Broadcast Address `0xC000 FFFF FFFF`
- Functional Address (only if a user specified a functional address)
- Group Address (only if a user specified a group address)

## **NDD\_CLEAR\_STATS**

The counters kept by the device will be zeroed.

## **NDD\_GET\_ALL\_STATS**

The `arg` parameter specifies the address of the `mon_all_stats_t` structure. This structure is defined in the `/usr/include/sys/cdli_tokuser.h` file.

The statistics that are returned contain statistics obtained from the device. If the device is inoperable, the statistics that are returned will not contain the current device statistics. The copy of the `ndd_flags` field can be checked to determine the state of the device.

## Trace Points and Error Log Templates for 8fc8 Token-Ring Device Driver

The Token-Ring device driver has three trace points. The IDs are defined in the `usr/include/sys/cdli_tokuser.h` file.

The Token-Ring error log templates are:

### **ERRID\_CTOK\_ADAP\_CHECK**

The microcode on the device performs a series of diagnostic checks when the device is idle. These checks can find errors and they are reported as adapter checks. If the device was connected to the network when this error occurred, the device driver will go into Network Recovery Mode in an attempt to recover from the error. The device is temporarily unavailable during the recovery procedure. User intervention is not required for this error unless the problem persists.

### **ERRID\_CTOK\_ADAP\_OPEN**

The device driver was unable to open the device. The device driver will go into Network Recovery Mode in an attempt to recover from the error. The device is temporarily unavailable during the recovery procedure. User intervention is not required for this error unless the problem persists.

### **ERRID\_CTOK\_AUTO\_RMV**

An internal hardware error following the beacon automatic removal process has been detected. The device driver will go into Network Recovery Mode in an attempt to recover from the error. The device is temporarily unavailable during the recovery procedure. User intervention is not required for this error unless the problem persists.

### **ERRID\_CONFIG**

The ring speed (or ring data rate) is probably wrong. Contact the network administrator to determine the speed of the ring. The device driver will only retry twice at 2 minute intervals after this error log entry has been generated.

### **ERRID\_CTOK\_DEVICE\_ERR**

The device detected an I/O channel error or an error in a command the device driver issued, an error occurred during a PIO operation, or the device has detected an error in a packet given to the device. The device driver will go into Network Recovery Mode in an attempt to recover from the error. The device is temporarily unavailable during the recovery procedure. User intervention is not required for this error unless the problem persists.

### **ERRID\_CTOK\_DOWNLOAD**

The download of the microcode to the device failed. User intervention is required to make the device available.

### **ERRID\_CTOK\_DUP\_ADDR**

The device has detected that another station on the ring has a device address that is the same as the device address being tested. Contact network administrator to determine why.

### **ERRID\_CTOK\_MEM\_ERR**

An error occurred while allocating memory or timer control block structures.

### **ERRID\_CTOK\_PERM\_HW**

The device driver could not reset the card. For example, did not receive status from the adapter within the retry period.

### **ERRID\_CTOK\_RCVRY\_EXIT**

The error that caused the device driver to go into error recovery mode has been corrected.

### **ERRID\_CTOK\_RMV\_ADAP**

The device has received a remove ring station MAC frame indicating that a network management function has directed this device to get off the ring. Contact network administrator to determine why.

## **ERRID\_CTOK\_WIRE\_FAULT**

There is probably a loose (or bad) cable between the device and the MAU. There is some chance that it might be a bad device. The device driver will go into Network Recovery Mode in an attempt to recover from the error. The device is temporarily unavailable during the recovery procedure. User intervention is required for this error.

---

## **High-Performance (8fa2) Token-Ring Device Driver**

**Note:** The information in this section is specific to AIX 5.1 and earlier.

The 8fa2 Token-Ring device driver is a dynamically loadable device driver. The device driver is automatically loaded into the system at device configuration time as part of the configuration process.

The interface to the device is through the kernel services known as Network Services.

Interfacing to the device driver is achieved by calling the device driver's entry points for opening the device, closing the device, transmitting data, doing a remote dump, and issuing device control commands.

The Token-Ring device driver interfaces with the Token-Ring High-Performance Network Adapter (8fa2). It provides a Micro Channel-based connection to a Token-Ring network. The adapter is IEEE 802.5 compatible and supports both 4 and 16 megabit per second networks. The adapter supports only a RJ-45 connection.

## **Configuration Parameters for 8fa2 Token-Ring Device Driver**

The following lists the configuration parameters necessary to use the device driver.

### **Ring Speed**

Indicates the Token-Ring speed. The speed is set at 4 or 16 megabits per second or autosense.

- 4** Specifies that the device driver will open the adapter with 4 Mbits. It will return an error if ring speed does not match the network speed.
- 16** Specifies that the device driver will open the adapter with 16 Mbits. It will return an error if ring speed does not match the network speed.

### **autosense**

Specifies that the adapter will open with the speed used determined as follows:

- If it is an open on an existing network, the speed will be the ring speed of the network.
- If it is an open on a new network:
- If the adapter is a new adapter, 16 Mbits is used.
- If the adapter had successfully opened, the ring speed will be the ring speed of the last successful open.

### **Software Transmit Queue**

Specifies a transmit request pointer that can be set to store between 32 and 2048 transmit request pointers. Each transmit request pointer corresponds to a transmit request which might be for several buffers of data.

### **Attention MAC frames**

Indicates if attention MAC frames should be received.

### **Beacon MAC frames**

Indicates if beacon MAC frames should be received.

### **Priority Data Transmission**

Specifies a request priority transmission of the data packets.

## Network Address

Specifies the use of the device's hardware address as the network address or an alternate network address configured through software. When an alternate address is used, any valid Individual Address can be used. The most significant bit of the address must be set to zero (definition of an Individual Address).

## Device Driver Configuration and Unconfiguration

The **tok\_config** entry point performs configuration functions Token-Ring device driver.

## Device Driver Open

The **tok\_open** function is called to open the specified network device.

The Token Ring device driver does a synchronous open. The device will be initialized at this time. When the resources have been successfully allocated, the device will start the process of attaching the device to the network.

If the connection is successful, the NDD\_RUNNING flag will be set in the ndd\_flags field and a NDD\_CONNECTED status block will be sent.

If the device connection fails, the NDD\_LIMBO flag will be set in the ndd\_flags field and a NDD\_LIMBO\_ENTRY status block will be sent.

If the device is eventually connected, the NDD\_LIMBO flag will be turned off and the NDD\_RUNNING flag will be set in the ndd\_flags field. Both NDD\_CONNECTED and NDD\_LIMBO\_EXIT status blocks will be set.

## Device Driver Close

The **tok\_close** function is called to close the specified network device. This function resets the device to a known state and frees system resources associated with the device.

The device will not be detached from the network until the device's transmit queue is allowed to drain.

## Data Transmission

The **tok\_output** function transmits data using the network device.

The device driver does *not* support mbufs from user memory (which have the M\_EXT flag set).

If the destination address in the packet is a broadcast address the M\_BCAST flag in the **p\_mbuf->m\_flags** field should be set prior to entering this routine. A broadcast address is defined as 0xFFFF FFFF FFFF or 0xC000 FFFF FFFF. If the destination address in the packet is a multicast address the M\_MCAST flag in the **p\_mbuf->m\_flags** field should be set prior to entering this routine. A multicast address is defined as a non-individual address other than a broadcast address. The device driver will keep statistics based upon the M\_BCAST and M\_MCAST flags.

If a packet is transmitted with a destination address which matches the adapter's address, the packet will be received. This is true for the adapter's physical address, broadcast addresses (0xC000 FFFF FFFF or 0xFFFF FFFF FFFF), enabled functional addresses, or an enabled group address.

## Data Reception

When the Token-Ring device driver receives a valid packet from the network device, the Token-Ring device driver calls the **nd\_receive** function that is specified in the ndd\_t structure of the network device. The **nd\_receive** function is part of a CDLI network demuxer. The packet is passed to the **nd\_receive** function in mbufs.

The Token-Ring device driver will pass only one packet to the **nd\_receive** function at a time.

The device driver will set the M\_BCAST flag in the p\_mbuf->m\_flags field when a packet is received which has an all stations broadcast address. This address is defined as 0xFFFF FFFF FFFF or 0xC000 FFFF FFFF.

The device driver will set the M\_MCAST flag in the p\_mbuf->m\_flags field when a packet is received which has a non-individual address which is different than the all-stations broadcast address.

The adapter will not pass invalid packets to the device driver.

## Asynchronous Status

When a status event occurs on the device, the Token-Ring device driver builds the appropriate status block and calls the **nd\_status** function that is specified in the ndd\_t structure of the network device. The **nd\_status** function is part of a CDLI network demuxer.

The following status blocks are defined for the Token-Ring device driver.

### Hard Failure

When a hard failure has occurred on the Token-Ring device, the following status blocks can be returned by the Token-Ring device driver. One of these status blocks indicates that a fatal error occurred.

#### NDD\_PIO\_FAIL

Indicates that when a PIO error occurs, it is retried 3 times. If the error persists, it is considered unrecoverable and the following status block is generated:

**code** Set to NDD\_HARD\_FAIL  
**option[0]** Set to NDD\_PIO\_FAIL  
**option[]** The remainder of the status block is used to return additional status information.

#### NDD\_HARD\_FAIL

Indicates that when a transmit error occurs it is retried. If the error is unrecoverable, the following status block is generated:

**code** Set to NDD\_HARD\_FAIL  
**option[0]** Set to NDD\_HARD\_FAIL  
**option[]** The remainder of the status block is used to return additional status information.

#### NDD\_ADAP\_CHECK

Indicates that when an adapter check has occurred, the following status block is generated:

**code** Set to NDD\_ADAP\_CHECK  
**option[]** The remainder of the status block is used to return additional status information.

#### NDD\_DUP\_ADDR

Indicates that the device detected a duplicated address in the network and the following status block is generated:

**code** Set to NDD\_DUP\_ADDR  
**option[]** The remainder of the status block is used to return additional status information.

#### NDD\_CMD\_FAIL

Indicates that the device detected an error in a command that the device driver issued. The following status block is generated:

**code** Set to NDD\_CMD\_FAIL

**option[0]** Set to the command code  
**option[]** Set to error information from the command.

### TOK\_RING\_SPEED

Indicates that when a ring speed error occurs while the device is being open, the following status block is generated:

**code** Set to NDD\_LIMBO\_ENTER  
**option[]** Set to error information.

### Enter Network Recovery Mode

Indicates that when the device driver has detected an error which requires initiating recovery logic that will make the device temporarily unavailable, the following status block is returned by the device driver.

**Note:** While the device driver is in this recovery logic, the device might not be fully functional. The device driver will notify users when the device is fully functional by way of an NDD\_LIMBO\_EXIT asynchronous status block.

**code** Set to NDD\_LIMBO\_ENTER  
**option[0]** Set to one of the following:

- NDD\_CMD\_FAIL
- TOK\_WIRE\_FAULT
- NDD\_BUS\_ERROR
- NDD\_ADAP\_CHECK
- NDD\_TX\_TIMEOUT
- TOK\_BEACONING

**option[]** The remainder of the status block is used to return additional status information by the device driver.

### Exit Network Recovery Mode

Indicates that when the device driver has successfully completed recovery logic from the error that made the device temporarily unavailable, the following status block is returned by the device driver. This status block indicates the device is now fully functional.

**code** Set to NDD\_LIMBO\_EXIT  
**option[]** N/A

### Device Connected

Indicates that when the device is successfully connected to the network the following status block is returned by the device driver:

**code** Set to NDD\_CONNECTED  
**option[]** N/A

## Device Control Operations

The **tok\_ctl** function is used to provide device control functions.

### NDD\_GET\_STATS

The user should pass in the **tok\_ddd\_stats\_t** structure as defined in **<sys/cdli\_tokuser.h>**. The driver will fail a call with a buffer smaller than the structure.

The structure must be in a kernel heap so that the device driver can copy the statistics into it; and it must be pinned.



### **NDD\_PROMISCUOUS\_ON**

Setting promiscuous mode will *not* cause non-LLC frames to be received by the driver unless the user also enables those filters (Attention MAC frames, Beacon MAC frames).

The driver will maintain a counter of requests.

### **NDD\_PROMISCUOUS\_OFF**

This command will release a request from a user to PROMISCUOUS\_ON; it will not exit the mode on the adapter if more requests are outstanding.

### **NDD\_MIB\_QUERY**

The *arg* parameter specifies the address of the **token\_ring\_all\_mib\_t** structure. This structure is defined in the **/usr/include/sys/tokenring\_mibs.h** file.

The device driver does *not* support any variables for read\_write or write only. If the syntax of a member of the structure is some integer type, the level of support flag will be stored in the whole field, regardless of the size of the field. For those fields which are defined as character arrays, the value will be returned only in the first byte in the field.

### **NDD\_MIB\_GET**

The *arg* parameter specifies the address of the **token\_ring\_all\_mib\_t** structure. This structure is defined in the **/usr/include/sys/tokenring\_mibs.h** file.

### **NDD\_ENABLE\_ADDRESS**

This command enables the receipt of packets with a functional or a group address. The functional address indicator (bit 0 "the MSB" of byte 2) indicates whether the address is a functional address (the bit is a 0) or a group address (the bit is a 1). The length field is not used because the address must be 6 bytes in length.

## **Functional Address**

The specified address is ORed with the currently specified functional addresses and the resultant address is set as the functional address for the device. Functional addresses are encoded in a bit-significant format, thereby allowing multiple individual groups to be designated by a single address.

The Token-Ring network architecture provides bit-specific functional addresses for widely used functions, such as configuration report server. Ring stations use functional address *masks* to identify these functions. For example, if function G is assigned a functional address of 0xC000 0008 0000, and function M is assigned a function address of 0xC000 0000 0040, then ring station Y, whose node contains function G and M, would have a mask of 0xC000 0008 0040. Ring station Y would receive packets addressed to either function G or M or to an address like 0xC000 0008 0048 because that address contains bits specified in the mask.

The NDD\_ALTADDRS and TOK\_RECEIVE\_FUNC flags in the **ndd\_flags** field are set.

Because functional addresses are encoded in a bit-significant format, reference counts are kept on each of the 31 least significant bits of the address.

## **Group Address**

The device support 256 general group addresses. The promiscuous mode will be turned on when the group addresses needed to be set are more than 256. The device driver will maintain a reference count on this operation.

The NDD\_ALTADDRS and TOK\_RECEIVE\_GROUP flags in the **ndd\_flags** field are set.

### **NDD\_DISABLE\_ADDRESS**

This command disables the receipt of packets with a functional or a group address. The functional address indicator (bit 0 "the MSB" of byte 2) indicates whether the address is a functional address (the bit is a 0) or a group address (the bit is a 1). The length field is not used because the address must be 6 bytes in length.

## Functional Address

The reference counts are decremented for those bits in the functional address that are one (meaning *on*). If the reference count for a bit goes to zero, the bit will be "turned off" in the functional address for the device.

If no functional addresses are active after receipt of this command, the TOK\_RECEIVE\_FUNC flag in the **ndd\_flags** field is reset. If no functional or group addresses are active after receipt of this command, the NDD\_ALTADDRS flag in the **ndd\_flags** field is reset.

## Group Address

If the number of group address enabled is less than 256, the driver sends a command to the device to disable the receipt of the packets with the specified group address. Otherwise, the driver just deletes the group address from the group address table.

If there are less than 256 group addresses enabled after the receipt of this command, the promiscuous mode is turned off.

If no group address is active after receipt of this command, the TOK\_RECEIVE\_GROUP flag in the **ndd\_flags** field is reset. If no functional or group addresses are active after receipt of this command, the NDD\_ALTADDRS flag in the **ndd\_flags** field is reset.

## NDD\_PRIORITY\_ADDRESS

The driver returns the address of the device driver's priority transmit routine.

## NDD\_MIB\_ADDR

The driver will return at least three addresses: device physical address (or alternate address specified by user) and two broadcast addresses (0xFFFF FFFF FFFF and 0xC000 FFFF FFFF). Additional addresses specified by the user, such as functional address and group addresses, might also be returned.

## NDD\_CLEAR\_STATS

The counters kept by the device are zeroed.

## NDD\_GET\_ALL\_STATS

The *arg* parameter specifies the address of the **mon\_all\_stats\_t** structure. This structure is defined in the **/usr/include/sys/cdli\_tokuser.h** file.

The statistics returned include statistics obtained from the device. If the device is inoperable, the statistics returned do not contain the current device statistics. The copy of the **ndd\_flags** field can be checked to determine the state of the device.

## Trace Points and Error Log Templates for 8fa2 Token-Ring Device Driver

The Token-Ring device driver has four trace points. The IDs are defined in the **/usr/include/sys/cdli\_tokuser.h** file.

The Token-Ring error log templates are :

### ERRID\_MPS\_ADAP\_CHECK

The microcode on the device performs a series of diagnostic checks when the device is idle. These checks can find errors and they are reported as adapter checks. If the device was connected to the network when this error occurred, the device driver goes into Network Recovery Mode to try to recover from the error. The device is temporarily unavailable during the recovery procedure. User intervention is not required unless the problem persists.

### ERRID\_MPS\_ADAP\_OPEN

The device driver was unable to open the device. The device driver goes into Network Recovery Mode to try to recover from the error. The device is temporarily unavailable during the recovery procedure. User intervention is not required unless the problem persists.

**ERRID\_MPS\_AUTO\_RMV**

An internal hardware error following the beacon automatic removal process has been detected. The device driver goes into Network Recovery Mode to try to recover from the error. The device is temporarily unavailable during the recovery procedure. User intervention is not required unless the problem persists.

**ERRID\_MPS\_RING\_SPEED**

The ring speed (or ring data rate) is probably wrong. Contact the network administrator to determine the speed of the ring. The device driver only retries twice at 2 minute intervals when this error log entry is generated.

**ERRID\_MPS\_DMAFAIL**

The device detected a DMA error in a TX or RX operation. The device driver goes into Network Recovery Mode to try to recover from the error. The device is temporarily unavailable during the recovery procedure. User intervention is not required unless the problem persists.

**ERRID\_MPS\_BUS\_ERR**

The device detected a Micro Channel bus error. The device driver goes into Network Recovery Mode to try to recover from the error. The device is temporarily unavailable during the recovery procedure. User intervention is not required unless the problem persists.

**ERRID\_MPS\_DUP\_ADDR**

The device has detected that another station on the ring has a device address which is the same as the device address being tested. Contact the network administrator to determine why.

**ERRID\_MPS\_MEM\_ERR**

An error occurred while allocating memory or timer control block structures.

**ERRID\_MPS\_PERM\_HW**

The device driver could not reset the card. For example, it did not receive status from the adapter within the retry period.

**ERRID\_MPS\_RCVRY\_EXIT**

The error that caused the device driver to go into error recovery mode has been corrected.

**ERRID\_MPS\_RMV\_ADAP**

The device has received a remove ring station MAC frame indicating that a network management function has directed this device to get off the ring. Contact the network administrator to determine why.

**ERRID\_MPS\_WIRE\_FAULT**

There is probably a loose (or bad) cable between the device and the MAU. There is some chance that it might be a bad device. The device driver goes into Network Recovery Mode to try to recover from the error. The device is temporarily unavailable during the recovery procedure. User intervention is required for this error.

**ERRID\_MPS\_RX\_ERR**

The device detected a receive error. The device driver goes into Network Recovery Mode to try to recover from the error. The device is temporarily unavailable during the recovery procedure. User intervention is not required unless the problem persists.

**ERRID\_MPS\_TX\_TIMEOUT**

The transmit watchdog timer expired before transmitting a frame is complete. The device driver goes into Network Recovery Mode to try to recover from the error. The device is temporarily unavailable during the recovery procedure. User intervention is not required unless the problem persists.

**ERRID\_MPS\_CTL\_ERR**

The IOCTL watchdog timer expired before the device driver received a response from the device. The device driver goes into Network Recovery Mode to try to recover from the error. The device is temporarily unavailable during the recovery procedure. User intervention is not required unless the problem persists.

---

## PCI Token-Ring Device Drivers

The following Token-Ring device drivers are dynamically loadable. The device driver is automatically loaded into the system at device configuration time as part of the configuration process.

- PCI Token-Ring High Performance Device Driver (14101800)
- PCI Token-Ring Device Driver (14103e00)

The interface to the device is through the kernel services known as *Network Services*. Interfacing to the device driver is achieved by calling the device driver's entry points to perform the following actions:

- Opening the device
- Closing the device
- Transmitting data
- Issuing device control commands

The PCI Token-Ring High Performance Device Driver (14101800) interfaces with the PCI Token-Ring High-Performance Network Adapter (14101800). The adapter is IEEE 802.5 compatible and supports both 4 and 16 Mbps networks. The adapter supports only an RJ-45 connection.

The PCI Token-Ring Device Driver (14103e00) interfaces with the PCI Token-Ring Network Adapter (14103e00). The adapter is IEEE 802.5 compatible and supports both 4 and 16 Mbps networks. The adapter supports both an RJ-45 and a 9 Pin connection.

## Configuration Parameters

The following configuration parameter is supported by all PCI Token-Ring Device Drivers:

### Ring Speed

The device driver supports a user-configurable parameter that indicates if the token-ring is to run at 4 or 16 Mbps.

The device driver supports a user-configurable parameter that selects the ring speed of the adapter. There are three options for the ring speed: 4, 16, or autosense.

1. If 4 is selected, the device driver opens the adapter with 4 Mbits. It returns an error if the ring speed does not match the network speed.
2. If 16 is selected, the device driver opens the adapter with 16 Mbits. It returns an error if the ring speed does not match the network speed.
3. If autosense is selected, the adapter guarantees a successful open, and the speed used to open is dependent on the following:
  - If the adapter is opened on an existing network the speed is determined by the ring speed of the network.
  - If the device is opened on a new network and the adapter is new, 16 Mbits is used. Or, if the adapter opened successfully, the ring speed is determined by the speed of the last successful open.

### Software Transmit Queue

The device driver supports a user-configurable transmit queue that can be set to store between 32 and 2048 transmit request pointers. Each transmit request pointer corresponds to a transmit request that might be for several buffers of data.

### Receive Queue

The device driver supports a user-configurable receive queue that can be set to store between 32 and 160 receive buffers. These buffers are **mbuf** clusters into which the device writes the received data.

### **Full Duplex**

Indicates whether the adapter is operating in full-duplex or half-duplex mode. If this field is set to yes, the device driver programs the adapter to be in full-duplex mode. The default value is half-duplex.

### **Attention MAC Frames**

The device driver supports a user-configurable parameter that indicates if attention MAC frames should be received.

### **Beacon MAC Frames**

The device driver supports a user-configurable parameter that indicates if beacon MAC frames should be received.

### **Network Address**

The driver supports the use of the device's hardware address as the network address or an alternate network address configured through software. When an alternate address is used, any valid individual address can be used. The most significant bit of the address must be set to zero.

In addition, the following configuration parameters are supported by the PCI Token-Ring High Performance Device Driver (14101800):

### **Priority Data Transmission**

The device driver supports a user option to request priority transmission of the data packets.

### **Software Priority Transmit Queue**

The device driver supports a user-configurable priority transmit queue that can be set to store between 32 and 160 transmit request pointers. Each transmit request pointer corresponds to a transmit request that might be for several buffers of data.

## **Device Driver Configuration and Unconfiguration**

The configuration entry points of the device drivers conform to the guidelines for kernel object file entry points. These configuration entry points are as follows:

- **tok\_config** for the PCI Token-Ring High Performance Device Driver (14101800).
- **cs\_config** for the PCI Token-Ring Device Driver (14103e00).

## **Device Driver Open**

The Token-Ring device driver performs a synchronous open. The device is initialized at this time. When the resources are successfully allocated, the device starts the process of attaching the device to the network.

If the connection is successful, the **NDD\_RUNNING** flag is set in the `ndd_flags` field, and an **NDD\_CONNECTED** status block is sent.

If the device connection fails, the **NDD\_LIMBO** flag is set in the `ndd_flags` field, and an **NDD\_LIMBO\_ENTRY** status block is sent.

If the device is eventually connected, the **NDD\_LIMBO** flag is turned off, and the **NDD\_RUNNING** flag is set in the `ndd_flags` field. Both **NDD\_CONNECTED** and **NDD\_LIMBO\_EXIT** status blocks are set.

The entry points are as follows:

- **tok\_open** for the PCI Token-Ring High Performance Device Driver (14101800).
- **cs\_open** for the PCI Token-Ring Device Driver (14103e00).

## Device Driver Close

This function resets the device to a known state and frees system resources associated with the device.

The device is not detached from the network until the device's transmit queue is allowed to drain.

The close entry points are as follows:

- **tok\_close** for the PCI Token-Ring High Performance Device Driver (14101800).
- **cs\_close** for the PCI Token-Ring Device Driver (14103e00).

## Data Transmission

The device drivers do not support **mbuf** structures from user memory that have the **M\_EXT** flag set.

If the destination address in the packet is a broadcast address, the **M\_BCAST** flag in the `p_mbuf->m_flags` field must be set prior to entering this routine. A broadcast address is defined as 0xFFFF FFFF FFFF or 0xC000 FFFF FFFF. If the destination address in the packet is a multicast address, the **M\_MCAST** flag in the `p_mbuf->m_flags` field must be set prior to entering this routine. A multicast address is defined as a non-individual address other than a broadcast address. The device driver keeps statistics based on the **M\_BCAST** and **M\_MCAST** flags.

If a packet is transmitted with a destination address that matches the adapter's address, the packet is received. This is true for the adapter's physical address, broadcast addresses (0xC000 FFFF FFFF or 0xFFFF FFFF FFFF), enabled functional addresses, or an enabled group address.

The output entry points are as follows:

- **tok\_output** for the PCI Token-Ring High Performance Device Driver (14101800).
- **cs\_close** for the PCI Token-Ring Device Driver (14103e00).

## Data Reception

When the Token-Ring device driver receives a valid packet from the network device, the Token-Ring device driver calls the **nd\_receive()** function specified in the **ndd\_t** structure of the network device. The **nd\_receive()** function is part of a CDLI network demuxer. The packet is passed to the **nd\_receive()** function in the **mbuf** structures.

The Token-Ring device driver passes only one packet to the **nd\_receive()** function at a time.

The device driver sets the **M\_BCAST** flag in the `p_mbuf->m_flags` field when a packet that has an all-stations broadcast address is received. This address is defined as 0xFFFF FFFF FFFF or 0xC000 FFFF FFFF.

The device driver sets the **M\_MCAST** flag in the `p_mbuf->m_flags` field when a packet is received that has a non-individual address that is different from the all-stations broadcast address.

The adapter does not pass invalid packets to the device driver.

## Asynchronous Status

When a status event occurs on the device, the Token-Ring device driver builds the appropriate status block and calls the **nd\_status()** function specified in the **ndd\_t** structure of the network device. The **nd\_status()** function is part of a CDLI network demuxer.

The following status blocks are defined for the Token-Ring device driver.

## Hard Failure

When a hard failure occurs on the Token-Ring device, the following status blocks are returned by the Token-Ring device driver. One of these status blocks indicates that a fatal error has occurred.

### NDD\_HARD\_FAIL

When a transmit error occurs, it tries to recover. If the error is unrecoverable, this status block is generated.

**code** Set to NDD\_HARD\_FAIL.

#### option[0]

Set to NDD\_HARD\_FAIL.

#### option[ ]

The remainder of the status block can be used to return additional status information.

## Enter Network Recovery Mode

When the device driver detects an error that requires initiating recovery logic to make the device temporarily unavailable, the following status block is returned by the device driver.

**Note:** While the device driver is in this recovery logic, the device might not be fully functional. The device driver notifies users when the device is fully functional by way of an NDD\_LIMBO\_EXIT asynchronous status block:

**code** Set to NDD\_LIMBO\_ENTER.

**option[0]** Set to one of the following:

- NDD\_CMD\_FAIL
- NDD\_ADAP\_CHECK
- NDD\_TX\_ERR
- NDD\_TX\_TIMEOUT
- NDD\_AUTO\_RMV
- TOK\_ADAP\_OPEN
- TOK\_ADAP\_INIT
- TOK\_DMA\_FAIL
- TOK\_RING\_SPEED
- TOK\_RMV\_ADAP
- TOK\_WIRE\_FAULT

**option[ ]** The remainder of the status block can be used to return additional status information by the device driver.

## Exit Network Recovery Mode

When the device driver has successfully completed recovery logic from the error that made the device temporarily unavailable, the following status block is returned by the device driver:

**code** Set to NDD\_LIMBO\_EXIT.

**option[ ]** The option fields are not used.

The device is now fully functional.

## Device Control Operations

The **ndd\_ctl** entry point is used to provide device control functions.

### NDD\_GET\_STATS

The user should pass in the **tok\_ndd\_stats\_t** structure as defined in the **sys/cdli\_tokuser.h** file. The driver fails a call with a buffer smaller than the structure.

The structure must be in kernel heap so that the device driver can copy the statistics into it. Also, it must be pinned.

#### **NDD\_PROMISCUOUS\_ON**

Setting promiscuous mode will *not* cause non-LLC frames to be received by the driver unless the user also enables those filters (Attention MAC frames, Beacon MAC frames).

The driver maintains a counter of requests.

#### **NDD\_PROMISCUOUS\_OFF**

This command releases a request from a user to **PROMISCUOUS\_ON**; it will not exit the mode on the adapter if more requests are outstanding.

#### **NDD\_MIB\_QUERY**

The **arg** parameter specifies the address of the **token\_ring\_all\_mib\_t** structure. This structure is defined in the **/usr/include/sys/tokenring\_mibs.h** file.

The device driver does *not* support any variables for read\_write or write only. If the syntax of a member of the structure is an integer type, the level of support flag is stored in the whole field, regardless of the size of the field. For those fields that are defined as character arrays, the value is returned only in the first byte in the field.

#### **NDD\_MIB\_GET**

The **arg** parameter specifies the address of the **token\_ring\_all\_mib\_t** structure. This structure is defined in the **/usr/include/sys/tokenring\_mibs.h** file.

#### **NDD\_ENABLE\_ADDRESS**

This command enables the receipt of packets with a functional or a group address. The functional address indicator (bit 0 "the MSB" of byte 2) indicates whether the address is a functional address (bit 0) or a group address (bit 1). The length field is not used because the address must be 6 bytes in length.

##### **functional address**

The specified address is ORed with the currently specified functional addresses, and the resultant address is set as the functional address for the device. Functional addresses are encoded in a bit-significant format, thereby allowing multiple individual groups to be designated by a single address.

The Token-Ring network architecture provides bit-specific functional addresses for widely used functions, such as configuration report server. Ring stations use functional address "masks" to identify these functions. For example, if function G is assigned a functional address of 0xC000 0008 0000, and function M is assigned a function address of 0xC000 0000 0040, then ring station Y, whose node contains function G and M, would have a mask of 0xC000 0008 0040. Ring station Y would receive packets addressed to either function G or M or to an address, such as 0xC000 0008 0048, because that address contains bits specified in the "mask."

The **NDD\_ALTADDRS** and **TOK\_RECEIVE\_FUNC** flags in the **ndd\_flags** field are set.

Because functional addresses are encoded in a bit-significant format, reference counts are kept on each of the 31 least significant bits of the address.

##### **group address**

The device supports 256 general group addresses. The promiscuous mode is turned on when the group addresses to be set is more than 256. The device driver maintains a reference count on this operation.

The device supports 256 general group addresses. The promiscuous mode is turned on when the group address needed to be set are more than 256. The device driver will maintain a reference count on this operation.

The **NDD\_ALTADDRS** and **TOK\_RECEIVE\_GROUP** flags in the **ndd\_flags** field are set.



## **NDD\_DISABLE\_ADDRESS**

This command disables the receipt of packets with a functional or a group address. The functional address indicator (bit 0 "the MSB" of byte 2) indicates whether the address is a functional address (bit 0) or a group address (bit 1). The length field is not used because the address must be 6 bytes in length.

### **functional address**

The reference counts are decremented for those bits in the functional address that are 1 (on). If the reference count for a bit goes to 0, the bit is "turned off" in the functional address for the device.

If no functional addresses are active after receipt of this command, the **TOK\_RECEIVE\_FUNC** flag in the `ndd_flags` field is reset. If no functional or group addresses are active after receipt of this command, the **NDD\_ALTADDRS** flag in the `ndd_flags` field is reset.

### **group address**

If group address enable is less than 256, the driver sends a command to the device to disable the receipt of the packets with the specified group address. Otherwise, the group address is deleted from the group address table.

If there are less than 256 group addresses enabled after the receipt of this command, the promiscuous mode is turned off.

If no group address is active after receipt of this command, the **TOK\_RECEIVE\_GROUP** flag in the `ndd_flags` field is reset. If no functional or group addresses are active after receipt of this command, the **NDD\_ALTADDRS** flag in the `ndd_flags` field is reset.

## **NDD\_PRIORITY\_ADDRESS**

The driver returns the address of the device driver's priority transmit routine.

## **NDD\_MIB\_ADDR**

The driver returns at least three addresses that are device physical addresses (or alternate addresses specified by the user), two broadcast addresses (0xFFFFFFFF and 0xC000 FFFF FFFF), and any additional addresses specified by the user, such as functional addresses and group addresses.

## **NDD\_CLEAR\_STATS**

The counters kept by the device are zeroed.

## **NDD\_GET\_ALL\_STATS**

Used to gather all statistics for the specified device. The **arg** parameter specifies the address of the statistics structure for this particular device type. The following structures are available:

- The **sky\_all\_stats\_t** structure is available for the PCI Token-Ring High Performance Device Driver (14101800), and is defined in the device-specific `/usr/include/sys/cdli_tokuser.h` include file.
- The **cs\_all\_stats\_t** structure is available for the PCI Token-Ring Device Driver (14103e00), and is defined in the device-specific `/usr/include/sys/cdli_tokuser.cstok.h` include file.

The statistics that are returned contain information obtained from the device. If the device is inoperable, the statistics returned are not the current device statistics. The copy of the `ndd_flags` field can be checked to determine the state of the device.

# **Reliability, Availability, and Serviceability (RAS)**

## **Trace**

For LAN device drivers, trace points enable error monitoring as well as tracking packets as they move through the driver. The drivers issue trace points for some or all of the following conditions:

- Beginning and ending of main functions in the main path
- Error conditions

- Beginning and ending of each function that is tracking buffers outside of the main path
- Debugging purposes (These trace points are only enabled when the driver is compiled with the **-DDEBUG** option turned, therefore, the driver can contain as many of these trace points as needed.)

Following is a list of trace hooks and location of definition files for the existing ethernet device drivers.

**The PCI Token-Ring High Performance Device Driver (14101800): Definition File:**  
**/sys/cdli\_tokuser.h**

#### Trace Hook IDs

- Transmit 2A7
- Receive 2A8
- Error 2A9
- Other 2AA

**The PCI Token-Ring (14103e00) Device Driver: Definition File:** /sys/cdli\_tokuser.cstok.h

#### Trace Hook IDs

- Transmit 2DA
- Receive 2DB
- General 2DC

### Error Logging

**PCI Token-Ring High Performance Device Driver (14101800):** The error IDs for the PCI Token-Ring High Performance Device Driver (14101800) are as follows:

#### **ERRID\_STOK\_ADAP\_CHECK**

The microcode on the device performs a series of diagnostic checks when the device is idle. These checks can find errors, and they are reported as adapter checks. If the device is connected to the network when this error occurs, the device driver goes into Network Recovery Mode in an attempt to recover from the error. The device is temporarily unavailable during the recovery procedure. User intervention is not required for this error unless the problem persists.

#### **ERRID\_STOK\_ADAP\_OPEN**

Enables the device driver to open the device. The device driver goes into Network Recovery Mode in an attempt to recover from the error. The device is temporarily unavailable during the recovery procedure. User intervention is not required for this error unless the problem persists.

#### **ERRID\_STOK\_AUTO\_RMV**

An internal hardware error following the beacon automatic removal process was detected. The device driver goes into Network Recovery Mode in an attempt to recover from the error. The device is temporarily unavailable during the recovery procedure. User intervention is not required for this error unless the problem persists.

#### **ERRID\_STOK\_RING\_SPEED**

The ring speed (or ring data rate) is probably wrong. Contact the network administrator to determine the speed of the ring. The device driver only retries twice at 2-minute intervals after this error log entry is generated.

#### **ERRID\_STOK\_DMAFAIL**

The device detected a DMA error in a TX or RX operation. The device driver goes into Network Recovery Mode in an attempt to recover from the error. The device is temporarily unavailable during the recovery procedure. User intervention is not required unless the problem persists.

#### **ERRID\_STOK\_BUS\_ERR**

The device detected a Micro Channel bus error. The device driver goes into Network Recovery

Mode in an attempt to recover from the error. The device is temporarily unavailable during the recovery procedure. User intervention is not required for this error unless the problem persists.

**Note:** Micro Channel is only supported on AIX 5.1 and earlier.

**ERRID\_STOK\_DUP\_ADDR**

The device detected that another station on the ring has a device address that is the same as the device address being tested. Contact the network administrator to determine why.

**ERRID\_STOK\_MEM\_ERR**

An error occurred while allocating memory or timer control block structures.

**ERRID\_STOK\_RCVRY\_EXIT**

The error that caused the device driver to go into error recovery mode was corrected.

**ERRID\_STOK\_RMV\_ADAP**

The device received a remove ring station MAC frame indicating that a network management function directed this device to get off the ring. Contact the network administrator to determine why.

**ERRID\_STOK\_WIRE\_FAULT**

There is a loose (or bad) cable between the device and the MAU. There is a chance that it might be a bad device. The device driver goes into Network Recover Mode in an attempt to recover from the error. The device is temporarily unavailable during the recovery procedure. User intervention is not required for this error unless the problem persists.

**ERRID\_STOK\_TX\_TIMEOUT**

The transmit watchdog timer expired before transmitting a frame. The device driver goes into Network Recovery Mode in an attempt to recover from the error. The device is temporarily unavailable during the recovery procedure. User intervention is not required for this error unless the problem persists.

**ERRID\_STOK\_CTL\_ERR**

The ioctl watchdog timer expired before the device driver received a response from the device. The device driver goes into Network Recovery Mode in an attempt to recover from the error. The device is temporarily unavailable during the recovery procedure. User intervention is not required for this error unless the problem persists.

**PCI Token-Ring Device Driver (14103e00):** The error IDs for the PCI Token-Ring Device Driver (14103e00) are as follows:

**ERRID\_CSTOK\_ADAP\_CHECK**

The microcode on the device performs a series of diagnostic checks when the device is idle on initialization. These checks find errors and they are reported as adapter checks. If the device was connected to the network when this error occurred, the device driver will go into Network Recovery Mode in an attempt to recover from the error. The device is temporarily unavailable during the recovery procedure. After this error log entry has been generated, the device driver will retry 3 times with no delay between retries. User intervention is not required for this error unless the problem persists.

**ERRID\_CSTOK\_ADAP\_OPEN**

The device driver was unable to open the device. The device driver will go into Network Recovery Mode in an attempt to recover from this error. The device is temporarily unavailable during the recovery procedure. The device driver will retry indefinitely with a 30 second delay between retries to recover. User intervention is not required for this error unless the problem persists.

**ERRID\_CSTOK\_AUTO\_RMV**

An internal hardware error following the beacon automatic removal process has been detected. The device driver will go into Network Recovery Mode in an attempt to recover from the error. The device is temporarily unavailable during the recovery procedure. User intervention is not required for this error unless the problem persists.

**ERRID\_CSTOK\_RING\_SPEED**

The ring speed or ring data rate is probably wrong. Contact the network administrator to determine the speed of the ring. The device driver will only retry twice at 2 minute intervals after this error log entry has been generated.

**ERRID\_CSTOK\_DMAFAIL**

The device detected a DMA error in a TX or RX operation. The device driver will go into Network Recovery Mode in an attempt to recover from this error. The device is temporarily unavailable during the recovery procedure. User intervention is not required for this error unless the problem persists.

**ERRID\_CSTOK\_BUS\_ERR**

The device detected a PCI bus error. The device driver will go into Network Recovery Mode in an attempt to recover from this error. The device is temporarily unavailable during the recovery procedure. User intervention is not required for this error unless the problem persists.

**ERRID\_CSTOK\_DUP\_ADDR**

The device has detected that another station on the ring has a device address which is the same as the device address being tested. Contact network administrator to determine why.

**ERRID\_CSTOK\_MEM\_ERR**

An error occurred while allocating memory or timer control block structures. This usually implies the system has run out of available memory. User intervention is required.

**ERRID\_CSTOK\_RCVRY\_ENTER**

An error has occurred which caused the device driver to go into network recovery.

**ERRID\_CSTOK\_RCVRY\_EXIT**

The error which caused the device driver to go into Network Recovery Mode has been corrected.

**ERRID\_CSTOK\_RMV\_ADAP**

The device has received a remove ring station MAC frame indicating that a network management function has directed this device to get off the ring. The device driver will only retry twice with 6 minute delay between retries after this error log entry has been generated. Contact network administrator to determine why.

**ERRID\_CSTOK\_WIRE\_FAULT**

There is probably a loose ( or bad ) cable between the device and the MAU. There is some chance that it might be a bad device. The device driver will go into Network Recovery Mode in an attempt to recover from this error. The device is temporarily unavailable during the recovery procedure. User intervention is not required for this error unless the problem persists.

**ERRID\_CSTOK\_RX\_ERR**

The device has detected a receive error. The device driver will go into Network Recovery Mode in an attempt to recover from this error. The device is temporarily unavailable during the recovery procedure. User intervention is not required for this error unless the problem persists.

**ERRID\_CSTOK\_TX\_ERR**

The device has detected a transmit error. The device driver will go into Network Recovery Mode in an attempt to recover from this error. The device is temporarily unavailable during the recovery procedure. User intervention is not required for this error unless the problem persists.

**ERRID\_CSTOK\_TX\_TMOUT**

The transmit watchdog timer has expired before the transmit of a frame has completed. The device driver will go into Network Recovery Mode in an attempt to recover from this error. The device is temporarily unavailable during the recovery procedure. User intervention is not required for this error unless the problem persists.

**ERRID\_CSTOK\_CMD\_TMOUT**

The ioctl watchdog timer has expired before the device driver received a response from the device. The device driver will go into Network Recovery Mode in an attempt to recover from this

error. The device is temporarily unavailable during the recovery procedure. User intervention is not required for this error unless the problem persists.

#### **ERRID\_CSTOK\_PIO\_ERR**

The driver has encountered a PIO operation error. The device driver will attempt to retry the operation 3 times before it will fail the command and return in the DEAD state to the user. User intervention is required.

#### **ERRID\_CSTOK\_PERM\_HW**

The microcode on the device performs a series of diagnostic checks on initialization. These checks can find errors and they are reported as adapter checks. If the error occurs 4 times during adapter initialization this error log will be generated and the device considered inoperable. User intervention is required.

#### **ERRID\_CSTOK\_ASB\_ERR**

The adapter has indicated that the processing of a TokenRing mac command failed.

#### **ERRID\_CSTOK\_AUTO\_FAIL**

The ring speed of the adapter is set to autosense, and open has failed because this adapter is the only one on the ring. User intervention is required.

#### **ERRID\_CSTOK\_EISR**

If the adapter detects a PCI Master or Target Abort, the Error Interrupt Status Register (EISR) will be set.

#### **ERRID\_CSTOK\_CMD\_ERR**

Adapter failed command due to a transient error and goes into limbo one time, if that fails the adapter goes into the dead state.

#### **ERRID\_CSTOK\_EEH\_ENTER**

The adapter encountered a Bus I/O Error, and is attempting to recover by using the EEH recovery process.

#### **ERRID\_CSTOK\_EEH\_EXIT**

The adapter successfully recovered from the I/O Error by using the EEH recovery process.

#### **ERRID\_CSTOK\_EEH\_HW\_ERR**

The adapter could not recover from the EEH error. The EEH error was the result of an adapter error, and not a bus error (logged by the kernel).

---

## **Ethernet Device Drivers**

The following Ethernet device drivers are dynamically loadable. The device drivers are automatically loaded into the system at device configuration time as part of the configuration process.

- PCI Ethernet Adapter Device Driver (22100020)
- 10/100Mbps Ethernet PCI Adapter Device Driver (23100020)
- 10/100Mbps Ethernet PCI Adapter II Device Driver (1410ff01)
- Gigabit Ethernet-SX PCI Adapter Device Driver (14100401)
- Gigabit Ethernet-SX PCI-X Adapter Device Driver (14106802)
- 10/100/1000 Base-T Ethernet PCI-X Adapter Device Driver (14106902)
- 2-Port Gigabit Ethernet-SX PCI-X Adapter (14108802)
- 2-Port 10/100/1000 Base-TX PCI-X Adapter (14108902)
- 10 Gigabit Ethernet-SR PCI-X Adapter (1410ba02)
- 10 Gigabit Ethernet-LR PCI\_X Adapter (1410bb02)
- 10 Gigabit Ethernet-SR PCI-X 2.0 DDR Adapter (1410eb02)
- 10 Gigabit Ethernet-LR PCI-X 2.0 DDR Adapter (1410ec02)
- Gigabit Ethernet-SX Adapter (e414a816)

- Gigabit Ethernet-SX Adapter (14101403)
- Gigabit Ethernet-SX Adapter (14106703)
- 4-Port 10/100/1000 Base-TX PCI-X Adapter (14101103)
- 4-Port 10/100/1000 Base-TX PCI-Express Adapter (14106803)
- 2-Port Gigabit Ethernet-SX PCI-Express Adapter (14103f03)
- 2-Port 10/100/1000 Base-TX PCI-Express Adapter (14104003)
- Host Ethernet Adapter Device Driver

The following information is provided about each of the Ethernet device drivers:

- Configuration Parameters
- Interface Entry Points
- Asynchronous Status
- Device Control Operations
- Trace
- Error Logging

For each Ethernet device, the interface to the device driver is achieved by calling the entry points for opening, closing, transmitting data, and issuing device control commands.

There are a number of Ethernet device drivers in use. All drivers provide PCI-based connections to an Ethernet network, and support both Standard and IEEE 802.3 Ethernet Protocols.

The PCI Ethernet Adapter Device Driver (22100020) supports the PCI Ethernet BNC/RJ-45 Adapter (feature 2985) and the PCI Ethernet BNC/AUI Adapter (feature 2987), as well as the integrated Ethernet port on certain systems.

The 10/100 Mbps Ethernet PCI Adapter Device Driver (23100020) supports the 10/100 Mbps Ethernet PCI Adapter (feature 2968) and the Four Port 10/100 Mbps Ethernet PCI Adapter (features 4951 and 4961), as well as the integrated Ethernet port on certain systems.

The 10/100 Mbps Ethernet PCI Adapter II Device Driver (1410ff01) supports the 10/100 Mbps Ethernet PCI Adapter II (feature 4962), as well as the integrated Ethernet port on certain systems.

The Gigabit Ethernet-SX PCI Adapter Device Driver (14100401) supports the Gigabit Ethernet-SX PCI Adapter (feature 2969) and the 10/100/1000 Base-T Ethernet Adapter (feature 2975).

The Gigabit Ethernet-SX PCI-X Adapter Device Driver (14106802) supports the Gigabit Ethernet-SX PCI-X Adapter (feature 5700).

The 10/100/1000 Base-TX Ethernet PCI-X Adapter Device Driver (14106902) supports the 10/100/1000 Base-TX Ethernet PCI-X Adapter (feature 5701).

The 2-Port Gigabit Ethernet-SX PCI-X Adapter Device Driver (14108802) supports the 2-Port Gigabit Ethernet-SX PCI-X Adapter (feature 5707).

The 2-Port 10/100/1000 Base-TX PCI-X Adapter Device Driver (14108902) supports the 2-Port 10/100/1000 Base-TX PCI-X Adapter (feature 5706).

The 10 Gigabit Ethernet-SR PCI-X Adapter Device Driver (1410ba02) supports the 10 Gigabit Ethernet-SR PCI-X Adapter (feature 5718).

The 10 Gigabit Ethernet-LR PCI-X Adapter Device Driver (1410bb02) supports the 10 Gigabit Ethernet-LR PCI-X Adapter (feature 5719).

The 10 Gigabit Ethernet-SR PCI-X 2.0 DDR Adapter Device Driver supports the 10 Gigabit Ethernet-SR PCI-X 2.0 DDR Adapter (feature 5721).

The 10 Gigabit Ethernet-LR PCI-X 2.0 DDR Adapter Device Driver supports the 10 Gigabit Ethernet-LR PCI-X 2.0 DDR Adapter (feature 5722).

The Gigabit Ethernet-SX Adapter Device Driver (e414a816) supports the eServer BladeCenter® JS20 Gigabit Ethernet-SX Adapter.

The Gigabit Ethernet-SX Adapter Device Driver (14101403) supports the eServer BladeCenter JS21 Gigabit Ethernet-SX Adapter.

The Gigabit Ethernet-SX Adapter Device Driver (14106703) supports the eServer BladeCenter Multiple Switch Interface Module Gigabit Ethernet-SX Adapter.

The 4-Port 10/100/1000 Base-TX Ethernet PCI-X Adapter Device Driver (14101103) supports the 4-Port 10/100/1000 Base-TX PCI-X Adapter (feature 5740).

The 4-Port 10/100/1000 Base-TX PCI-Express Adapter Device Driver (14106803) supports the 4-Port 10/100/1000 Ethernet PCI-E Adapter (feature 5717).

The 2-Port Gigabit Ethernet-SX PCI-Express Adapter Device Driver (14103f03) supports the 2-Port Gigabit Ethernet-SX PCI-Express Adapter (feature 5768).

The 2-Port 10/100/1000 Base-TX Ethernet PCI-Express Adapter Device Driver (14104003) supports the 2-Port 10/100/1000 Base-TX PCI-Express Adapter (feature 5767).

The Host Ethernet Adapter Device Driver supports the Host Ethernet Adapter (feature 5636, 5637, and 5639).

## Configuration Parameters

The following configuration parameter is supported by all Ethernet device drivers:

### Alternate Ethernet Addresses

The device drivers support the device's hardware address as the network address or an alternate network address configured through software. When an alternate address is used, any valid Individual Address can be used. The least significant bit of an Individual Address must be set to zero. A multicast address cannot be defined as a network address. Two configuration parameters are provided to provide the alternate Ethernet address and enable the alternate address.

### PCI Ethernet Device Driver (22100020)

The PCI Ethernet Device Driver (22100020) supports the following additional configuration parameters:

#### Full Duplex

Indicates whether the adapter is operating in full-duplex or half-duplex mode. If this field is set to yes, the device driver programs the adapter to be in full-duplex mode.

#### Hardware Transmit Queue

Specifies the actual queue size the adapter uses to transmit packets. Each element corresponds to an Ethernet packet. It is configurable at 16, 32, 64, 128, and 256 elements.

#### Hardware Receive Queue

Specifies the actual queue size the adapter uses to receive packets. Each element corresponds to an Ethernet packet. It is configurable at 16, 32, 64, 128, and 256 elements.

### 10/100 Mbps Ethernet PCI Adapter Device Driver (23100020)

The 10/100 Mbps Ethernet PCI Adapter Device Driver (23100020) supports the following additional configuration parameters:

### Software Transmit Queue

Indicates the number of transmit requests that can be queued for transmission by the device driver. Valid values range from 16 through 16 384.

### Hardware Receive Queue

The 10/100 Mbps Ethernet PCI Adapter Device Driver (23100020) supports a user-configurable receive queue for the adapter. This is the actual queue the adapter uses to receive packets. Each element corresponds to an Ethernet packet. It is configurable at 16, 32, 64, 128, and 256 elements.

### Receive Buffer Pool

The 10/100 Mbps Ethernet PCI Adapter Device Driver (23100020) implements a private pool of receive memory buffers in order to enhance driver performance. The number of private receive buffers reserved by the driver is configurable from 16 to 2048 elements.

### Media Speed

The 10/100 Mbps Ethernet PCI Adapter Device Driver (23100020) supports a user-configurable media speed for the adapter. The media speed attribute indicates the speed at which the adapter attempts to operate. The available speeds are 10 Mbps half-duplex, 10 Mbps full-duplex, 100 Mbps half-duplex, 100 Mbps full-duplex and autonegotiation, with a default of autonegotiation. Select autonegotiate when the adapter should use autonegotiation across the network to determine the speed. When the network does not support autonegotiation, select the specific speed.

**Note:** If autonegotiation is selected, the remote link device must also be set to autonegotiate or the link might not function properly.

### Inter Packet Gap

The 10/100 Mbps Ethernet PCI Adapter Device Driver (23100020) supports a user-configurable inter packet gap for the adapter. The inter packet gap attribute controls the aggressiveness of the adapter on the network. A small number increases the aggressiveness of the adapter, but a large number decreases the aggressiveness (and increase the fairness) of the adapter. A small number (more aggressive) could cause the adapter to capture the network by forcing other less aggressive nodes to defer. A larger number (less aggressive) might cause the adapter to defer more often than normal. If the statistics for other nodes on the network show a large number of collisions and deferrals, then try increasing this number. The default is 96, which results in IPG of 9.6 micro seconds for 10 Mbps and 0.96 microseconds for 100 Mbps media speed. Each unit of bit rate introduces an IPG of 100 nsec at 10 Mbps, and 10 nsec at 100 Mbps media speed.

### Link Polling Timer

The 10/100 Mbps Ethernet PCI Adapter Device Driver (23100020) implements a polling function (**Enable Link Polling**) that periodically queries the adapter to determine whether the Ethernet link is up or down. The **Enable Link Polling** attribute is disabled by default. If this function is enabled, the link polling timer value indicates how often the driver should poll the adapter for link status. This value can range from 100 to 1000 milliseconds. If the adapter's link goes down, the device driver disables its **NDD\_RUNNING** flag. When the device driver finds that the link has come back up, it enables this **NDD\_RUNNING** flag. In order for this to work successfully, protocol layer implementations, such as Etherchannel, need notification if the link has gone down. Enable the **Enable Link Polling** attribute to obtain this notification. Because of the additional PIO calls that the device driver makes, enabling this attribute can decrease the performance of this adapter.

## 10/100 Mbps Ethernet PCI Adapter II Device Driver (1410ff01)

The 10/100 Mbps Ethernet PCI Adapter II Device Driver (1410ff01) supports the following additional configuration parameters:

### Software Transmit Queue

Indicates the number of transmit requests that can be queued for transmission by the device driver. Valid values range from 512 through 16 384.



### **Hardware Transmit Queue**

The 10/100 Mbps Ethernet PCI Adapter II Device Driver (1410ff01) supports a user-configurable transmit queue for the adapter. This is the actual queue the adapter uses to transmit packets. Each element corresponds to an Ethernet packet. It is configurable from 100 to 1024 elements.

### **Hardware Receive Queue**

The 10/100 Mbps Ethernet PCI Adapter II Device Driver (1410ff01) supports a user-configurable receive queue for the adapter. This is the actual queue the adapter uses to receive packets. Each element corresponds to an Ethernet packet. It is configurable from 100 to 1024 elements.

### **Receive Buffer Pool**

The 10/100 Mbps Ethernet PCI Adapter II Device Driver (1410ff01) implements a private pool of receive memory buffers in order to enhance driver performance. The number of private receive buffers reserved by the driver is configurable from 512 to 2048 elements.

### **Media Speed**

The 10/100 Mbps Ethernet PCI Adapter II Device Driver (1410ff01) supports a user-configurable media speed for the adapter. The media speed attribute indicates the speed at which the adapter attempts to operate. The available speeds are 10 Mbps half-duplex, 10 Mbps full-duplex, 100 Mbps half-duplex, 100 Mbps full-duplex and autonegotiation, with a default of autonegotiation. Select autonegotiate when the adapter should use autonegotiation across the network to determine the speed. When the network does not support autonegotiation, select the specific speed.

**Note:** If autonegotiation is selected, the remote link device must also be set to autonegotiate or the link might not function properly.

### **Link Polling Timer**

The 10/100 Mbps Ethernet PCI Adapter II Device Driver (1410ff01) implements a polling function which periodically queries the adapter to determine whether the Ethernet link is up or down. If this function is enabled, the link polling timer value indicates how often the driver should poll the adapter for link status. This value can range from 100 to 1000 milliseconds.

### **Checksum Offload**

The 10/100 Mbps Ethernet PCI Adapter II Device Driver (1410ff01) supports the capability of the adapter to calculate TCP checksums in hardware. If this capability is enabled, the TCP checksum calculation is performed on the adapter instead of the host, which can increase system performance. Possible values are Yes and No.

### **Transmit TCP Resegmentation Offload**

Permits the adapter to perform resegmentation of transmitted TCP segments in hardware. With this capability, the host can use TCP segments that are larger than the actual MTU size of the Ethernet link, which can increase system performance. You can specify the **Yes** and **No** values.

### **IPsec Offload**

The 10/100 Mbps Ethernet PCI Adapter II Device Driver (1410ff01) supports the capability of the adapter to perform IPsec cryptographic algorithms for data encryption and authentication in hardware. This capability enables the host to offload processor-intensive cryptographic processing to the adapter, which can increase system performance. Possible values are Yes and No.

## **Gigabit Ethernet-SX PCI Adapter Device Driver (14100401)**

The Gigabit Ethernet-SX PCI Adapter Device Driver (14100401) supports the following additional configuration parameters:

### **Software Transmit Queue Size**

Indicates the number of transmit requests that can be queued for transmission by the device driver. Valid values range from 512 through 16 384.

### **Transmit Jumbo Frames**

When you set the attribute value to **Yes**, frames up to 9018 bytes in length can be transmitted on

this adapter. When you set the attribute value to **No**, the maximum size of frames transmitted is 1518 bytes. Frames up to 9018 bytes in length can always be received on this adapter.

### **Enable Hardware Checksum Offload**

When you set the attribute value to **Yes**, the adapter calculates the checksum for transmitted and received TCP frames. When you set the attribute value to **No**, the checksum is calculated by appropriate software.

**Note:** The **mbuf** structure, which describes a transmitted frame, contains a flag that indicates whether the adapter calculates the checksum for the frame.

### **Media Speed**

The Gigabit Ethernet-SX PCI Adapter Device Driver (14100401) supports a user-configurable media speed only for the IBM 10/100/1000 Base-T Ethernet PCI adapter (feature 2975). For the Gigabit Ethernet-SX PCI Adapter (feature 2969), the only possible choice is autonegotiation. The media speed attribute indicates the speed at which the adapter attempts to operate. The available speeds are 10 Mbps half-duplex, 10 Mbps full-duplex, 100 Mbps half-duplex, 100 Mbps full-duplex and autonegotiation, with a default of autonegotiation. Select autonegotiate when the adapter should use autonegotiation across the network to determine the speed. When the network does not support autonegotiation, select the specific speed.

**Note:** The autonegotiation setting must be selected in order for the adapter to run at 1000 M-bit/s.

### **Enable Hardware Transmit TCP Resegmentation**

When you set the attribute value to **Yes**, the adapter performs TCP resegmentation on transmitted TCP segments. With this capability, TCP/IP can send larger datagrams to the adapter, which can increase performance. When you set the attribute value to **No**, TCP resegmentation is not performed.

**Note:** The default values for the Gigabit Ethernet-SX PCI Adapter Device Driver (14100401) configuration parameters were chosen for optimal performance. Do not change these default values.

The following configuration parameters for the Gigabit Ethernet-SX PCI Adapter Device Driver (14100401) are not accessible using SMIT. You can modify them only using the **chdev** command.

#### **stat\_ticks**

Indicates the number of microseconds that the adapter waits before updating the adapter statistics (through a DMA write) and generating an interrupt to the host. Valid values range from 1000 through 1 000 000. The default value is 1 000 000.

#### **receive\_ticks**

Indicates the number of microseconds that the adapter waits before updating the producer index of the receive return ring (through a DMA write) and generating an interrupt to the host. Valid values range from 0 through 1000, the default value is 50.

#### **receive\_bds**

Indicates the number of receive buffers that the adapter transfers to host memory before updating the producer index of the receive return ring (through a DMA write) and generating an interrupt to the host. Valid values range from 0 through 128. The default value is 6.

#### **tx\_done\_ticks**

Indicates the number of microseconds that the adapter waits before updating the send consumer index (through a DMA write) and generating an interrupt to the host. Valid values range from 0 through 1 000 000. The default value is 1 000 000.

#### **tx\_done\_count**

Indicates the number of transmit buffers that the adapter transfers from host memory before updating the send consumer index (through a DMA write) and generating an interrupt to the host. Valid values range from 0 through 128. The default value is 64.

**receive\_proc**

When this number of receive buffer descriptors is processed by the device driver (or all packets are received), the device driver returns this number of receive buffer descriptors to the adapter through an MMIO write. Valid values range from 1 through 64. The default value is 16.

**rxdesc\_count**

When this number of receive buffer descriptors is processed by the device driver (or all packets were received), the device driver exits the **rx\_handler()** routine and continues processing other adapter events, such as transmit completions and adapter status changes. Valid values range from 1 through 1 000 000. The default value is 1000.

**slih\_hog**

Indicates the number of adapter events (such as receive completions, transmit completions and adapter status changes) that are processed by the device driver per interrupt. Valid values range from 1 through 1 000 000. The default value is 10.

**copy\_bytes**

When the number of data bytes in a transmit **mbuf** structure exceeds this value, the device driver maps the **mbuf** data area into DMA memory and updates the transmit descriptor so that it points to this DMA memory area. When the number of data bytes in a transmit **mbuf** structure does not exceed this value, the data is copied from the **mbuf** structure into a preallocated transmit buffer that is already mapped into DMA memory. The device driver also attempts to coalesce transmit data in an **mbuf** chain into a single preallocated transmit buffer until the total transmit data size exceeds that of the preallocated buffer (2048 bytes). Valid values range from 64 through 2048. The default value is 2048.

**Gigabit Ethernet-SX PCI-X Adapter Device Driver (14106802)**

The Gigabit Ethernet-SX PCI-X Adapter Device Driver (14106802) supports the following additional configuration parameters:

**Transmit Descriptor Queue Size**

Indicates the number of transmit requests that can be queued for transmission by the adapter. Valid values range from 128 through 1024.

**Receive Descriptor Queue Size**

Indicates the maximum number of received Ethernet packets that the adapter can hold in its buffer. Valid values range from 128 through 1024.

**Software Transmit Queue Size**

Indicates the number of transmit requests that can be queued for transmission by the device driver. Valid values range from 512 through 16 384.

**Media Speed**

Indicates the speed at which the adapter attempts to operate. The available speeds are 1000 Mbps full-duplex and autonegotiation. The default is autonegotiation. Select autonegotiate when the adapter should use autonegotiation across the network to determine the duplexity. When the network does not support autonegotiation, select 1000 Mbps full-duplex.

**Transmit Jumbo Frames**

When you set the attribute value to **Yes**, frames up to 9018 bytes in length can be transmitted on this adapter. When you set the attribute value to **No**, the maximum size of frames transmitted is 1518 bytes. Frames up to 9018 bytes in length can always be received on this adapter.

**Enable Hardware Transmit TCP Resegmentation**

Permits the adapter to perform resegmentation of transmitted TCP segments in hardware. With this capability, the host can use TCP segments that are larger than the actual MTU size of the Ethernet link, which can increase system performance. You can specify the values of **Yes** and **No**.

**Note:** The **mbuf** structure, which describes a transmitted frame, contains a flag that indicates whether the adapter performs TCP resegmentation for the frame.

### Enable Hardware Transmit and Receive Checksum

When you set the attribute value to **Yes**, the adapter calculates the checksum for transmitted and received TCP frames. When you set the attribute value to **No**, the checksum is calculated by appropriate software.

**Note:** The **mbuf** structure, which describes a transmitted frame, contains a flag that indicates whether the adapter calculates the TCP checksum for the frame.

The following configuration parameters for the Gigabit Ethernet-SX PCI-X Adapter Device Driver (14106802) are not accessible using SMIT. You can modify them only using the **chdev** command.

#### **rx\_hog**

When this number of receive buffer descriptors is processed by the device driver (or all packets were received), the device driver exits the **rx\_handler()** routine and continues processing other adapter events, such as transmit completions and adapter status changes. Valid values range from 1 through 1 000 000. The default value is 1000.

#### **slih\_hog**

Indicates the number of adapter events (such as receive completions, transmit completions, and adapter status changes) that are processed by the device driver per interrupt. Valid values range from 1 through 1 000 000. The default value is 10.

#### **copy\_bytes**

When the number of data bytes in a transmit **mbuf** structure exceeds this value, the device driver maps the **mbuf** data area into DMA memory and updates the transmit descriptor so that it points to this DMA memory area. When the number of data bytes in a transmit **mbuf** structure does not exceed this value, the data is copied from the **mbuf** structure into a preallocated transmit buffer that is already mapped into DMA memory. The device driver also attempts to coalesce transmit data in an **mbuf** chain into a single preallocated transmit buffer until the total transmit data size exceeds that of the preallocated buffer (2048 bytes). Valid values range from 64 through 2048. The default value is 2048.

#### **delay\_open**

When you set the attribute value to **Yes**, the adapter device driver delays its open completion until the Ethernet link status is determined to be either up or down. This prevents applications from sending data before the Ethernet link is established. Commands (for example, the **ifconfig** command) however might take longer to complete, especially when an active Ethernet link is not present. You can specify the values of **Yes** and **No**. The default value is **No**.

### 10/100/1000 Base-T Ethernet PCI-X Adapter Device Driver (14106902)

The 10/100/1000 Base-T Ethernet PCI-X Adapter Device Driver (14106902) supports the following additional configuration parameters:

#### **Transmit Descriptor Queue Size**

Indicates the number of transmit requests that can be queued for transmission by the adapter. Valid values range from 128 through 1024.

#### **Receive Descriptor Queue Size**

Indicates the maximum number of received Ethernet packets that the adapter can buffer. Valid values range from 128 through 1024.

#### **Software Transmit Queue Size**

Indicates the number of transmit requests that can be queued for transmission by the device driver. Valid values range from 512 through 16 384.

#### **Media Speed**

Indicates the speed at which the adapter attempts to operate. The available speeds are 10 Mbps half-duplex, 10 Mbps full-duplex, 100 Mbps half-duplex, 100 Mbps full-duplex, 1000 Mbps full-duplex, and autonegotiation, with a default of autonegotiation. Select autonegotiate when the

adapter should use autonegotiation across the network to determine the speed. When the network does not support autonegotiation, select the specific speed.

### **Transmit Jumbo Frames**

When you set the attribute value to **Yes**, frames up to 9018 bytes in length can be transmitted on this adapter. When you set the attribute value to **No**, the maximum size of frames transmitted is 1518 bytes. Frames up to 9018 bytes in length can always be received on this adapter.

### **Enable Hardware Transmit TCP Resegmentation**

Permits the adapter to perform resegmentation of transmitted TCP segments in hardware. With this capability, the host can use TCP segments that are larger than the actual MTU size of the Ethernet link, which can increase system performance. You can specify the values of **Yes** and **No**.

**Note:** The **mbuf** structure, which describes a transmitted frame, contains a flag that indicates whether the adapter performs TCP resegmentation for the frame.

### **Enable Hardware Transmit and Receive Checksum**

When you set the attribute value to **Yes**, the adapter calculates the checksum for transmitted and received TCP frames. When you set the attribute value to **No**, the checksum is calculated by appropriate software.

**Note:** The **mbuf** structure, which describes a transmitted frame, contains a flag that indicates whether the adapter calculates the TCP checksum for the frame.

The following configuration parameters for the 10/100/1000 Base-T Ethernet PCI-X Adapter Device Driver (14106902) are not accessible using SMIT. You can modify them only using the **chdev** command.

#### **rx\_hog**

When this number of receive buffer descriptors is processed by the device driver (or all packets were received), the device driver exits the **rx\_handler()** routine and continues processing other adapter events, such as transmit completions and adapter status changes. Valid values range from 1 through 1 000 000. The default value is 1000.

#### **slih\_hog**

Indicates the number of adapter events (such as receive completions, transmit completions, and adapter status changes) that are processed by the device driver per interrupt. Valid values range from 1 through 1 000 000. The default value is 10.

#### **copy\_bytes**

When the number of data bytes in a transmit **mbuf** structure exceeds this value, the device driver maps the **mbuf** data area into DMA memory and updates the transmit descriptor so that it points to this DMA memory area. When the number of data bytes in a transmit **mbuf** structure does not exceed this value, the data is copied from the **mbuf** structure into a preallocated transmit buffer that is already mapped into DMA memory. The device driver also attempts to coalesce transmit data in an **mbuf** chain into a single preallocated transmit buffer until the total transmit data size exceeds that of the preallocated buffer (2048 bytes). Valid values range from 64 through 2048. The default value is 2048.

#### **delay\_open**

When you set the attribute value to **Yes**, the adapter device driver delays its open completion until the Ethernet link status is determined to be either up or down. This prevents applications from sending data before the Ethernet link is established. Commands (for example, the **ifconfig** command) however might take longer to complete, especially when an active Ethernet link is not present. You can specify the values of **Yes** and **No**. The default value is **No**.

#### **compat\_mode**

Setting this attribute to Yes forces the adapter to implement an early version of the IEEE 802.3z autonegotiation protocol. Use the yes value only if the adapter is unable to establish a link with your older Gigabit Ethernet-TX adapters or switches. Valid values are yes and No. The default value is No.

**Note:** If this option is enabled, the adapter cannot establish a link with newer Gigabit Ethernet-TX hardware. Enable this option only if you cannot establish a link using autonegotiation, but can force a link at a slower speed (for example, 100 full-duplex).

## 2-Port Gigabit Ethernet-SX PCI-X Adapter Device Driver (14108802)

The 2-Port Gigabit Ethernet-SX PCI-X Adapter Device Driver (14108802) supports the following additional configuration parameters:

### Transmit Descriptor Queue Size

Indicates the number of transmit requests that can be queued for transmission by the adapter. Valid values range from 128 through 1024.

### Receive Descriptor Queue Size

Indicates the maximum number of received Ethernet packets that the adapter can hold in its buffer. Valid values range from 128 through 1024.

### Software Transmit Queue Size

Indicates the number of transmit requests that can be queued for transmission by the device driver. Valid values range from 512 through 16 384.

### Media Speed

Indicates the speed at which the adapter attempts to operate. The available speeds are 1000 Mbps full-duplex and autonegotiation. The default is autonegotiation. Select autonegotiate when the adapter should use autonegotiation across the network to determine the duplexity. When the network does not support autonegotiation, select 1000 Mbps full-duplex.

### Transmit Jumbo Frames

When you set the attribute value to **Yes**, frames up to 9018 bytes in length can be transmitted on this adapter. When you set the attribute value to **No**, the maximum size of frames transmitted is 1518 bytes. Frames up to 9018 bytes in length can always be received on this adapter.

### Enable Hardware Transmit TCP Resegmentation

Permits the adapter to perform resegmentation of transmitted TCP segments in hardware. With this capability, the host can use TCP segments that are larger than the actual MTU size of the Ethernet link, which can increase system performance. You can specify the values of **Yes** and **No**.

**Note:** The **mbuf** structure, which describes a transmitted frame, contains a flag that indicates whether the adapter performs TCP resegmentation for the frame.

### Enable Hardware Transmit and Receive Checksum

When you set the attribute value to **Yes**, the adapter calculates the checksum for transmitted and received TCP frames. When you set the attribute value to **No**, the checksum is calculated by appropriate software.

**Note:** The **mbuf** structure, which describes a transmitted frame, contains a flag that indicates whether the adapter calculates the TCP checksum for the frame.

### Enable Failover Mode

Indicates the requested failover configuration for the port. You can specify the values of **primary**, **backup**, and **disable**.

<b>primary</b>	Indicates the port is to act as the primary port in a failover configuration for a 2-port gigabit adapter.
<b>backup</b>	Indicates the port is to act as the backup port in a failover configuration for a 2-port gigabit adapter.
<b>disable</b>	Indicates the port is not a member of a failover configuration. This is the default value for failover.

The following configuration parameters for the 2-Port Gigabit Ethernet-SX PCI-X Adapter Device Driver (14108802) are not accessible using SMIT. You can modify them only using the **chdev** command.

### rx\_hog

When this number of receive buffer descriptors is processed by the device driver (or all packets

were received), the device driver exits the **rx\_handler()** routine and continues processing other adapter events, such as transmit completions and adapter status changes. Valid values range from 1 through 1 000 000. The default value is 1000.

#### **slih\_hog**

Indicates the number of adapter events (such as receive completions, transmit completions, and adapter status changes) that are processed by the device driver per interrupt. Valid values range from 1 through 1 000 000. The default value is 10.

#### **copy\_bytes**

When the number of data bytes in a transmit **mbuf** structure exceeds this value, the device driver maps the **mbuf** data area into DMA memory and updates the transmit descriptor so that it points to this DMA memory area. When the number of data bytes in a transmit **mbuf** structure does not exceed this value, the data is copied from the **mbuf** structure into a preallocated transmit buffer that is already mapped into DMA memory. The device driver also attempts to coalesce transmit data in an **mbuf** chain into a single preallocated transmit buffer until the total transmit data size exceeds that of the preallocated buffer (2048 bytes). Valid values range from 64 through 2048. The default value is 2048.

#### **delay\_open**

When you set the attribute value to **Yes**, the adapter device driver delays its open completion until the Ethernet link status is determined to be either up or down. This prevents applications from sending data before the Ethernet link is established. Commands (for example, the **ifconfig** command) however might take longer to complete, especially when an active Ethernet link is not present. You can specify the values of **Yes** and **No**. The default value is **No**.

#### **failback**

This attribute is used with the **Failover Mode** attribute. If the **Failover Mode** attribute is enabled, setting this attribute to the **Yes** value causes the adapter to automatically fail back to the primary port if the primary port recovers. You can specify the values of **Yes** and **No**. The default value is **Yes**.

#### **failback\_delay**

This attribute is used with the **failback** attribute. If the **failback** attribute is enabled, the **failback\_delay** attribute specifies the number of seconds that the adapter waits before failing back to the primary port, after the primary port recovers. This delay is useful for ensuring that the primary port has fully recovered and for allowing switch protocols (for example, Spanning Tree Protocol) to complete. Valid values range from 0 through 300 seconds. Setting the **failback\_delay** attribute to 0 seconds disables the delay timer, causing failback to occur immediately. The default value is 15 seconds.

### **2-Port 10/100/1000 Base-TX PCI-X Adapter (14108902)**

The 2-Port 10/100/1000 Base-TX PCI-X Adapter Device Driver (14108902) supports the following additional configuration parameters:

#### **Transmit Descriptor Queue Size**

Indicates the number of transmit requests that can be queued for transmission by the adapter. Valid values range from 128 through 1024.

#### **Receive Descriptor Queue Size**

Indicates the maximum number of received Ethernet packets that the adapter can hold in its buffer. Valid values range from 128 through 1024.

#### **Software Transmit Queue Size**

Indicates the number of transmit requests that can be queued for transmission by the device driver. Valid values range from 512 through 16 384.

#### **Media Speed**

Indicates the speed at which the adapter attempts to operate. The available speeds are 10 Mbps half-duplex, 10 Mbps full-duplex, 100 Mbps half-duplex, 100 Mbps full-duplex, 1000 Mbps full-duplex, and autonegotiation. The default is autonegotiation. Select autonegotiate when the

adapter should use autonegotiation across the network to determine the speed. When the network does not support autonegotiation, select the specific speed.

### Transmit Jumbo Frames

When you set the attribute value to **Yes**, frames up to 9018 bytes in length can be transmitted on this adapter. When you set the attribute value to **No**, the maximum size of frames transmitted is 1518 bytes. Frames up to 9018 bytes in length can always be received on this adapter.

### Enable Hardware Transmit TCP Resegmentation

Permits the adapter to perform resegmentation of transmitted TCP segments in hardware. With this capability, the host can use TCP segments that are larger than the actual MTU size of the Ethernet link, which can increase system performance. You can specify the values of **Yes** and **No**.

**Note:** The **mbuf** structure, which describes a transmitted frame, contains a flag that indicates whether the adapter performs TCP resegmentation for the frame.

### Enable Hardware Transmit and Receive Checksum

When you set the attribute value to **Yes**, the adapter calculates the checksum for transmitted and received TCP frames. When you set the attribute value to **No**, the checksum is calculated by appropriate software.

**Note:** The **mbuf** structure, which describes a transmitted frame, contains a flag that indicates whether the adapter calculates the TCP checksum for the frame.

### Failover Mode (failover)

Indicates the requested failover configuration for the port. You can specify the values of **primary**, **backup**, and **disable**.

<b>primary</b>	Indicates the port is to act as the primary port in a failover configuration for a 2-port gigabit adapter.
<b>backup</b>	Indicates the port is to act as the backup port in a failover configuration for a 2-port gigabit adapter.
<b>disable</b>	Indicates the port is not a member of a failover configuration. This is the default value for failover.

The following configuration parameters for the 2-Port 10/100/1000 Base-TX PCI-X Adapter Device Driver (14108902) are not accessible using SMIT. You can modify them only using the **chdev** command.

### rx\_hog

When this number of receive buffer descriptors is processed by the device driver (or all packets were received), the device driver exits the **rx\_handler()** routine and continues processing other adapter events, such as transmit completions and adapter status changes. Valid values range from 1 through 1 000 000. The default value is 1000.

### slih\_hog

Indicates the number of adapter events (such as receive completions, transmit completions, and adapter status changes) that are processed by the device driver per interrupt. Valid values range from 1 through 1 000 000. The default value is 10.

### copy\_bytes

When the number of data bytes in a transmit **mbuf** structure exceeds this value, the device driver maps the **mbuf** data area into DMA memory and updates the transmit descriptor so that it points to this DMA memory area. When the number of data bytes in a transmit **mbuf** structure does not exceed this value, the data is copied from the **mbuf** structure into a preallocated transmit buffer that is already mapped into DMA memory. The device driver also attempts to coalesce transmit data in an **mbuf** chain into a single preallocated transmit buffer until the total transmit data size exceeds that of the preallocated buffer (2048 bytes). Valid values range from 64 through 2048. The default value is 2048.

### delay\_open

When you set the attribute value to **Yes**, the adapter device driver delays its open completion until the Ethernet link status is determined to be either up or down. This prevents applications from



sending data before the Ethernet link is established. Commands (for example, the **ifconfig** command) however might take longer to complete, especially when an active Ethernet link is not present. You can specify the values of **Yes** and **No**. The default value is **No**.

#### **failback**

This attribute is used with the **Failover Mode** attribute. If the **Failover Mode** attribute is enabled, setting this attribute to the **Yes** value causes the adapter to automatically fail back to the primary port if the primary port recovers. You can specify the values of **Yes** and **No**. The default value is **Yes**.

#### **failback\_delay**

This attribute is used with the **failback** attribute. If the **failback** attribute is enabled, the **failback\_delay** attribute specifies the number of seconds that the adapter waits before failing back to the primary port, after the primary port recovers. This delay is useful for ensuring that the primary port has fully recovered and for allowing switch protocols (for example, Spanning Tree Protocol) to complete. Valid values range from 0 through 300 seconds. Setting the **failback\_delay** attribute to 0 seconds disables the delay timer, causing failback to occur immediately. The default value is 15 seconds.

#### **compat\_mode**

When you set the attribute value to **Yes**, the adapter is forced to implement an early version of the IEEE 802.3z autonegotiation protocol. Use the **Yes** value only if the adapter is unable to establish a link with your older Gigabit Ethernet-TX adapters or switches. You can specify the values of **Yes** and **No**. The default value is **No**.

**Note:** If this option is enabled, the adapter cannot establish a link with newer Gigabit Ethernet-TX hardware. Enable this option only if you cannot establish a link using autonegotiation, but can force a link at a slower speed (for example, 100 full-duplex).

### **10 Gigabit Ethernet-SR PCI-X Adapter (1410ba02) and 10 Gigabit Ethernet-LR PCI\_X Adapter (1410bb02)**

The 10 Gigabit Ethernet-SR PCI-X Adapter (1410ba02) and 10 Gigabit Ethernet-LR PCI\_X Adapter (1410bb02) support the following configuration parameters:

#### **Transmit Descriptor Queue Size**

Indicates the number of transmit requests that can be queued for transmission by the adapter. Valid values range from 128 through 1024.

#### **Receive Descriptor Queue Size**

Indicates the maximum number of received Ethernet packets that the adapter can hold in its buffer. Valid values range from 128 through 1024.

#### **Software Transmit Queue**

Indicates the number of transmit requests that can be queued for transmission by the device driver. Valid values range from 512 through 16 384.

#### **Transmit Jumbo Frames**

When you set the attribute value to **Yes**, frames up to 9018 bytes in length can be transmitted on this adapter. When you set the attribute value to **No**, the maximum size of frames transmitted is 1518 bytes. Frames up to 9018 bytes in length can always be received on this adapter.

#### **Transmit TCP Resegmentation Offload**

Permits the adapter to perform resegmentation of transmitted TCP segments in hardware. With this capability, the host can use TCP segments that are larger than the actual MTU size of the Ethernet link, which can increase system performance. You can specify the **Yes** and **No** values.

#### **Enable Hardware Checksum Offload**

When you set the attribute value to **Yes**, the adapter calculates the checksum for transmitted and received TCP frames. When you set the attribute value to **No**, the checksum is calculated by appropriate software.

**Note:** The **mbuf** structure, which describes a transmitted frame, contains a flag that indicates whether the adapter calculates the checksum for the frame.

## 10 Gigabit Ethernet-SR PCI-X 2.0 DDR Adapter (1410eb02) and 10 Gigabit Ethernet-LR PCI-X 2.0 DDR Adapter (1410ec02)

The 10 Gigabit Ethernet-SR PCI-X 2.0 DDR Adapter (1410eb02) and the 10 Gigabit Ethernet-SR PCI-X 2.0 DDR Adapter (1410ec02) support the following configuration parameters:

### Transmit Jumbo Frames

When you set the attribute value to **Yes**, frames up to 9018 bytes in length can be transmitted on this adapter. The TCP/IP settings for the interface associated with the adapter are automatically initialized to maximum transmission unit MTU 9000 when the **Yes** value is selected. When you set the attribute value to **No**, the maximum size of frames transmitted is 1518 bytes and the MTU is 1500. When the jumbo frame setting is enabled, it might be possible to communicate only with other network nodes that are also jumbo-enabled and that have the same MTU. This feature can result in a considerable performance improvement. Frames up to 9018 bytes in length can always be received on this adapter.

### Enable Hardware Transmit TCP Resegmentation

Permits the adapter to perform resegmentation of transmitted TCP packets that are transmitted over IPv4 and IPv6 in hardware. With this capability, the host can create TCP segments that are larger than the actual MTU size of the Ethernet link (packet sizes of up to 64 KB can be created). The adapter then subdivides these very large packets into multiple MTU-size packets. This offloading of packet creation is recommended for increased system performance. You can specify the values of **Yes** and **No**.

### Enable Hardware Transmit and Receive Checksum Offload

When you set this attribute to **Yes**, the adapter calculates the checksum for TCP frames transmitted and received over IPv4 and IPv6. This setting is suggested for improved system performance. When you specify the **No** value, the checksum is calculated by the appropriate system software.

**Note:** The **mbuf** structure, which describes a transmitted frame, contains a flag that indicates whether the adapter must calculate the checksum for the frame.

### Enable Hardware Receive UDP Checksum Offload

Setting this attribute to the Yes value indicates that the adapter calculates the checksum for UDP standard and fragmented frames received over IPv4 and IPv6. This setting is suggested for improved system performance. If you specify the No value, the checksum is calculated by the appropriate system software.

**Note:** The **mbuf** structure that describes a received frame contains a flag that indicates whether the adapter calculated the checksum for the frame.

## Gigabit Ethernet-SX Adapter Device Driver (e414a816)

The Gigabit Ethernet-SX Adapter Device Driver (e414a816) supports the following additional configuration parameters:

### Transmit Jumbo Frames

When you set the attribute value to **Yes**, frames up to 9018 bytes in length can be transmitted on this adapter. When you set the attribute value to **No**, the maximum size of frames transmitted is 1518 bytes. Frames up to 9018 bytes in length can always be received on this adapter.

### Enable Hardware Checksum Offload

When you set the attribute value to **Yes**, the adapter calculates the checksum for transmitted and received TCP frames. When you set the attribute value to **No**, the checksum is calculated by appropriate software.

**Note:** The **mbuf** structure, which describes a transmitted frame, contains a flag that indicates whether the adapter calculates the checksum for the frame.

### Media Speed

The Gigabit Ethernet-SX Adapter Device Driver (e414a816) supports a user-configurable media speed for 1000 Mbps full-duplex and autonegotiation. The media speed attribute indicates the speed at which the adapter attempts to operate. Select autonegotiate when the adapter should use autonegotiation across the network to determine the speed. When the network does not support autonegotiation, select the specific speed.

**Note:** The default values for the Gigabit Ethernet-SX Adapter Device Driver (e414a816) configuration parameters were chosen for optimal performance. Do not change these default values.

The following configuration parameters for the Gigabit Ethernet-SX Adapter Device Driver (e414a816) are not accessible using SMIT. You can modify them only using the **chdev** command.

#### stat\_ticks

Indicates the number of microseconds that the adapter waits before updating the adapter statistics (through a DMA write) and generating an interrupt to the host. Valid values range from 1000 through 1 000 000. The default value is 1 000 000.

#### receive\_ticks

Indicates the number of microseconds that the adapter waits before updating the producer index of the receive return ring (through a DMA write) and generating an interrupt to the host. Valid values range from 0 through 1000, the default value is 50.

#### receive\_bds

Indicates the number of receive buffers that the adapter transfers to host memory before updating the producer index of the receive return ring (through a DMA write) and generating an interrupt to the host. Valid values range from 0 through 128. The default value is 6.

#### tx\_done\_ticks

Indicates the number of microseconds that the adapter waits before updating the send consumer index (through a DMA write) and generating an interrupt to the host. Valid values range from 0 through 1 000 000. The default value is 1 000 000.

#### tx\_done\_count

Indicates the number of transmit buffers that the adapter transfers from host memory before updating the send consumer index (through a DMA write) and generating an interrupt to the host. Valid values range from 0 through 128. The default value is 64.

#### receive\_proc

When this number of receive buffer descriptors is processed by the device driver (or all packets are received), the device driver returns this number of receive buffer descriptors to the adapter through an MMIO write. Valid values range from 1 through 64. The default value is 16.

#### rxdesc\_count

When this number of receive buffer descriptors is processed by the device driver (or all packets were received), the device driver exits the **rx\_handler()** routine and continues processing other adapter events, such as transmit completions and adapter status changes. Valid values range from 1 through 1 000 000. The default value is 1000.

#### slih\_hog

Indicates the number of adapter events (such as receive completions, transmit completions and adapter status changes) that are processed by the device driver per interrupt. Valid values range from 1 through 1 000 000. The default value is 10.

#### copy\_bytes

When the number of data bytes in a transmit **mbuf** structure exceeds this value, the device driver maps the **mbuf** data area into DMA memory and updates the transmit descriptor so that it points to this DMA memory area. When the number of data bytes in a transmit **mbuf** structure does not

exceed this value, the data is copied from the **mbuf** structure into a preallocated transmit buffer that is already mapped into DMA memory. The device driver also attempts to coalesce transmit data in an **mbuf** chain into a single preallocated transmit buffer until the total transmit data size exceeds that of the preallocated buffer (2048 bytes). Valid values range from 64 through 2048. The default value is 2048.

## Gigabit Ethernet-SX Adapter Device Driver (14101403)

The Gigabit Ethernet-SX Adapter Device Driver (14101403) supports the following additional configuration parameters:

### Transmit Jumbo Frames

When you set the attribute value to **Yes**, frames up to 9018 bytes in length can be transmitted on this adapter. When you set the attribute value to **No**, the maximum size of frames transmitted is 1518 bytes. Frames up to 9018 bytes in length can always be received on this adapter.

### Enable Hardware Checksum Offload

When you set the attribute value to **Yes**, the adapter calculates the checksum for transmitted and received TCP frames. When you set the attribute value to **No**, the checksum is calculated by appropriate software.

**Note:** The **mbuf** structure, which describes a transmitted frame, contains a flag that indicates whether the adapter calculates the checksum for the frame.

### Media Speed

The Gigabit Ethernet-SX Adapter Device Driver (14101403) supports a user-configurable media speed for 1000 Mbps full-duplex and autonegotiation. The media speed attribute indicates the speed at which the adapter attempts to operate. Select autonegotiate when the adapter should use autonegotiation across the network to determine the speed. When the network does not support autonegotiation, select the specific speed.

**Note:** The default values for the Gigabit Ethernet-SX Adapter Device Driver (14101403) configuration parameters were chosen for optimal performance. Do not change these default values.

The following configuration parameters for the Gigabit Ethernet-SX Adapter Device Driver (14101403) SMIT. You can modify them only using the **chdev** command.

### receive\_ticks

Indicates the number of microseconds that the adapter waits before updating the producer index of the receive return ring (through a DMA write) and generating an interrupt to the host. Valid values range from 0 through 1000, the default value is 50.

### receive\_bds

Indicates the number of receive buffers that the adapter transfers to host memory before updating the producer index of the receive return ring (through a DMA write) and generating an interrupt to the host. Valid values range from 0 through 128. The default value is 6.

### tx\_done\_ticks

Indicates the number of microseconds that the adapter waits before updating the send consumer index (through a DMA write) and generating an interrupt to the host. Valid values range from 0 through 1 000 000. The default value is 1 000 000.

### tx\_done\_count

Indicates the number of transmit buffers that the adapter transfers from host memory before updating the send consumer index (through a DMA write) and generating an interrupt to the host. Valid values range from 0 through 128. The default value is 64.

### receive\_proc

When this number of receive buffer descriptors is processed by the device driver (or all packets are received), the device driver returns this number of receive buffer descriptors to the adapter through an MMIO write. Valid values range from 1 through 64. The default value is 16.

### **rx\_hog**

When this number of receive buffer descriptors is processed by the device driver (or all packets were received), the device driver exits the **rx\_handler()** routine and continues processing other adapter events, such as transmit completions and adapter status changes. Valid values range from 1 through 1 000 000. The default value is 1000.

### **slih\_hog**

Indicates the number of adapter events (such as receive completions, transmit completions, and adapter status changes) that are processed by the device driver per interrupt. Valid values range from 1 through 1 000 000. The default value is 10.

### **copy\_bytes**

When the number of data bytes in a transmit **mbuf** structure exceeds this value, the device driver maps the **mbuf** data area into DMA memory and updates the transmit descriptor so that it points to this DMA memory area. When the number of data bytes in a transmit **mbuf** structure does not exceed this value, the data is copied from the **mbuf** structure into a preallocated transmit buffer that is already mapped into DMA memory. The device driver also attempts to coalesce transmit data in an **mbuf** chain into a single preallocated transmit buffer until the total transmit data size exceeds that of the preallocated buffer (2048 bytes). Valid values range from 64 through 2048. The default value is 2048.

### **delay\_open**

When you set the attribute value to **Yes**, the adapter device driver delays its open completion until the Ethernet link status is determined to be either up or down. This prevents applications from sending data before the Ethernet link is established. Commands (for example, the **ifconfig** command) however might take longer to complete, especially when an active Ethernet link is not present. You can specify the values of **Yes** and **No**. The default value is **No**.

## **Gigabit Ethernet-SX Adapter Device Driver (14106703)**

The Gigabit Ethernet-SX Adapter Device Driver (14106703) supports the following additional configuration parameters:

### **Transmit Jumbo Frames**

When you set the attribute value to **Yes**, frames up to 9018 bytes in length can be transmitted on this adapter. When you set the attribute value to **No**, the maximum size of frames transmitted is 1518 bytes. Frames up to 9018 bytes in length can always be received on this adapter.

### **Enable Hardware Checksum Offload**

When you set the attribute value to **Yes**, the adapter calculates the checksum for transmitted and received TCP frames. When you set the attribute value to **No**, the checksum is calculated by appropriate software.

**Note:** The **mbuf** structure, which describes a transmitted frame, contains a flag that indicates whether the adapter calculates the checksum for the frame.

### **Media Speed**

Indicates the speed at which the adapter attempts to operate. You can specify the values of 1000 Mbps full-duplex and autonegotiation. The default is autonegotiation. Select autonegotiate when the adapter should use autonegotiation across the network to determine the duplexity. When the network does not support autonegotiation, select 1000 Mbps full-duplex.

**Note:** The default values for the Gigabit Ethernet-SX Adapter Device Driver (14106703) configuration parameters are chosen for optimal performance. Do not change these default values.

The following configuration parameters for the Gigabit Ethernet-SX Adapter Device Driver (14106703) are not accessible using SMIT. You can modify them only using the **chdev** command.

### **receive\_ticks**

Indicates the number of microseconds that the adapter waits before updating the producer index of

the receive return ring (through a DMA write) and generating an interrupt to the host. Valid values range from 0 through 1000. The default value is 50.

#### **receive\_bds**

Indicates the number of receive buffers that the adapter transfers to host memory before updating the producer index of the receive return ring (through a DMA write) and generating an interrupt to the host. Valid values range from 0 through 128. The default value is 6.

#### **tx\_done\_ticks**

Indicates the number of microseconds that the adapter waits before updating the send consumer index (through a DMA write) and generating an interrupt to the host. Valid values range from 0 through 1 000 000. The default value is 1 000 000.

#### **tx\_done\_count**

Indicates the number of transmit buffers that the adapter transfers from host memory before updating the send consumer index (through a DMA write) and generating an interrupt to the host. Valid values range from 0 through 128. The default value is 64.

#### **receive\_proc**

When this number of receive buffer descriptors is processed by the device driver (or all packets are received), the device driver returns this number of receive buffer descriptors to the adapter through an MMIO write. Valid values range from 1 through 64. The default value is 16.

#### **rx\_hog**

When this number of receive buffer descriptors is processed by the device driver (or all packets were received), the device driver exits the **rx\_handler()** routine and continues processing other adapter events, such as transmit completions and adapter status changes. Valid values range from 1 through 1 000 000. The default value is 1000.

#### **slih\_hog**

Indicates the number of adapter events (such as receive completions, transmit completions, and adapter status changes) that is processed by the device driver per interrupt. Valid values range from 1 through 1 000 000. The default value is 10.

#### **copy\_bytes**

When the number of data bytes in a transmit **mbuf** structure exceeds this value, the device driver maps the **mbuf** data area into DMA memory and updates the transmit descriptor so that it points to this DMA memory area. When the number of data bytes in a transmit **mbuf** structure does not exceed this value, the data is copied from the **mbuf** structure into a preallocated transmit buffer that is already mapped into DMA memory. The device driver also attempts to coalesce transmit data in an **mbuf** chain into a single preallocated transmit buffer until the total transmit data size exceeds that of the preallocated buffer (2048 bytes). Valid values range from 64 through 2048. The default value is 2048.

#### **delay\_open**

When you set the attribute value to **Yes**, the adapter device driver delays its open completion until the Ethernet link status is determined to be either up or down. This prevents applications from sending data before the Ethernet link is established. Commands, such as the **ifconfig** command, however, might take longer to complete, especially when an active Ethernet link is not present. You can specify the values of **Yes** and **No**. The default value is **No**.

### **4-Port 10/100/1000 Base-TX PCI-X Adapter (14101103)**

The 4-Port 10/100/1000 Base-TX PCI-X Adapter (14101103) supports the following additional configuration parameters:

#### **Transmit Descriptor Queue Size**

Indicates the number of transmit requests that can be queued for transmission by the adapter. Valid values range from 128 through 1024.

### Receive Descriptor Queue Size

Indicates the maximum number of received Ethernet packets that the adapter can hold in its buffer. Valid values range from 128 through 1024.

### Software Transmit Queue Size

Indicates the number of transmit requests that can be queued for transmission by the device driver. Valid values range from 512 through 16 384.

### Media Speed

Indicates the speed at which the adapter attempts to operate. The available speeds are 10 Mbps half-duplex, 10 Mbps full-duplex, 100 Mbps half-duplex, 100 Mbps full-duplex, 1000 Mbps full-duplex, and autonegotiation. The default is autonegotiation. Select autonegotiate when the adapter should use autonegotiation across the network to determine the speed. When the network does not support autonegotiation, select the specific speed.

**Note:** 1000 Mbps half-duplex is not a valid value. The IEEE 802.3z specification dictates that the gigabit speeds for half-duplex must be autonegotiated for copper (TX)-based adapters. Select autonegotiation if this speed is required.

### Transmit Jumbo Frames

When you set the attribute value to **Yes**, frames up to 9018 bytes in length can be transmitted on this adapter. When you set the attribute value to **No**, the maximum size of frames transmitted is 1518 bytes. Frames up to 9018 bytes in length can always be received on this adapter.

### Transmit TCP Resegmentation Offload

Permits the adapter to perform resegmentation of transmitted TCP segments in hardware. With this capability, the host can use TCP segments that are larger than the actual MTU size of the Ethernet link, which can increase system performance. You can specify the **Yes** and **No** values.

### Enable Hardware Checksum Offload

When you set the attribute value to **Yes**, the adapter calculates the checksum for transmitted and received TCP frames. When you set the attribute value to **No**, the checksum is calculated by appropriate software.

**Note:** The **mbuf** structure, which describes a transmitted frame, contains a flag that indicates whether the adapter calculates the checksum for the frame.

### Gigabit Backward Compatibility

Older gigabit TX equipment might not be able to communicate with this adapter. If the adapter is unable to communicate with your older gigabit equipment, enabling this option forces the adapter to implement the IEEE 802.3z incorrectly. As such, this option should be enabled if the adapter is unable to communicate with your older gigabit equipment.

**Note:** Enabling this option forces the adapter to implement the IEEE 802.3z incorrectly. If this option is enabled, the adapter cannot communicate with newer equipment. Enable this option only if you cannot obtain a link using autonegotiation, but can force a link at a slower speed (for example, 100 full-duplex).

### Failover Mode (failover)

Indicates the requested failover configuration for the port. You can specify the values of **primary**, **backup**, and **disable**. You can change this attribute using SMIT.

<b>primary</b>	Indicates the port is to act as the primary port in a failover configuration for a 2-port gigabit adapter.
<b>backup</b>	Indicates the port is to act as the backup port in a failover configuration for a 2-port gigabit adapter.
<b>disable</b>	Indicates the port is not a member of a failover configuration. This is the default value for failover.

## 4-Port 10/100/1000 Base-TX PCI-Express Adapter (14106803)

The 4-Port 10/100/1000 Base-TX PCI-Express Adapter (14106803) supports the following additional configuration parameters:

### Transmit Descriptor Queue Size

Indicates the number of transmit requests that can be queued for transmission by the adapter. Valid values range from 128 through 1024.

### Receive Descriptor Queue Size

Indicates the maximum number of received Ethernet packets that the adapter can hold in its buffer. Valid values range from 128 through 1024.

### Software Transmit Queue Size

Indicates the number of transmit requests that can be queued for transmission by the device driver. Valid values range from 512 through 16 384.

### Media Speed

Indicates the speed at which the adapter attempts to operate. The available speeds are 10 Mbps half-duplex, 10 Mbps full-duplex, 100 Mbps half-duplex, 100 Mbps full-duplex, 1000 Mbps full-duplex, and autonegotiation. The default is autonegotiation. Select autonegotiate when the adapter must use autonegotiation across the network to determine the speed. When the network does not support autonegotiation, select the specific speed.

**Note:** 1000 Mbps half-duplex is not a valid value. The IEEE 802.3z specification dictates that the gigabit speeds for half-duplex must be autonegotiated for copper (TX)-based adapters. Select autonegotiation if this speed is required.

### Transmit Jumbo Frames

When you set the attribute value to **Yes**, frames up to 9018 bytes in length can be transmitted on this adapter. When you set the attribute value to **No**, the maximum size of frames transmitted is 1518 bytes. Frames up to 9018 bytes in length can always be received on this adapter.

### Transmit TCP Resegmentation Offload

Permits the adapter to perform resegmentation of transmitted TCP segments in hardware. With this capability, the host can use TCP segments that are larger than the actual MTU size of the Ethernet link, which can increase system performance. You can specify the **Yes** and **No** values.

### Enable Hardware Checksum Offload

When you set the attribute value to **Yes**, the adapter calculates the checksum for transmitted and received TCP frames. When you set the attribute value to **No**, the checksum is calculated by appropriate software.

**Note:** The **mbuf** structure, which describes a transmitted frame, contains a flag that indicates whether the adapter calculates the checksum for the frame.

### Gigabit Backward Compatibility

Forces the adapter to implement the IEEE 802.3z incorrectly. Older gigabit TX equipment might not be able to communicate with the adapter. Enable the option if the adapter is unable to communicate with your older gigabit equipment.

**Important:** If the option is enabled, the adapter cannot communicate with newer equipment. Enable the option only if you cannot obtain a link using autonegotiation, but can force a link at a slower speed (for example, 100 full-duplex).

### Failover Mode (failover)

Indicates the requested failover configuration for the port. You can specify the values of **primary**, **backup**, and **disable**. You can change this attribute using SMIT.

<b>primary</b>	Indicates the port is to act as the primary port in a failover configuration for a 2-port gigabit adapter.
<b>backup</b>	Indicates the port is to act as the backup port in a failover configuration for a 2-port gigabit adapter.
<b>disable</b>	Indicates the port is not a member of a failover configuration. This is the default value for failover.



## 2-Port Gigabit Ethernet-SX PCI-Express Adapter (14103f03)

The 2-Port Gigabit Ethernet-SX PCI-Express Adapter (14103f03) supports the following additional configuration parameters:

### Transmit Descriptor Queue Size

Indicates the number of transmit requests that can be queued for transmission by the adapter. Valid values range from 128 through 1024.

### Receive Descriptor Queue Size

Indicates the maximum number of received Ethernet packets that the adapter can hold in its buffer. Valid values range from 128 through 1024.

### Software Transmit Queue Size

Indicates the number of transmit requests that can be queued for transmission by the device driver. Valid values range from 512 through 16 384.

### Media Speed

Indicates the speed at which the adapter attempts to operate. You can specify the values of 1000 Mbps full-duplex and autonegotiation. The default is autonegotiation. Select autonegotiate when the adapter should use autonegotiation across the network to determine the duplexity. When the network does not support autonegotiation, select 1000 Mbps full-duplex.

### Transmit Jumbo Frames

When you set the attribute value to **Yes**, frames up to 9018 bytes in length can be transmitted on this adapter. When you set the attribute value to **No**, the maximum size of frames transmitted is 1518 bytes. Frames up to 9018 bytes in length can always be received on this adapter.

### Transmit TCP Resegmentation Offload

Permits the adapter to perform resegmentation of transmitted TCP segments in hardware. With this capability, the host can use TCP segments that are larger than the actual MTU size of the Ethernet link, which can increase system performance. You can specify the **Yes** and **No** values.

### Enable Hardware Checksum Offload

When you set the attribute value to **Yes**, the adapter calculates the checksum for transmitted TCP frames and received TCP frames. When you set the attribute value to **No**, the checksum is calculated by appropriate software.

**Note:** The **mbuf** structure, which describes a transmitted frame, contains a flag that indicates whether the adapter calculates the checksum for the frame.

### Failover Mode (failover)

Indicates the requested failover configuration for the port. You can specify the values of **primary**, **backup**, and **disable**. You can change this attribute using SMIT.

<b>primary</b>	Indicates the port is to act as the primary port in a failover configuration for a 2-port gigabit adapter.
<b>backup</b>	Indicates the port is to act as the backup port in a failover configuration for a 2-port gigabit adapter.
<b>disable</b>	Indicates the port is not a member of a failover configuration. This is the default value for failover.

The following configuration parameters for the 2-Port Gigabit Ethernet-SX PCI-Express Adapter Device Driver (14103f03) are not accessible using SMIT. You can modify them only using the **chdev** command.

### rx\_hog

When this number of receive buffer descriptors is processed by the device driver (or all packets were received), the device driver exits the **rx\_handler()** routine and continues processing other adapter events, such as transmit completions and adapter status changes. Valid values range from 1 through 1 000 000. The default value is 1000.

### slih\_hog

Indicates the number of adapter events (such as receive completions, transmit completions, and

adapter status changes) that are processed by the device driver per interrupt. Valid values range from 1 through 1 000 000. The default value is 10.

#### **copy\_bytes**

When the number of data bytes in a transmit **mbuf** structure exceeds this value, the device driver maps the **mbuf** data area into DMA memory and updates the transmit descriptor so that it points to this DMA memory area. When the number of data bytes in a transmit **mbuf** structure does not exceed this value, the data is copied from the **mbuf** structure into a preallocated transmit buffer that is already mapped into DMA memory. The device driver also attempts to coalesce transmit data in an **mbuf** chain into a single preallocated transmit buffer until the total transmit data size exceeds that of the preallocated buffer (2048 bytes). Valid values range from 64 through 2048. The default value is 2048.

#### **delay\_open**

When you set the attribute value to **Yes**, the adapter device driver delays its open completion until the Ethernet link status is determined to be either up or down. This prevents applications from sending data before the Ethernet link is established. Commands (for example, the **ifconfig** command) however might take longer to complete, especially when an active Ethernet link is not present. You can specify the values of **Yes** and **No**. The default value is **No**.

#### **failback**

The attribute is used with the **Failover Mode** attribute. When the **Failover Mode** attribute is enabled, setting this attribute value to **Yes** causes the adapter to automatically fail back to the primary port if the primary port recovers. You can specify the values of **Yes** and **No**. The default value is **Yes**.

#### **failback\_delay**

The attribute is used with the **failback** attribute. When the **failback** attribute is enabled, the **failback\_delay** attribute specifies the number of seconds that the adapter waits before failing back to the primary port, after the primary port recovers. This delay is useful for ensuring that the primary port has fully recovered and for allowing switch protocols (for example, Spanning Tree Protocol) to complete. Valid values range from 0 through 300 seconds. Setting the **failback\_delay** attribute to 0 seconds disables the delay timer, causing failback to occur immediately. The default value is 15 seconds.

### **2-Port 10/100/1000 Base-TX PCI-Express Adapter (14104003)**

The 2-Port 10/100/1000 Base-TX PCI-Express Adapter Device Driver (14104003) supports the following additional configuration parameters:

#### **Transmit Descriptor Queue Size**

Indicates the number of transmit requests that can be queued for transmission by the adapter. Valid values range from 128 through 1024.

#### **Receive Descriptor Queue Size**

Indicates the maximum number of received Ethernet packets that the adapter can hold in its buffer. Valid values range from 128 through 1024.

#### **Software Transmit Queue Size**

Indicates the number of transmit requests that can be queued for transmission by the device driver. Valid values range from 512 through 16 384.

#### **Media Speed**

Indicates the speed at which the adapter attempts to operate. The available speeds are 10 Mbps half-duplex, 10 Mbps full-duplex, 100 Mbps half-duplex, 100 Mbps full-duplex, 1000 Mbps full-duplex, and autonegotiation. The default is autonegotiation. Select autonegotiate when the adapter should use autonegotiation across the network to determine the speed. When the network does not support autonegotiation, select a specific speed.

**Note:** 1000 Mbps half-duplex is not a valid value. The IEEE 802.3z specification dictates that the gigabit speeds for half-duplex must be autonegotiated for copper (TX)-based adapters. Select autonegotiation if this speed is required.

### Transmit Jumbo Frames

When you set the attribute value to **Yes**, frames up to 9018 bytes in length can be transmitted on this adapter. When you set the attribute value to **No**, the maximum size of frames transmitted is 1518 bytes. Frames up to 9018 bytes in length can always be received on this adapter.

### Enable Hardware Transmit TCP Resegmentation

Permits the adapter to perform resegmentation of transmitted TCP segments in hardware. With this capability, the host can use TCP segments that are larger than the actual MTU size of the Ethernet link, which can increase system performance. You can specify the values of **Yes** and **No**.

**Note:** The **mbuf** structure, which describes a transmitted frame, contains a flag that indicates whether the adapter performs TCP resegmentation for the frame.

### Enable Hardware Transmit and Receive Checksum

When you set the attribute value to **Yes**, the adapter calculates the checksum for transmitted and received TCP frames. When you set the attribute value to **No**, the checksum is calculated by appropriate software.

**Note:** The **mbuf** structure, which describes a transmitted frame, contains a flag that indicates whether the adapter calculates the TCP checksum for the frame.

### Failover Mode (failover)

Indicates the requested failover configuration for the port. You can specify the attribute values of **primary**, **backup**, and **disable**.

<b>primary</b>	Indicates the port is to act as the primary port in a failover configuration for a 2-port gigabit adapter.
<b>backup</b>	Indicates the port is to act as the backup port in a failover configuration for a 2-port gigabit adapter.
<b>disable</b>	Indicates the port is not a member of a failover configuration. This is the default value for failover.

The following configuration parameters for the 2-Port 10/100/1000 Base-TX PCI-Express Adapter Device Driver (14104003) are not accessible using SMIT. You can modify them only using the **chdev** command.

### rx\_hog

When this number of receive buffer descriptors is processed by the device driver (or all packets were received), the device driver exits the **rx\_handler()** routine and continues processing other adapter events, such as transmit completions and adapter status changes. Valid values range from 1 through 1 000 000. The default value is 1000.

### slih\_hog

Indicates the number of adapter events (such as receive completions, transmit completions, and adapter status changes) that is processed by the device driver per interrupt. Valid values range from 1 through 1 000 000. The default value is 10.

### copy\_bytes

When the number of data bytes in a transmit **mbuf** structure exceeds this value, the device driver maps the **mbuf** data area into DMA memory and updates the transmit descriptor so that it points to this DMA memory area. When the number of data bytes in a transmit **mbuf** structure does not exceed this value, the data is copied from the **mbuf** structure into a preallocated transmit buffer that is already mapped into DMA memory. The device driver also attempts to coalesce transmit data in an **mbuf** chain into a single preallocated transmit buffer until the total transmit data size exceeds that of the preallocated buffer (2048 bytes). Valid values range from 64 through 2048. The default value is 2048.

### delay\_open

When you set the attribute value to **Yes**, the adapter device driver delays its open completion until

the Ethernet link status is determined to be either up or down. This prevents applications from sending data before the Ethernet link is established. Commands (for example, the **ifconfig** command), however, might take longer to complete, especially when an active Ethernet link is not present. You can specify the values of **Yes** and **No**. The default value is **No**.

#### **failback**

This attribute is used with the **Failover Mode** attribute. If the **Failover Mode** attribute is enabled, setting this attribute to the **Yes** value causes the adapter to automatically fail back to the primary port if the primary port recovers. You can specify the values of **Yes** and **No**. The default value is **Yes**.

#### **failback\_delay**

This attribute is used with the **failback** attribute. If the **failback** attribute is enabled, the **failback\_delay** attribute specifies the number of seconds that the adapter waits before failing back to the primary port, after the primary port recovers. This delay is useful for ensuring that the primary port has fully recovered and for allowing switch protocols (for example, Spanning Tree Protocol) to complete. Valid values range from 0 through 300 seconds. Setting the **failback\_delay** attribute to 0 seconds disables the delay timer, causing failback to occur immediately. The default value is 15 seconds.

#### **compat\_mode**

When you set the attribute value to **Yes**, the adapter is forced to implement an early version of the IEEE 802.3z autonegotiation protocol. Use the **Yes** value only if the adapter is unable to establish a link with your older Gigabit Ethernet-TX adapters or switches. You can specify the values of **Yes** and **No**. The default value is **No**.

**Note:** If this option is enabled, the adapter cannot establish a link with newer Gigabit Ethernet-TX hardware. Enable this option only if you cannot establish a link using autonegotiation, but can force a link at a slower speed (for example, 100 full-duplex).

### **Host Ethernet Adapter Device Driver**

The Host Ethernet Adapter Device Driver supports the following additional configuration parameters:

#### **Enable Alternate Ethernet Address (use\_alt\_addr)**

Specifies whether the adapter uses the MAC address that is specified in the **Alternate Ethernet Address** attribute. You can set the attribute value to **Yes** or **No**. The default value is **No**.

#### **Enable RX and TX Checksum Offload of TCP Segments (tx\_cksum and rx\_cksum)**

Specifies whether the adapter is to perform receive checksum calculation and transmit checksum calculation for TCP segments. You can set the attribute value to **Yes** or **No**. The default value is **Yes**.

#### **Flow Control (flow\_ctrl)**

Specifies whether the adapter enables transmit flow control and receive flow control. You can set the attribute value to **Yes** or **No**. The default value is **Yes**. The setting is only applicable to a system that is not managed by Hardware Management Console (HMC) or by Integrated Virtualization Manager (IVM).

#### **Jumbo Frames (jumbo\_frames)**

Specifies whether to transmit and receive jumbo frames, which range from 1522 bytes through 9022 bytes in size. You can set the attribute value to **Yes** or **No**. The default value is **No**. To enable the attribute, authorization from HMC or IVM is required. Disabling of this attribute is not dependent on HMC or IVM.

#### **Receive Processing Count (rx\_proc)**

Indicates the number of receive descriptor that is processed before posting a work-queue element. The default value is determined per performance tuning.

#### **Enable TX TCP Resegmentation Offload of TCP Segments (large\_send)**

Specifies whether the adapter is to perform transmit TCP resegmentation for TCP segments. You can set the attribute value to **Yes** or **No**. The default value is **Yes**.

### **Transmit Immediate Data Copy Threshold (tx\_immd\_copy)**

Specifies the maximum size for copying data into the immediate data area inside the transmit work-queue element. This value is tunable for performance. Possible values range from 136 through 224. The default value is 224.

### **Requested Media Speed (media\_speed)**

Indicates the speed at which the adapter attempts to operate. The available speeds are 10 Mbps full-duplex, 100 Mbps full-duplex, 1000 Mbps full-duplex, 10 000 Mbps full-duplex, and autonegotiation. The default is autonegotiation. Select autonegotiate when the adapter must use autonegotiation across the network to determine the speed. When the network does not support autonegotiation, select a specific speed.

**Note:** Half-duplex is not a valid value. This attribute is only applicable to a system that is not managed by HMC or by IVM.

### **Enable Multicore Scaling (Multicore)**

Enables or disables driver multithreading. You can set the attribute value to **Yes** or **No**. The default value is **Yes**. The multithreading value is specified by the HMC *MCS* parameter.

### **Use Transmit Interface Specific Buffers (tx\_isb)**

Enables or disables transmit copy optimizations. You can set the attribute value to **Yes** or **No**. The default value is **Yes**.

### **Enable Receive TCP Segment Aggregation (large\_receive)**

Enables coalescing receive packets into a larger packet before passing to the upper layer for enhanced performance. You can set the attribute value to **Yes** or **No**. The default value is **No**.

**Note:** Disable this value when the adapter is used in IP forwarding.

### **Receive Packet Coalescing (rx\_coalesce)**

Allows tuning of the RX packets that are coalesced before passing to a stack. Increasing this value impacts latency under certain scenarios. Decreasing this value might enhance large reception under limited scenarios. You can set the attribute to a value ranging from 1 through 64. The default value is 16, which demonstrates the best balance of performance versus latency.

## **Interface Entry Points**

### **Device Driver Configuration and Unconfiguration**

The configuration entry points of the device drivers conform to the guidelines for kernel object file entry points. These configuration entry points are as follows:

- **kent\_config** for the PCI Ethernet Device Driver (22100020)
- **phxent\_config** for the 10/100 Mbps Ethernet PCI Adapter Device Driver (23100020)
- **scent\_config** for the 10/100 Mbps Ethernet PCI Adapter II Device Driver (1410ff01)
- **gxent\_config** for the Gigabit Ethernet-SX PCI Adapter Device Driver (14100401)
- **goent\_config** for the Gigabit Ethernet-SX PCI-X Adapter Device Driver (14106802), the 10/100/1000 Base-T Ethernet PCI-X Adapter Device Driver (14106902), the 2-Port Gigabit Ethernet-SX PCI-X Adapter Device Driver (14108802), the 2-Port 10/100/1000 Base-TX PCI-X Adapter Device Driver (14108902), the 4-Port 10/100/1000 Base-TX PCI-X Adapter Device Driver (14101103), and 4-Port 10/100/1000 Base-TX PCI-Express Adapter Device Driver (14106803), the 2-Port Gigabit Ethernet-SX PCI-Express Adapter Device Driver (14103f03), and the 2-Port 10/100/1000 Base-TX PCI-Express Adapter Device Driver (14104003).
- **vent\_config** for the 10 Gigabit Ethernet-SR PCI-X Adapter Device Driver (1410ba02) and the 10 Gigabit Ethernet-LR PCI-X Adapter Device Driver (1410bb02).
- **bent\_config** for the Gigabit Ethernet-SX Adapter Device Driver (e414a816).
- **ment\_config** for the Gigabit Ethernet-SX Adapter Device Driver (14101403) and the Gigabit Ethernet-SX PCI-X Adapter Device Driver (14106703).

- **kn gent\_config** for the 10 Gigabit Ethernet-SR PCI-X 2.0 DDR Adapter Device Driver (1410eb02) and the 10 Gigabit Ethernet-LR PCI\_X 2.0 DDR Adapter Device Driver (1410ec02).
- **hea\_config** for Host Ethernet Adapter Device Driver.

## Device Driver Open

The open entry point for the device drivers perform a synchronous open of the specified network device.

The device driver issues commands to start the initialization of the device. The state of the device now is **OPEN\_PENDING**. The device driver invokes the open process for the device. The open process involves a sequence of events that are necessary to initialize and configure the device. The device driver does the sequence of events in an orderly fashion to make sure that one step is finished executing on the adapter before the next step is continued. Any error during these sequence of events makes the open fail. The device driver requires about 2 seconds to open the device. When the whole sequence of events is done, the device driver verifies the open status and then returns to the caller of the open with a return code to indicate open success or open failure.

After the device has been successfully configured and connected to the network, the device driver sets the device state to **OPENED**, the **NDD\_RUNNING** flag in the **NDD** flags field is turned on. In the case of unsuccessful open, both the **NDD\_UP** and **NDD\_RUNNING** flags in the **NDD** flags field are off and a non-zero error code is returned to the caller.

The open entry points are as follows:

- **kent\_open** for the PCI Ethernet Device Driver (22100020)
- **phxent\_open** for the 10/100 Mbps Ethernet PCI Adapter Device Driver (23100020)
- **scent\_open** for the 10/100 Mbps Ethernet PCI Adapter II Device Driver (1410ff01)
- **gxent\_open** for the Gigabit Ethernet-SX PCI Adapter Device Driver (14100401)
- **goent\_open** for the Gigabit Ethernet-SX PCI-X Adapter Device Driver (14106802), the 10/100/1000 Base-T Ethernet PCI-X Adapter Device Driver (14106902), the 2-Port Gigabit Ethernet-SX PCI-X Adapter Device Driver (14108802), the 2-Port 10/100/1000 Base-TX PCI-X Adapter Device Driver (14108902), the 4-Port 10/100/1000 Base-TX PCI-X Adapter Device Driver (14101103), and 4-Port 10/100/1000 Base-TX PCI-Express Adapter Device Driver (14106803), the 2-Port Gigabit Ethernet-SX PCI-Express Adapter Device Driver (14103f03), and the 2-Port 10/100/1000 Base-TX PCI-Express Adapter Device Driver (14104003).
- **vent\_open** for the 10 Gigabit Ethernet-SR PCI-X Adapter Device Driver (1410ba02) and the 10 Gigabit Ethernet-LR PCI\_X Adapter Device Driver (1410bb02).
- **bent\_open** for the Gigabit Ethernet-SX Adapter Device Driver (e414a816).
- **ment\_open** for the Gigabit Ethernet-SX Adapter Device Driver (14101403) and the Gigabit Ethernet-SX PCI-X Adapter Device Driver (14106703).
- **kn gent\_open** for the 10 Gigabit Ethernet-SR PCI-X 2.0 DDR Adapter Device Driver (1410eb02) and the 10 Gigabit Ethernet-LR PCI\_X 2.0 DDR Adapter Device Driver (1410ec02).
- **hea\_open** for Host Ethernet Adapter Device Driver.

## Device Driver Close

The close entry point for the device drivers is called to close the specified network device. This function resets the device to a known state and frees system resources associated with the device.

The device will not be detached from the network until the device's transmit queue drains. That is, the close entry point will not return until all packets have been transmitted or timed out. If the device is inoperable at the time of the close, the device's transmit queue does not have to drain.

At the beginning of the close entry point, the device state is set to be **CLOSE\_PENDING**. The **NDD\_RUNNING** flag in the **n dd\_flags** is turned off. After the outstanding transmit queue is all done, the

device driver starts a sequence of operations to deactivate the adapter and to free up resources. Before the close entry point returns to the caller, the device state is set to **CLOSED**.

The close entry points are as follows:

- **kent\_close** for the PCI Ethernet Device Driver (22100020)
- **phxent\_close** for the 10/100 Mbps Ethernet PCI Adapter Device Driver (23100020)
- **scent\_close** for the 10/100 Mbps Ethernet PCI Adapter II Device Driver (1410ff01)
- **gxent\_close** for the Gigabit Ethernet-SX PCI Adapter Device Driver (14100401)
- **goent\_close** for the Gigabit Ethernet-SX PCI-X Adapter Device Driver (14106802), the 10/100/1000 Base-T Ethernet PCI-X Adapter Device Driver (14106902), the 2-Port Gigabit Ethernet-SX PCI-X Adapter Device Driver (14108802), the 2-Port 10/100/1000 Base-TX PCI-X Adapter Device Driver (14108902), the 4-Port 10/100/1000 Base-TX PCI-X Adapter Device Driver (14101103), and 4-Port 10/100/1000 Base-TX PCI-Express Adapter Device Driver (14106803), the 2-Port Gigabit Ethernet-SX PCI-Express Adapter Device Driver (14103f03), and the 2-Port 10/100/1000 Base-TX PCI-Express Adapter Device Driver (14104003).
- **vent\_close** for the 10 Gigabit Ethernet-SR PCI-X Adapter Device Driver (1410ba02) and the 10 Gigabit Ethernet-LR PCI\_X Adapter Device Driver (1410bb02).
- **bent\_close** for the Gigabit Ethernet-SX Adapter Device Driver (e414a816).
- **ment\_close** for the Gigabit Ethernet-SX Adapter Device Driver (14101403) and the Gigabit Ethernet-SX PCI-X Adapter Device Driver (14106703).
- **kngent\_close** for the 10 Gigabit Ethernet-SR PCI-X 2.0 DDR Adapter Device Driver (1410eb02) and the 10 Gigabit Ethernet-LR PCI\_X 2.0 DDR Adapter Device Driver (1410ec02).
- **hea\_close** for Host Ethernet Adapter Device Driver.

## Data Transmission

The output entry point transmits data using the specified network device.

The data to be transmitted is passed into the device driver by way of **mbuf** structures. The first **mbuf** structure in the chain must be of **M\_PKTHDR** format. Multiple **mbuf** structures can be used to hold the frame. Link the **mbuf** structures using the **m\_next** field of the **mbuf** structure.

Multiple packet transmits are supported with the mbufs being chained using the **m\_nextpkt** field of the **mbuf** structure. The **m\_pkthdr.len** field must be set to the total length of the packet. The device driver does *not* support mbufs from user memory (which have the **M\_EXT** flag set).

On successful transmit requests, the device driver is responsible for freeing all the mbufs associated with the transmit request. If the device driver returns an error, the caller is responsible for the mbufs. If any of the chained packets can be transmitted, the transmit is considered successful and the device driver is responsible for all of the mbufs in the chain.

If the destination address in the packet is a broadcast address the **M\_BCAST** flag in the **m\_flags** field should be set prior to entering this routine. A broadcast address is defined as 0xFFFF FFFF FFFF. If the destination address in the packet is a multicast address the **M\_MCAST** flag in the **m\_flags** field should be set prior to entering this routine. A multicast address is defined as a non-individual address other than a broadcast address. The device driver keeps statistics based upon the **M\_BCAST** and **M\_MCAST** flags.

For packets that are shorter than the Ethernet minimum MTU size (60 bytes), the device driver pads them by adjusting the transmit length to the adapter so they can be transmitted as valid Ethernet packets.

Users are not notified by the device driver about the status of the transmission. Various statistics about data transmission are kept by the driver in the **ndd** structure. These statistics are part of the data returned by the **NDD\_GET\_STATS** control operation.

The output entry points are as follows:

- **kent\_output** for the PCI Ethernet Device Driver (22100020)
- **phxent\_output** for the 10/100 Mbps Ethernet PCI Adapter Device Driver (23100020)
- **scent\_output** for the 10/100 Mbps Ethernet PCI Adapter II Device Driver (1410ff01)
- **gxent\_output** for the Gigabit Ethernet-SX PCI Adapter Device Driver (14100401)
- **goent\_output** for the Gigabit Ethernet-SX PCI-X Adapter Device Driver (14106802), the 10/100/1000 Base-T Ethernet PCI-X Adapter Device Driver (14106902), the 2-Port Gigabit Ethernet-SX PCI-X Adapter Device Driver (14108802), the 2-Port 10/100/1000 Base-TX PCI-X Adapter Device Driver (14108902), the 4-Port 10/100/1000 Base-TX PCI-X Adapter Device Driver (14101103), and 4-Port 10/100/1000 Base-TX PCI-Express Adapter Device Driver (14106803), the 2-Port Gigabit Ethernet-SX PCI-Express Adapter Device Driver (14103f03), and the 2-Port 10/100/1000 Base-TX PCI-Express Adapter Device Driver (14104003).
- **vent\_output** for the 10 Gigabit Ethernet-SR PCI-X Adapter Device Driver (1410ba02) and the 10 Gigabit Ethernet-LR PCI\_X Adapter Device Driver (1410bb02).
- **bent\_output** for the Gigabit Ethernet-SX Adapter Device Driver (e414a816).
- **ment\_output** for the Gigabit Ethernet-SX Adapter Device Driver (14101403) and the Gigabit Ethernet-SX PCI-X Adapter Device Driver (14106703).
- **kgent\_output** for the 10 Gigabit Ethernet-SR PCI-X 2.0 DDR Adapter Device Driver (1410eb02) and the 10 Gigabit Ethernet-LR PCI\_X 2.0 DDR Adapter Device Driver (1410ec02).
- **hea\_output** for Host Ethernet Adapter Device Driver.

## Data Reception

When the Ethernet device drivers receive a valid packet from the network device, the device drivers call the **nd\_receive** function that is specified in the **ndd\_t** structure of the network device. The **nd\_receive** function is part of a CDLI network demultiplexer. The packet is passed to the **nd\_receive** function in the form of a mbuf.

The Ethernet device drivers can pass multiple packets to the **nd\_receive** function by chaining the packets together using the **m\_nextpkt** field of the **mbuf** structure. The **m\_pkthdr.len** field must be set to the total length of the packet. If the source address in the packet is a broadcast address the **M\_BCAST** flag in the **m\_flags** field should be set. If the source address in the packet is a multicast address the **M\_MCAST** flag in the **m\_flags** field should be set.

When the device driver initially configures the device to discard all invalid frames. A frame is considered to be invalid for the following reasons:

- The packet is too short.
- The packet is too long.
- The packet contains a CRC error.
- The packet contains an alignment error.

If the asynchronous status for receiving invalid frames has been issued to the device driver, the device driver configures the device to receive bad packets as well as good packets. Whenever a bad packet is received by the driver, an asynchronous status block **NDD\_BAD\_PKTS** is created and delivered to the appropriate user. The user must copy the contents of the mbuf to another memory area. The user must not modify the contents of the mbuf or free the mbuf. The device driver has the responsibility of releasing the mbuf upon returning from **nd\_status**.

Various statistics about data reception on the device are kept by the driver in the **ndd** structure. These statistics are part of the data returned by the **NDD\_GET\_STATS** and **NDD\_GET\_ALL\_STATS** control operations.



There is no specified entry point for this function. The device informs the device driver of a received packet using an interrupt. Upon determining that the interrupt was the result of a packet reception, the device driver's interrupt handler invoke the **rx\_handler** completion routine to perform the tasks mentioned above.

## Asynchronous Status

When a status event occurs on the device, the Ethernet device drivers build the appropriate status block and call the **nd\_status** function that is specified in the **ndd\_t** structure of the network device. The **nd\_status** function is part of a CDLI network demuxer.

The following status blocks are defined for the Ethernet device drivers.

**Note:** The following device drivers support the Device Connected status block:

- Gigabit Ethernet-SX PCI Adapter Device Driver (14100401)
- Gigabit Ethernet-SX PCI-X Adapter Device Driver (14106802)
- 10/100/1000 Base-T Ethernet PCI-X Adapter Device Driver (14106902)
- 2-Port Gigabit Ethernet-SX PCI-X Adapter (14108802)
- 2-Port 10/100/1000 Base-TX PCI-X Adapter (14108902)
- 4-Port 10/100/1000 Base-TX PCI-X Adapter (14101103)
- 2-Port Gigabit Ethernet-SX PCI-Express Adapter Device Driver (14103f03)
- 2-Port 10/100/1000 Base-TX PCI-Express Adapter Device Driver (14104003)
- 10 Gigabit Ethernet-SR PCI-X 2.0 DDR Adapter (1410eb02)
- 10 Gigabit Ethernet-LR PCI-X 2.0 DDR Adapter (1410ec02)
- 4-Port 10/100/1000 Base-TX PCI-Express Adapter (14106803)

The PCI Ethernet Device Driver (22100020) supports the Bad Packets status block.

### Bad Packets

When the a bad packet has been received by a device driver (and a user has requested bad packets), the following status block is returned by the device driver.

**code** Set to NDD\_BAD\_PKTS.

#### option[0]

Specifies the error status of the packet. These error numbers are defined in **<sys/cdli\_entuser.h>**.

#### option[1]

Pointer to the mbuf containing the bad packet.

#### option[]

The remainder of the status block can be used to return additional status information by the device driver.

**Note:** The user does *not* own the mbuf containing the bad packet. The user must copy the mbuf (and the status block information if necessary). The device driver frees the mbuf upon return from the **nd\_status** function.

### Device Connected

When the device is successfully connected to the network the following status block is returned by the device driver.

**code** Set to NDD\_CONNECTED.

#### option[]

The option fields are not used.

## Device Control Operations

The `ndd_ctl` entry point is used to provide device control functions.

### NDD\_GET\_STATS Device Control Operation

The `NDD_GET_STATS` command returns statistics concerning the network device. General statistics are maintained by the device driver in the `ndd_genstats` field in the `ndd_t` structure. The `ndd_specstats` field in the `ndd_t` structure is a pointer to media-specific and device-specific statistics maintained by the device driver. Both sets of statistics are directly readable at any time by those users of the device that can access them. This command provides a way for any of the users of the device to access the general and media-specific statistics.

The `arg` and `length` parameters specify the address and length in bytes of the area where the statistics are to be written. The length specified *must* be the exact length of the general and media-specific statistics.

**Note:** The `ndd_speclen` field in the `ndd_t` structure plus the length of the `ndd_genstats_t` structure is the required length. The device-specific statistics might change with each new release of the operating system, but the general and media-specific statistics are not expected to change.

The user should pass in the `ent_ndd_stats_t` structure as defined in `sys/cdli_entuser.h`. The driver fails a call with a buffer smaller than the structure.

The statistics that are returned contain statistics obtained from the device. If the device is inoperable, the statistics that are returned do not contain the current device statistics. The copy of the `ndd_flags` field can be checked to determine the state of the device.

### NDD\_MIB\_QUERY Device Control Operation

The `NDD_MIB_QUERY` operation is used to determine which device-specific MIBs are supported on the network device. The `arg` and `length` parameters specify the address and length in bytes of a device-specific `MIB` structure. The device driver fills every member of that structure with a flag indicating the level of support for that member. The individual `MIB` variables that are not supported on the network device are set to `MIB_NOT_SUPPORTED`. The individual `MIB` variables that can only be read on the network device are set to `MIB_READ_ONLY`. The individual `MIB` variables that can be read and set on the network device are set to `MIB_READ_WRITE`. The individual `MIB` variables that can only be set (not read) on the network device are set to `MIB_WRITE_ONLY`. These flags are defined in the `/usr/include/sys/ndd.h` file.

The `arg` parameter specifies the address of the `Ethernet_all_mib` structure. This structure is defined in the `/usr/include/sys/Ethernet_mibs.h` file.

### NDD\_MIB\_GET Device Control Operation

The `NDD_MIB_GET` operation is used to get all MIBs on the specified network device. The `arg` and `length` parameters specify the address and length in bytes of the device specific MIB structure. The device driver sets any unsupported variables to zero (nulls for strings).

If the device supports the RFC 1229 receive address object, the corresponding variable is set to the number of receive addresses currently active.

The `arg` parameter specifies the address of the `Ethernet_all_mib` structure. This structure is defined in the `/usr/include/sys/Ethernet_mibs.h` file.

### NDD\_ENABLE\_ADDRESS Device Control Operation

The `NDD_ENABLE_ADDRESS` command enables the receipt of packets with an alternate (for example, multicast) address. The `arg` and `length` parameters specify the address and length in bytes of the alternate address to be enabled. The `NDD_ALTADDRS` flag in the `ndd_flags` field is set.

The device driver verifies that if the address is a valid multicast address. If the address is not a valid multicast address, the operation fails with an **EINVAL** error. If the address is valid, the driver adds it to its multicast table and enable the multicast filter on the adapter. The driver keeps a reference count for each individual address. Whenever a duplicate address is registered, the driver simply increments the reference count of that address in its multicast table, no update of the adapter's filter is needed. There is a hardware limitation on the number of multicast addresses in the filter.

### **NDD\_DISABLE\_ADDRESS Device Control Operation**

The **NDD\_DISABLE\_ADDRESS** command disables the receiving packets with a specified alternate (for example, multicast) address. The *arg* and *length* parameters specify the address and length in bytes of the alternate address to be disabled. The **NDD\_ALTADDRS** flag in the **ndd\_flags** field is reset if this is the last alternate address.

The device driver verifies that if the address is a valid multicast address. If the address is not a valid multicast address, the operation fails with an **EINVAL** error. The device driver makes sure that the multicast address can be found in its multicast table. Whenever a match is found, the driver decrements the reference count of that individual address in its multicast table. If the reference count becomes 0, the driver deletes the address from the table and update the multicast filter on the adapter.

### **NDD\_MIB\_ADDR Device Control Operation**

The **NDD\_MIB\_ADDR** operation is used to get all the addresses for which the specified device accepts packets or frames. The *arg* parameter specifies the address of the **ndd\_mib\_addr\_t** structure. The length parameter specifies the length of the structure with the appropriate number of **ndd\_mib\_addr\_t.mib\_addr** elements. This structure is defined in the **/usr/include/sys/ndd.h** file. If the length is less than the length of the **ndd\_mib\_addr\_t** structure, the device driver returns **EINVAL**. If the structure is not large enough to hold all the addresses, the addresses that fit are still placed in the structure. The **ndd\_mib\_addr\_t.count** field is set to the number of addresses returned and **E2BIG** is returned.

One of the following address types is returned:

- Device physical address (or alternate address specified by user)
- Broadcast addresses
- Multicast addresses

### **NDD\_CLEAR\_STATS Device Control Operation**

The counters kept by the device are zeroed.

### **NDD\_GET\_ALL\_STATS Device Control Operation**

The **NDD\_GET\_ALL\_STATS** operation is used to gather all the statistics for the specified device. The **arg** parameter specifies the address of the statistics structure for the particular device type. The following structures are available:

- The **kent\_all\_stats\_t** structure is available for the PCI Ethernet Adapter Device Driver (22100020), and is defined in the **cdli\_entuser.h** include file.
- The **phxent\_all\_stats\_t** structure is available for the 10/100 Mbps Ethernet PCI Adapter Device Driver (23100020), and is defined in the device-specific **cdli\_entuser.phxent.h** include file.
- The **scent\_all\_stats\_t** structure is available for the 10/100 Mbps Ethernet PCI Adapter II Device Driver (1410ff01), and is defined in the device-specific **cdli\_entuser.scent.h** include file.
- The **gxent\_all\_stats\_t** structure is available for the Gigabit Ethernet-SX PCI Adapter Device Driver (14100401), and is defined in the device-specific **cdli\_entuser.gxent.h** include file.
- The **goent\_all\_stats\_t** structure is available for the Gigabit Ethernet-SX PCI-X Adapter Device Driver (14106802), the 10/100/1000 Base-T Ethernet PCI-X Adapter Device Driver (14106902), the 2-Port Gigabit Ethernet-SX PCI-X Adapter Device Driver (14108802), the 2-Port 10/100/1000 Base-T Ethernet PCI-X Adapter Device Driver (14108902), the 4-Port 10/100/1000 Base-TX PCI-X Adapter Device Driver (14101103), and the 4-Port 10/100/1000 Base-TX PCI-Express Adapter Device Driver (14106803), the

2-Port Gigabit Ethernet-SX PCI-Express Adapter Device Driver (14103f03), and the 2-Port 10/100/1000 Base-TX PCI-Express Adapter Device Driver (14104003), and is defined in the device-specific **cdli\_entuser.goent.h** include file.

- The **vent\_all\_stats\_t** structure is available for the 10 Gigabit Ethernet-SR PCI-X Adapter Device Driver (1410ba02) and the 10 Gigabit Ethernet-LR PCI-X Adapter Device Driver (1410bb02), and is defined in the device-specific **cdli\_entuser.vent.h** include file.
- The **bent\_all\_stats\_t** structure is available for the Gigabit Ethernet-SX Adapter Device Driver (e414a816), and is defined in the device-specific **cdli\_entuser.bent.h** include file.
- The **ment\_all\_stats\_t** structure is available for the Gigabit Ethernet-SX Adapter Device Driver (14101403) and the Gigabit Ethernet-SX PCI-X Adapter Device Driver (14106703), and is defined in the device-specific **cdli\_entuser.ment.h** include file.
- The **kngent\_all\_stats\_t** structure is available for the 10 Gigabit Ethernet-SR Adapter Device Driver (1410eb02) and the 10 Gigabit Ethernet-LR PCI-X 2.0 DDR Adapter (1410ec02), and is defined in the device-specific **cdli\_entuser.kngent.h** include file.
- The **hea\_all\_stats\_t** structure is available for the Host Ethernet Adapter Device Driver, and is defined in the **cdli\_entuser.hea.h** include file.

The statistics that are returned contain statistics obtained from the device. If the device is inoperable, the statistics that are returned do not contain the current device statistics. The copy of the **ndd\_flags** field can be checked to determine the state of the device.

### **NDD\_ENABLE\_MULTICAST Device Control Operation**

The **NDD\_ENABLE\_MULTICAST** command enables the receipt of packets with any multicast (or group) address. The *arg* and *length* parameters are not used. The **NDD\_MULTICAST** flag in the **ndd\_flags** field is set.

### **NDD\_DISABLE\_MULTICAST Device Control Operation**

The **NDD\_DISABLE\_MULTICAST** command disables the receipt of *all* packets with multicast addresses and only receives those packets whose multicast addresses were specified using the **NDD\_ENABLE\_ADDRESS** command. The *arg* and *length* parameters are not used. The **NDD\_MULTICAST** flag in the **ndd\_flags** field is reset only after the reference count for multicast addresses has reached zero.

### **NDD\_PROMISCUOUS\_ON Device Control Operation**

The **NDD\_PROMISCUOUS\_ON** command turns on promiscuous mode. The *arg* and *length* parameters are not used.

When the device driver is running in promiscuous mode, all network traffic is passed to the network demultiplexer. When the Ethernet device driver receives a valid packet from the network device, the Ethernet device driver calls the **nd\_receive** function that is specified in the **ndd\_t** structure of the network device. The **NDD\_PROMISC** flag in the **ndd\_flags** field is set. Promiscuous mode is considered to be valid packets only. See the **NDD\_ADD\_STATUS** command for information about how to request support for bad packets.

The device driver maintains a reference count on this operation. The device driver increments the reference count for each operation. When this reference count is equal to one, the device driver issues commands to enable the promiscuous mode. If the reference count is greater than one, the device driver does not issue any commands to enable the promiscuous mode.

### **NDD\_PROMISCUOUS\_OFF Device Control Operation**

The **NDD\_PROMISCUOUS\_OFF** command terminates promiscuous mode. The *arg* and *length* parameters are not used. The **NDD\_PROMISC** flag in the **ndd\_flags** field is reset.

The device driver maintains a reference count on this operation. The device driver decrements the reference count for each operation. When the reference count is not equal to zero, the device driver does

not issue commands to disable the promiscuous mode. After the reference count for this operation is equal to zero, the device driver issues commands to disable the promiscuous mode.

### **NDD\_DISABLE\_ADAPTER Device Control Operation**

**Note:** This device control operation is not supported on the PCI Ethernet Adapter Device Driver (22100020).

The **NDD\_DISABLE\_ADAPTER** operation is used by Etherchannel to disable the adapter so that it cannot transmit or receive data. During this operation the **NDD\_RUNNING** and **NDD\_LIMBO** flags are cleared and the adapter is reset. The *arg* and *len* parameters are not used.

### **NDD\_ENABLE\_ADAPTER Device Control Operation**

**Note:** This device control operation is not supported on the PCI Ethernet Adapter Device Driver (22100020).

The **NDD\_ENABLE\_ADAPTER** operation is used by Etherchannel to return the adapter to a running state so it can transmit and receive data. During this operation the adapter is started and the **NDD\_RUNNING** flag is set. The *arg* and *len* parameters are not used.

### **NDD\_SET\_LINK\_STATUS Device Control Operation**

**Note:** This device control operation is not supported on the PCI Ethernet Adapter Device Driver (22100020).

The **NDD\_SET\_LINK\_STATUS** operation is used by Etherchannel to pass the driver a function pointer and argument that will subsequently be called by the driver whenever the link status changes. The *arg* parameter contains a pointer to a **ndd\_sls\_t** structure, and the *len* parameter contains the length of the **ndd\_sls\_t** structure.

### **NDD\_SET\_MAC\_ADDR Device Control Operation**

**Note:** This device control operation is not supported on the PCI Ethernet Adapter Device Driver (22100020).

The **NDD\_SET\_MAC\_ADDR** operation is used by Etherchannel to set the adapters MAC address at runtime. The MAC address set by this ioctl is valid until another **NDD\_SET\_MAC\_ADDR** call is made with a new address or when the adapter is closed. If the adapter is closed, the previously configured MAC address. The MAC address configured with the ioctl supersedes any alternate address that might have been configured.

The *arg* argument is char [6], representing the MAC address that is configured on the adapter. The *len* argument is 6.

## **Trace**

For LAN device drivers, trace points enable error monitoring as well as tracking packets as they move through the driver. The drivers issue trace points for some or all of the following conditions:

- Beginning and ending of main functions in the main path
- Error conditions
- Beginning and ending of each function that is tracking buffers outside of the main path
- Debugging purposes (These trace points are only enabled when the driver is compiled with **-DDEBUG** turned on, and therefore the driver can contain as many of these trace points as necessary.)

The existing Ethernet device drivers each have either three or four trace points. The Trace Hook IDs the PCI Ethernet Adapter Device Driver (22100020) is defined in the **sys/cdli\_entuser.h** file. Other drivers have defined local **cdli\_entuser.driver.h** files with the Trace Hook definitions. For more information, see “Debug and Performance Tracing” on page 347.

Following is a list of trace hooks (and location of definition file) for the existing Ethernet device drivers.

### **PCI Ethernet Adapter Device Driver (22100020)**

Definition file: **cdli\_entuser.h**

Trace Hook IDs:

Transmit	-2A4
Receive	-2A5
Other	-2A6

### **10/100 Mbps Ethernet PCI Adapter Device Driver (23100020)**

Definition file: **cdli\_entuser.phxent.h**

Trace Hook IDs:

Transmit	-2E6
Receive	-2E7
Other	-2E8

### **10/100 Mbps Ethernet PCI Adapter II Device Driver (1410ff01)**

Definition file: **cdli\_entuser.scent.h**

Trace Hook IDs:

Transmit	-470
Receive	-471
Other	-472

### **Gigabit Ethernet-SX PCI Adapter Device Driver (14100401)**

Definition file: **cdli\_entuser.gxent.h**

Trace Hook IDs:

Transmit	-2EA
Receive	-2EB
Other	-2EC

### **Gigabit Ethernet-SX PCI-X Adapter Device Driver (14106802), 10/100/1000 Base-T Ethernet PCI-X Adapter Device Driver (14106902), 2-Port Gigabit Ethernet-SX PCI-X Adapter (14108802), 2-Port 10/100/1000 Base-TX PCI-X Adapter (14108902), 4-Port 10/100/1000 Base-TX PCI-X Adapter (14101103), 4-Port 10/100/1000 Base-TX PCI-Express Adapter (14106803), 2-Port Gigabit Ethernet-SX PCI-Express Adapter (14103f03), and 2-Port 10/100/1000 Base-TX PCI-Express Adapter (14104003)**

Definition file: **cdli\_entuser.goent.h**

Trace Hook IDs:

Transmit	-473
----------	------

Receive	-474
Other	-475

The device driver also has the following trace points to support the **netpmon** program:

<b>WQUE</b>	An output packet has been queued for transmission.
<b>WEND</b>	The output of a packet is complete.
<b>RDAT</b>	An input packet has been received by the device driver.
<b>RNOT</b>	An input packet has been given to the demuxer.
<b>REND</b>	The demultiplexer has returned.

## 10 Gigabit Ethernet-SR PCI-X Adapter (1410ba02) and 10 Gigabit Ethernet-LR PCI\_X Adapter (1410bb02)

Definition file: **cdli\_entuser.vent.h**

Trace Hook IDs:

Transmit	-598
Receive	-599
Other	-59A

The device driver also has the following trace points to support the **netpmon** program:

<b>WQUE</b>	An output packet has been queued for transmission.
<b>WEND</b>	The output of a packet is complete.
<b>RDAT</b>	An input packet has been received by the device driver.
<b>RNOT</b>	An input packet has been given to the demuxer.
<b>REND</b>	The demultiplexer has returned.

## Gigabit Ethernet-SX Adapter Device Driver (e414a816)

Definition file: **cdli\_entuser.bent.h**

Trace Hook IDs:

Transmit	-5B2
Receive	-5B3
Other	-5B4

Definition file: **cdli\_entuser.kngent.h**

Trace Hook IDs:

Transmit	-4a1
Receive	-4a2
Other	-4a3

The device driver also has the following trace points to support the **netpmon** program:

<b>WQUE</b>	An output packet has been queued for transmission.
<b>WEND</b>	The output of a packet is complete.

**RDAT** An input packet has been received by the device driver.

**RNOT** An input packet has been given to the demuxer.

**REND** The demultiplexer has returned.

## **Gigabit Ethernet-SX Adapter Device Driver (14101403) and Gigabit Ethernet-SX PCI-X Adapter Device Driver (14106703)**

Definition file: **cdli\_entuser.ment.h**

Trace Hook IDs:

Transmit	-4C5
Receive	-4C6
Other	-4C7

## **Host Ethernet Adapter Device Driver**

Definition file: **cdli\_entuser.heh.h**

Trace Hook IDs:

Transmit	0x5df
Receive	0x5e0
Other	0x5e1

## **Error Logging**

For error logging information, see “Error Logging” on page 344.

### **PCI Ethernet Adapter Device Driver (22100020)**

The Error IDs for the PCI Ethernet Adapter Device Driver (22100020) are as follows:

#### **ERRID\_KENT\_ADAP\_ERR**

Indicates that the adapter is not responding to initialization commands. User intervention is necessary to fix the problem.

#### **ERRID\_KENT\_RCVRY**

Indicates that the device driver detected a temporary adapter error requiring that it enter network recovery mode. It has reset the adapter in an attempt to fix the problem.

#### **ERRID\_KENT\_TX\_ERR**

Indicates that the device driver has detected a transmission error. User intervention is not required unless the problem persists.

#### **ERRID\_KENT\_PIO**

Indicates that the device driver has detected a program IO error. The device driver was unable to fix the problem. User intervention is necessary to fix the problem.

#### **ERRID\_KENT\_DOWN**

Indicates that the device driver has shut down the adapter due to an unrecoverable error. The adapter is no longer functional due to the error. The error that caused the device to shut down is error logged immediately before this error log entry. User intervention is necessary to fix the problem.

### **10/100 Mbps Ethernet PCI Adapter Device Driver (23100020)**

The Error IDs for the 10/100 Mbps Ethernet PCI Adapter Device Driver (23100020) are as follows:

#### **ERRID\_PHXENT\_ADAP\_ERR**

Indicates that the adapter is not responding to initialization commands. User-intervention is necessary to fix the problem.



**ERRID\_PHXENT\_ERR\_RCVRY**

Indicates that the device driver detected a temporary adapter error requiring that it enter network recovery mode. It has reset the adapter in an attempt to fix the problem.

**ERRID\_PHXENT\_TX\_ERR**

Indicates that the device driver has detected a transmission error. User-intervention is not required unless the problem persists.

**ERRID\_PHXENT\_PIO**

Indicates that the device driver has detected a program IO error. The device driver was unable to fix the problem. User intervention is necessary to fix the problem.

**ERRID\_PHXENT\_DOWN**

Indicates that the device driver has shutdown the adapter due to an unrecoverable error. The adapter is no longer functional due to the error. The error that caused the device shutdown is error logged immediately before this error log entry. User intervention is necessary to fix the problem.

**ERRID\_PHXENT\_EEPROM\_ERR**

Indicates that the device driver is in a defined state due to an invalid or bad EEPROM. The device driver does not become available. Hardware support should be contacted.

**ERRID\_PHXENT\_EEPROM2\_ERR**

Indicates that the device driver is in a defined state due to an invalid or bad EEPROM. The device driver does not become available. Hardware support should be contacted.

**ERRID\_PHXENT\_CLOSE\_ERR**

Indicates that an application is holding a private receive mbuf owned by the device driver during a close operation. User intervention is not required.

**ERRID\_PHXENT\_LINK\_ERR**

Indicates that the link between the adapter and the network switch is down. The device driver attempts to reestablish the connection after the physical link is reestablished. When the link is again established, the device driver logs **ERRID\_PHXENT\_ERR\_RCVRY**. User intervention is necessary to fix the problem.

**Gigabit Ethernet-SX PCI Adapter Device Driver (14100401)**

The Error IDs for the Gigabit Ethernet-SX PCI Adapter Device Driver (14100401) are as follows:

**ERRID\_GXENT\_ADAP\_ERR**

Indicates that the adapter failed initialization commands. User intervention is necessary to fix the problem.

**ERRID\_GXENT\_CMD\_ERR**

Indicates that the device driver has detected an error while issuing commands to the adapter. The device driver enters an adapter recovery mode where it attempts to recover from the error. If the device driver is successful, it logs **ERRID\_GXENT\_RCVRY\_EXIT**. User intervention is not necessary for this error unless the problem persists.

**ERRID\_GXENT\_DOWNLOAD\_ERR**

Indicates that an error occurred while downloading firmware to the adapter. User intervention is necessary to fix the problem.

**ERRID\_GXENT\_EEPROM\_ERR**

Indicates that an error occurred while reading the adapter EEPROM. User intervention is necessary to fix the problem.

**ERRID\_GXENT\_LINK\_DOWN**

Indicates that the link between the adapter and the network switch is down. The device driver attempts to reestablish the connection after the physical link is reestablished. When the link is again established, the device driver logs **ERRID\_GXENT\_RCVRY\_EXIT**. User intervention is necessary to fix the problem.

**ERRID\_GXENT\_RCVRY\_EXIT**

Indicates that a temporary error (link down, command error, or transmission error) has been corrected.

**ERRID\_GXENT\_TX\_ERR**

Indicates that the device driver has detected a transmission error. The device driver enters an adapter recovery mode in an attempt to recover from the error. If the device driver is successful, it logs **ERRID\_GXENT\_RCVRY\_EXIT**. User intervention is not necessary for this error unless the problem persists.

**ERRID\_GXENT\_EEH\_SERVICE\_ERR**

Indicates that the device driver has detected an error during an attempt to recover from a PCI bus error. User intervention is necessary to fix the problem.

**10/100 Mbps Ethernet PCI Adapter II Device Driver (1410ff01)**

The Error IDs for the 10/100 Mbps Ethernet PCI Adapter II Device Driver (1410ff01) are as follows:

**ERRID\_SCENT\_ADAP\_ERR**

Indicates that the adapter failed initialization commands. User intervention is necessary to fix the problem.

**ERRID\_SCENT\_PIO\_ERR**

Indicates that the device driver has detected a program IO error. The device driver was unable to fix the problem. User intervention is necessary to fix the problem.

**ERRID\_SCENT\_EEPROM\_ERR**

Indicates that an error occurred while reading the adapter EEPROM. User intervention is necessary to fix the problem.

**ERRID\_SCENT\_LINK\_DOWN**

Indicates that the link between the adapter and the network switch is down. The device driver attempts to reestablish the connection after the physical link is reestablished. When the link is again established, the device driver logs **ERRID\_SCENT\_RCVRY\_EXIT**. User intervention is necessary to fix the problem.

**ERRID\_SCENT\_RCVRY\_EXIT**

Indicates that a temporary error (link down, command error, or transmission error) has been corrected.

**ERRID\_SCENT\_TX\_ERR**

Indicates that the device driver has detected a transmission error. The device driver enters an adapter recovery mode in an attempt to recover from the error. If the device driver is successful, it logs **ERRID\_SCENT\_RCVRY\_EXIT**. User intervention is not necessary for this error unless the problem persists.

**ERRID\_SCENT\_EEH\_SERVICE\_ERR**

Indicates that the device driver has detected an error during an attempt to recover from a PCI bus error. User intervention is necessary to fix the problem.

**Gigabit Ethernet-SX PCI-X Adapter Device Driver (14106802), 10/100/1000 Base-T Ethernet PCI-X Adapter Device Driver (14106902), 2-Port Gigabit Ethernet-SX PCI-X Adapter (14108802), 2-Port 10/100/1000 Base-TX PCI-X Adapter (14108902), 4-Port 10/100/1000 Base-TX PCI-X Adapter Device Driver (14101103), 4-Port 10/100/1000 Base-TX PCI-Express Adapter Device Driver (14106803), 2-Port Gigabit Ethernet-SX PCI-Express Adapter (14103f03), and 2-Port 10/100/1000 Base-TX PCI-Express Adapter (14104003)**

The Error IDs for the Gigabit Ethernet-SX PCI-X Adapter Device Driver (14106802), the 10/100/1000 Base-T Ethernet PCI-X Adapter Device Driver (14106902), the 2-Port Gigabit Ethernet-SX PCI-X Adapter Device Driver (14108802), 2-Port 10/100/1000 Base-TX PCI-X Adapter Device Driver (14108902), the 4-Port 10/100/1000 Base-TX PCI-X Adapter Device Driver (14101103), the 4-Port 10/100/1000 Base-TX

PCI-Express Adapter Device Driver (14106803), the 2-Port Gigabit Ethernet-SX PCI-Express Adapter (14103f03), and the 2-Port 10/100/1000 Base-TX PCI-Express Adapter (14104003) are as follows:

**ERRID\_GOENT\_ADAP\_ERR**

Indicates that the adapter failed initialization commands. User intervention is necessary to fix the problem.

**ERRID\_GOENT\_PIO\_ERR**

Indicates that the device driver has detected a program I/O error. The device driver was unable to fix the problem. User intervention is necessary to fix the problem.

**ERRID\_GOENT\_EEPROM\_ERR**

Indicates that an error occurred while reading the adapter EEPROM. User intervention is necessary to fix the problem.

**ERRID\_GOENT\_LINK\_DOWN**

Indicates that the link between the adapter and the network switch is down. The device driver attempts to reestablish the connection after the physical link is reestablished. When the link is again established, the device driver logs **ERRID\_GOENT\_RCVRY\_EXIT**. User intervention is necessary to fix the problem.

**ERRID\_GOENT\_RCVRY\_EXIT**

Indicates that a temporary error (link down, command error, or transmission error) has been corrected.

**ERRID\_GOENT\_TX\_ERR**

Indicates that the device driver has detected a transmission error. The device driver enters an adapter recovery mode in an attempt to recover from the error. If the device driver is successful, it logs **ERRID\_GOENT\_RCVRY\_EXIT**. User intervention is not necessary for this error unless the problem persists.

**ERRID\_GOENT\_EEH\_SERVICE\_ERR**

Indicates that the device driver has detected a error during an attempt to recover from a PCI bus error. User intervention is necessary to fix the problem.

**10 Gigabit Ethernet-SR PCI-X Adapter (1410ba02) and 10 Gigabit Ethernet-LR PCI\_X Adapter (1410bb02)**

The error IDs for the 10 Gigabit Ethernet-SR PCI-X Adapter (1410ba02) and 10 Gigabit Ethernet-LR PCI\_X Adapter (1410bb02) are as follows:

**ERRID\_VENT\_ADAP\_ERR**

Indicates that the adapter failed initialization commands. User intervention is necessary to fix the problem.

**ERRID\_VENT\_PIO\_ERR**

Indicates that the device driver has detected a program I/O error. The device driver was unable to fix the problem. User intervention is necessary to fix the problem.

**ERRID\_VENT\_EEPROM\_ERR**

Indicates that an error occurred while reading the adapter EEPROM. User intervention is necessary to fix the problem.

**ERRID\_VENT\_LINK\_DOWN**

Indicates that the link between the adapter and the network switch is down. The device driver attempts to reestablish the connection after the physical link is reestablished. When the link is again established, the device driver logs **ERRID\_VENT\_RCVRY\_EXIT**. User intervention is necessary to fix the problem.

**ERRID\_VENT\_RCVRY\_EXIT**

Indicates that a temporary error (link down, command error, or transmission error) has been corrected.

**ERRID\_VENT\_TX\_ERR**

Indicates that the device driver has detected a transmission error. The device driver enters an adapter recovery mode in an attempt to recover from the error. If the device driver is successful, it logs **ERRID\_VENT\_RCVRY\_EXIT**. User intervention is not necessary for this error unless the problem persists.

**ERRID\_VNT\_EEH\_SERVICE\_ERR**

Indicates that the device driver has detected an error during an attempt to recover from a PCI bus error. User intervention is necessary to fix the problem.

**Gigabit Ethernet-SX Adapter Device Driver (e414a816)**

The Error IDs for the Gigabit Ethernet-SX Adapter Device Driver (e414a816) are as follows:

**ERRID\_BENT\_ADAP\_ERR**

Indicates that the adapter failed initialization commands. User intervention is necessary to fix the problem.

**ERRID\_BENT\_DOWNLOAD\_ERR**

Indicates that an error occurred while downloading firmware to the adapter. User intervention is necessary to fix the problem.

**ERRID\_BENT\_EEPROM\_ERR**

Indicates that an error occurred while reading the adapter EEPROM. User intervention is necessary to fix the problem.

**ERRID\_BENT\_LINK\_DOWN**

Indicates that the link between the adapter and the network switch is down. The device driver attempts to reestablish the connection after the physical link is reestablished. When the link is again established, the device driver logs **ERRID\_BENT\_RCVRY\_EXIT**. User intervention is necessary to fix the problem.

**ERRID\_BENT\_RCVRY\_EXIT**

Indicates that a temporary error (link down, or transmission error) was corrected.

**ERRID\_BENT\_TX\_ERR**

Indicates that the device driver has detected a transmission error. The device driver enters an adapter recovery mode in an attempt to recover from the error. If the device driver is successful, it logs **ERRID\_BENT\_RCVRY\_EXIT**. User intervention is not necessary for this error unless the problem persists.

**Gigabit Ethernet-SX Adapter Device Driver (14101403) and Gigabit Ethernet-SX PCI-X Adapter Device Driver (14106703)**

The Error IDs for the Gigabit Ethernet-SX Adapter Device Driver (14101403) and the Gigabit Ethernet-SX PCI-X Adapter Device Driver (14106703) are as follows:

**ERRID\_MENT\_ADAP\_ERR**

Indicates that the adapter failed initialization commands. User intervention is necessary to fix the problem.

**ERRID\_MENT\_DOWNLOAD\_ERR**

Indicates that an error occurred while downloading firmware to the adapter. User intervention is necessary to fix the problem.

**ERRID\_MENT\_EEPROM\_ERR**

Indicates that an error occurred while reading the adapter EEPROM. User intervention is necessary to fix the problem.

**ERRID\_MENT\_LINK\_DOWN**

Indicates that the link between the adapter and the network switch is down. The device driver attempts to reestablish the connection after the physical link is reestablished. When the link is again established, the device driver logs **ERRID\_MENT\_RCVRY\_EXIT**. User intervention is necessary to fix the problem.

**ERRID\_MENT\_RCVRY\_EXIT**

Indicates that a temporary error (link down, or transmission error) was corrected.

**ERRID\_MENT\_TX\_ERR**

Indicates that the device driver has detected a transmission error. The device driver enters an adapter recovery mode in an attempt to recover from the error. If the device driver is successful, it logs **ERRID\_MENT\_RCVRY\_EXIT**. User intervention is not necessary for this error unless the problem persists.

**10 Gigabit Ethernet-SR PCI-X 2.0 DDR Adapter (1410eb02) and 10 Gigabit Ethernet-LR PCI\_X 2.0 DDR Adapter (1410ec02)**

The error IDs for the 10 Gigabit Ethernet-SR PCI-X 2.0 DDR Adapter (1410eb02) and 10 Gigabit Ethernet-LR PCI\_X 2.0 DDR Adapter (1410ec02) are as follows:

**ERRID\_KNGENT\_ADAP\_ERR**

Indicates that the adapter failed initialization commands. User intervention is necessary to fix the problem.

**ERRID\_KNGENT\_PIO\_ERR**

Indicates that the device driver has detected a program I/O error. The device driver was unable to fix the problem. User intervention is necessary to fix the problem.

**ERRID\_KNGENT\_EEPROM\_ERR**

Indicates that an error occurred while reading the adapter **EEPROM**. User intervention is necessary to fix the problem.

**ERRID\_KNGENT\_LINK\_DOWN**

Indicates that the link between the adapter and the network switch is down. The device driver attempts to reestablish the connection after the physical link is reestablished. When the link is again established, the device driver logs **ERRID\_KNGENT\_RCVRY\_EXIT**. User intervention is necessary to fix the problem.

**ERRID\_KNGENT\_RCVRY\_EXIT**

Indicates that a temporary error (link down, command error, or transmission error) has been corrected.

**ERRID\_KNGENT\_TX\_ERR**

Indicates that the device driver has detected a transmission error. The device driver enters an adapter recovery mode in an attempt to recover from the error. If the device driver is successful, it logs **ERRID\_KNGENT\_RCVRY\_EXIT**. User intervention is not necessary for this error unless the problem persists.

**ERRID\_KNGENT\_EEH\_SERVICE\_ERR**

Indicates that the device driver has detected an error during an attempt to recover from a PCI bus error. User intervention is necessary to fix the problem.

**Host Ethernet Adapter device driver**

The error IDs for the Host Ethernet adapter device driver are as follows:

**ERRID\_KNGENT\_EEH\_SERVICE\_ERR**

Indicates that the adapter had a permanent hardware failure. This error is irrecoverable. This error is logged when a notification-event-queue event, for example, **ADAPTER\_MALFUNCTION**, occurred. User intervention is necessary to fix the problem.

**ERRID\_HEA\_ENS\_PORT\_ERR**

Indicates that the external network logical switch has a temporary port error or a link down. The adapter and other ports of the adapter, if any, remain available. This event is notified to the Host Ethernet Adapter device driver through a notification-event-queue event of the **HEA\_ENS\_PORT\_MALFUNCTION** type. User intervention is necessary to fix the problem.

**ERRID\_HEA\_LINK\_DOWN**

Indicates that the adapter has detected that the Ethernet link is down. This event is notified to the

Host Ethernet Adapter device driver through a notification-event-queue event that is `PORT_STATE_CHANGE`. When the link is re-established, the device driver logs the `ERRID_HEA_ERR_RCVRY_EXIT` error. User intervention is necessary to fix the problem.

#### **ERRID\_HEA\_ERR\_RCVRY\_EXIT**

Indicates that the adapter encountered a temporary error (`TX_ERR` or `LINK_DOWN`) that halted network activity and the problem has been resolved.

#### **ERRID\_HEA\_TX\_ERR**

Indicates that the device driver has detected a transmission error (transmit timeout). If the recovery is successful, the driver logs the `ERRID_HEA_ERR_RCVRY_EXIT` error; otherwise the `ERRID_HEA_ADAP_ERR` error is logged. User intervention is not required unless the `ERRID_HEA_ADAP_ERR` error is also logged.

---

## **Related Information**

“Common Communications Status and Exception Codes” on page 109.

“Logical File System Kernel Services” on page 63.

System Management Interface Tool (SMIT) in *AIX Version 6.1 General Programming Concepts: Writing and Debugging Programs*.

Error Logging Overview in *AIX Version 6.1 General Programming Concepts: Writing and Debugging Programs*.

Status Blocks for the Serial Optical Link Device Driver, Sense Data for the Serial Optical Link Device Driver in *AIX Version 6.1 Technical Reference: Kernel and Subsystems Volume 2*.

## **Subroutine References**

The `readx` subroutine in *AIX Version 6.1 Technical Reference: Base Operating System and Extensions Volume 2*.

## **Commands References**

The `entstat` Command in *AIX Version 6.1 Commands Reference, Volume 1*.

The `lecstat` Command, `mpcstat` Command in *AIX Version 6.1 Commands Reference, Volume 3*.

The `tokstat` Command in *AIX Version 6.1 Commands Reference, Volume 5*.

## **Technical References**

The `ddwrite` entry point, `ddselect` entry point in *AIX Version 6.1 Technical Reference: Kernel and Subsystems Volume 2*.

The `CIO_GET_STAT` operation, `CIO_HALT` operation, `CIO_START` operation in *AIX Version 6.1 Technical Reference: Kernel and Subsystems Volume 2*.

The `mpconfig Multiprotocol (MPQP) Device Handler Entry Point`, `mpwrite Multiprotocol (MPQP) Device Handler Entry Point`, `mpread Multiprotocol (MPQP) Device Handler Entry Point`, `mpmpx Multiprotocol (MPQP) Device Handler Entry Point`, `mpopen Multiprotocol (MPQP) Device Handler Entry Point`, `mpselect Multiprotocol (MPQP) Device Handler Entry Point`, `mpclose Multiprotocol (MPQP) Device Handler Entry Point`, `mpioctl Multiprotocol (MPQP) Device Handler Entry Point` in *AIX Version 6.1 Technical Reference: Kernel and Subsystems Volume 2*.

---

## Chapter 8. Graphic Input Devices Subsystem

The graphic input devices subsystem includes the keyboard/sound, mouse, tablet, dials, and lighted programmable-function keys (LPGK) devices. These devices provide operator input primarily to graphic applications. However, the keyboard can provide system input by means of the console.

The program interface to the input device drivers is described in the **inputdd.h** header file. This header file is available as part of the **bos.adt.graphics** filesset.

---

### open and close Subroutines

An **open** subroutine call is used to create a channel between the caller and a graphic input device driver. The keyboard supports two such channels. The most recently created channel is considered the active channel. All other graphic input device drivers support only one channel. The **open** subroutine call is processed normally, except that the *OFLAG* and *MODE* parameters are ignored. The keyboard provides support for the **fp\_open** subroutine call; however, only one kernel mode channel can be open at any given time. The **fp\_open** subroutine call returns EACCES for all other graphic input devices.

The **close** subroutine is used to remove a channel created by the **open** subroutine call.

---

### read and write Subroutines

The graphic input device drivers do not support read or write operations. A read or write to a graphic input device special file behaves as if a read or write was made to **/dev/null**.

---

### ioctl Subroutines

The ioctl operations provide run-time services. The special files support the following ioctl operations:

- Keyboard
- Mouse
- Tablet
- GIO (Graphics I/O) Adapter
- Dials
- LPGK

### Keyboard

<b>IOCINFO</b>	Returns the <b>devinfo</b> structure.
<b>KSQUERYID</b>	Queries the keyboard device identifier.
<b>KSQUERYSV</b>	Queries the keyboard service vector.
<b>KSREGRING</b>	Registers the input ring.
<b>KSRFLUSH</b>	Flushes the input ring.
<b>KSLED</b>	Sets and resets the keyboard LEDs.
<b>KSCFGCLICK</b>	Configures the clicker.
<b>KSVOLUME</b>	Sets the alarm volume.
<b>KSALARM</b>	Sounds the alarm.
<b>KSTRATE</b>	Sets the repeat rate.
<b>KSTDELAY</b>	Sets the delay before repeat.
<b>KSKAP</b>	Enables and disables the keep-alive poll.
<b>KSKAPACK</b>	Acknowledges the keep-alive poll.
<b>KSDIAGMODE</b>	Enables and disables the diagnostics mode.

## Note:

1. A nonactive channel processes only **IOCINFO**, **KSQUERYID**, **KSQUERYSV**, **KSREGRING**, **KSRFLUSH**, **KSKAP**, and **KSKAPACK**. All other ioctl subroutine calls are ignored without error.
2. The **KSLED**, **KSCFGCLICK**, **KSVOLUME**, **KSALARM**, **KSTRATE**, and **KSTDELAY** ioctl subroutine calls return an **EBUSY** error in the **errno** global variable when the keyboard is in diagnostics mode.
3. The **KSQUERYSV** ioctl subroutine call is only available when the channel is open from kernel mode (with the **fp\_open** kernel service).
4. The **KSKAP**, **KSKAPACK**, **KSDIAGMODE** ioctl subroutine calls are only available when the channel is open from user mode.

## Mouse

<b>IOCINFO</b>	Returns the <b>devinfo</b> structure.
<b>MQUERYID</b>	Queries the mouse device identifier.
<b>MREGRING</b>	Registers the input ring.
<b>MRFLUSH</b>	Flushes the input ring.
<b>MTHRESHOLD</b>	Sets the mouse reporting threshold.
<b>MRESOLUTION</b>	Sets the mouse resolution.
<b>MSCALE</b>	Sets the mouse scale.
<b>MSAMPLERATE</b>	Sets the mouse sample rate.

## Tablet

<b>IOCINFO</b>	Returns the <b>devinfo</b> structure.
<b>TABQUERYID</b>	Queries the tablet device identifier.
<b>TABREGRING</b>	Registers the input ring.
<b>TABFLUSH</b>	Flushes the input ring.
<b>TABCONVERSION</b>	Sets the tablet conversion mode.
<b>TABRESOLUTION</b>	Sets the tablet resolution.
<b>TABORIGIN</b>	Sets the tablet origin.
<b>TABSAMPLERATE</b>	Sets the tablet sample rate.
<b>TABDEADZONE</b>	Sets the tablet dead zones.

## GIO (Graphics I/O) Adapter

<b>IOCINFO</b>	Returns the <b>devinfo</b> structure.
<b>GIOQUERYID</b>	Returns the ID of the attached devices.

## Dials

<b>IOCINFO</b>	Returns the <b>devinfo</b> structure.
<b>DIALREGRING</b>	Registers the input ring.
<b>DIALRFLUSH</b>	Flushes the input ring.
<b>DIALSETGRAND</b>	Sets the dial granularity.

## LPFK

<b>IOCINFO</b>	Returns the <b>devinfo</b> structure.
<b>LPFKREGRING</b>	Registers the input ring.
<b>LPFKRFLUSH</b>	Flushes the input ring.



---

## Input Ring

Data is obtained from graphic input devices by way of a circular First-In First-Out (FIFO) queue or input ring, rather than with a **read** subroutine call. The memory address of the input ring is registered with an **ioctl** (or **fp\_ioctl**) subroutine call. The program that registers the input ring is the owner of the ring and is responsible for allocating, initializing, and freeing the storage associated with the ring. The same input ring can be shared by multiple devices.

The input ring consists of the input ring header followed by the reporting area. The input ring header contains the reporting area size, the head pointer, the tail pointer, the overflow flag, and the notification type flag. Before registering an input ring, the ring owner must ensure that the head and tail pointers contain the starting address of the reporting area. The overflow flag must also be cleared and the size field set equal to the number of bytes in the reporting area. After the input ring has been registered, the owner can modify only the head pointer and the notification type flag.

Data stored on the input ring is structured as one or more event reports. Event reports are placed at the tail of the ring by the graphic input device drivers. Previously queued event reports are taken from the head of the input ring by the owner of the ring. The input ring is empty when the head and tail locations are the same. An overflow condition exists if placement of an event on the input ring would overwrite data that has not been processed. Following an overflow, new event reports are not placed on the input ring until the input ring is flushed via an **ioctl** subroutine or service vector call.

The owner of the input ring is notified when an event is available for processing via a SIGMSG signal or via callback if the channel was created by an **fp\_open** subroutine call. The notification type flag in the input ring header specifies whether the owner should be notified each time an event is placed on the ring or only when an event is placed on an empty ring.

## Management of Multiple Keyboard Input Rings

When multiple keyboard channels are opened, keyboard events are placed on the input ring associated with the most recently opened channel. When this channel is closed, the alternate channel is activated and keyboard events are placed on the input ring associated with that channel.

## Event Report Formats

Each event report consists of an identifier followed by the report size in bytes, a time stamp (system time in milliseconds), and one or more bytes of device-dependent data. The value of the identifier is specified when the input ring is registered. The program requesting the input-ring registration is responsible for identifier uniqueness within the input-ring scope.

**Note:** Event report structures are placed on the input-ring without spacing. Data wraps from the end to the beginning of the reporting area. A report can be split on any byte boundary into two non-contiguous sections.

The event reports are as follows:

### Keyboard

ID	Specifies the report identifier.
Length	Specifies the report length.
Time stamp	Specifies the system time (in milliseconds).
Key position code	Specifies the key position code.
Key scan code	Specifies the key scan code.
Status flags	Specifies the status flags.

## Tablet

ID	Specifies the report identifier.
Length	Specifies the report length.
Time stamp	Specifies the system time (in milliseconds).
Absolute X	Specifies the absolute X coordinate.
Absolute Y	Specifies the absolute Y coordinate.

## LPFK

ID	Specifies the report identifier.
Length	Specifies the report length.
Time stamp	Specifies the system time (in milliseconds).
Number of key pressed	Specifies the number of the key pressed.

## Dials

ID	Specifies the report identifier.
Length	Specifies the report length.
Time stamp	Specifies the system time (in milliseconds).
Number of dial changed	Specifies the number of the dial changed.
Delta change	Specifies delta dial rotation.

## Mouse (Standard Format)

ID	Specifies the report identifier.
Length	Specifies the report length.
Time stamp	Specifies the system time (in milliseconds).
Delta X	Specifies the delta mouse motion along the X axis.
Delta Y	Specifies the delta mouse motion along the Y axis.
Button status	Specifies the button status.

## Mouse (Extended Format)

ID	Specifies the report identifier.
Length	Specifies the report length.
Time stamp	Specifies the system time (in milliseconds).
Format	Specifies the format of additional fields.

### Format 1:

- **Status:** Specifies the extended button status
- **Delta Wheel:** Specifies the delta wheel movement

### Format 2:

- **Button Status:** Specifies the button status.
- **Delta X:** Specifies the delta mouse motion along the X axis.
- **Delta Y:** Specifies the delta mouse motion along the Y axis.
- **Delta Wheel:** Specifies the delta wheel movement

## Keyboard Service Vector

The keyboard service vector provides a limited set of keyboard-related and sound-related services for kernel extensions. The following services are available:

- Sound alarm
- Enable and disable secure attention key (SAK)
- Flush input queue

The address of the service vector is obtained with the `fp_ioctl` subroutine call during a non-critical period. The kernel extension can later invoke the service using an indirect call as follows:

```
(*ServiceVector[ServiceNumber]) (dev_t DeviceNumber, caddr_t Arg);
```

where:

- The service vector is a pointer to the service vector obtained by the **KSQUERYSV** `fp_ioctl` subroutine call.
- The *ServiceNumber* parameter is defined in the **inputdd.h** file.
- The *DeviceNumber* parameter specifies the major and minor numbers of the keyboard.
- The *Arg* parameter points to a **ksalarm** structure for alarm requests and a **uint** variable for SAK enable and disable requests. The *Arg* parameter is NULL for flush queue requests.

If successful, the function returns a value of 0 is returned. Otherwise, the function returns an error number defined in the **errno.h** file. Flush-queue and enable/disable-SAK requests are always processed, but alarm requests are ignored if the kernel extension's channel is inactive.

The following example uses the service vector to sound the alarm:

```
/* pinned data structures */
/* This example assumes that pinning is done elsewhere. */
int (**ksvtbl) ();
struct ksalarm alarm;
dev_t devno;

/* get address of service vector */
/* This should be done in a noncritical section */
if (fp_ioctl(fp, KSQUERYSV, &ksvtbl, 0)) {
    /* error recovery */
}
.
.
.

/* critical section */
/* sound alarm for 1 second using service vector */
alarm.duration = 128;
alarm.frequency = 100;

if ((*ksvtbl[KSVALARMS]) (devno, &alarm)) {
    /* error recovery */
}
}
```

## Special Keyboard Sequences

Special keyboard sequences are provided for the Secure Attention Key (SAK) and the Keep Alive Poll (KAP).

### Secure Attention Key

The user requests a secure shell by keying a secure attention. The keyboard driver interprets the key sequence CTRL x r as the SAK. An indirect call using the keyboard service vector enables and disables the detection of this key sequence. If detection of the SAK is enabled, a SAK causes the SAK callback to

be invoked. The SAK callback is invoked even if the input ring is inactive due to a user process issuing an open to the keyboard special file. The SAK callback runs within the interrupt environment.

### **Keep Alive Poll**

The keyboard device driver supports a special key sequence that kills the process that owns the keyboard. This sequence must first be defined with a **KSKAP** ioctl operation. After this sequence is defined, the keyboard device driver sends a **SIGKAP** signal to the process that owns the keyboard when the special sequence is entered on the keyboard. The process that owns the keyboard must acknowledge the **KSKAP** signal with a **KSKAPACK** ioctl within 30 seconds or the keyboard driver will terminate the process with a **SIGKILL** signal. The KAP is enabled on a per-channel basis and is unavailable if the channel is owned by a kernel extension.

---

## Chapter 9. Low Function Terminal Subsystem

This chapter discusses the following topics:

- Low Function Terminal Interface Functional Description
- Components Affected by the Low Function Terminal Interface
- Accented Characters

The low function terminal (lft) interface is a pseudo-device driver that interfaces with device drivers for the system keyboard and display adapters. The lft interface adheres to all standard requirements for pseudo-device drivers and has all the entry points and configuration code as required by the device driver architecture. This section gives a high-level description of the various configuration methods and entry points provided by the lft interface.

All the device drivers controlled by the lft interface are also used by AIXwindows. Consequently, along with the functions required for the tty subsystem interface, the lft interface provides the functions required by AIXwindows interfaces with display device driver adapters.

---

### Low Function Terminal Interface Functional Description

This section covers the lft interface functional description:

- Configuration
- Terminal Emulation
- IOCTLs Needed for AIXwindows Support
- Low Function Terminal to System Keyboard Interface
- Low Function Terminal to Display Device Driver Interface
- Low Function Terminal Device Driver Entry Points

### Configuration

The lft interface uses the common define, undefine, and unconfiguration methods standard for most devices.

**Note:** The lft interface does not support any change method for dynamically changing the lft configuration. Instead, use the **-P** flag with the **chdev** command. The changes become effective the next time the lft interface is configured.

The configuration process for the lft opens all display device drivers. To define the default display and console, select the default display and console during the console configuration process. If a graphics display is chosen as the system console, it automatically becomes the default display. The lft interface displays text on the default display.

The configuration process for the lft interface queries the ODM database for the available fonts and software keyboard map for the current session.

### Terminal Emulation

The lft interface is a stream-based tty subsystem. The lft interface provides VT100 (or IBM 3151) terminal emulation for the standard part of the ANSI 3.64 data stream. All line discipline handling is performed in the layers above the lft interface. The lft interface does not support virtual terminals.

The lft interface supports multiple fonts to handle the different screen sizes and resolutions necessary in providing a 25x80 character display on various display adapters.

**Note:** Applications requiring hft extensions need to use aixterm.

## IOCTLS Needed for AIXwindows Support

AIXwindows and the lft interface share the system keyboard and display device drivers. To prevent screen and keyboard inconsistencies, a set of ioctl coordinates usage between AIXwindows and the lft interface. On a system with multiple displays, the lft interface can still use the default display as long as AIXwindows is using another display.

**Note:** The lft interface provides ioctl support to set and change the default display.

## Low Function Terminal to System Keyboard Interface

The lft interface with the system keyboard uses an input ring mechanism. The details of the keyboard driver ioctls, as well as the format and description of this input ring, are provided in Chapter 8, “Graphic Input Devices Subsystem,” on page 197. The keyboard device driver passes raw keystrokes to the lft interface. These keystrokes are converted to the appropriate code point using keyboard tables. The use of keyboard-language-dependent keyboard tables ensures that the lft interface provides National Language Support.

## Low Function Terminal to Display Device Driver Interface

The lft uses a device independent interface known as the virtual display driver (vdd) interface. Because the lft interface has no virtual terminal or monitor mode support, some of the vdd entry points are not used by the lft.

The display drivers might enqueue font request through the font process started during lft initialization. The font process pins and unpins the requested fonts for **DMA** to the display adapter.

## Low Function Terminal Device Driver Entry Points

The lft interface supports the open, close, read, write, ioctl, and configuration entry points.

---

## Components Affected by the Low Function Terminal Interface

The lft interface impacts the following components:

- Configuration User Commands
- Keyboard Device Driver (Information about this is contained in Graphic Input Device Driver Programming Interface.)
- Display Device Driver
- Rendering Context Manager

## Configuration User Commands

The lft interface is a pseudo-device driver. Consequently, the system configuration process does not detect the lft interface as it does an adapter. The system provides for pseudo-device drivers to be started through **Config\_Rules**. To start the lft interface, use the **startlft** program.

Supported commands include:

- **lsfont**
- **mkfont**
- **chfont**
- **lskbd**
- **chkbd**
- **lsdisp** (see note)
- **chdisp** (see note)

**Note:**

1. *lstdisp* outputs the logical device name instead of the instance number.
2. *chdisp* uses the *ioctl* interface to the *lft* to set the requested display.

## Display Device Driver

Beginning with AIX 4.1, a display device driver is required for each supported display adapter.

The display device drivers provide all the standard interfaces (such as *config*, *initialize*, *terminate*, and *so forth*) required in any AIX 4.1 (or later) device drivers. The only device switch table entries supported are *open*, *close*, *config*, and *ioctl*. All other device switch table entries are set to *nodev*. In addition, the display device drivers provide a set of *ioctls* for use by *AIXwindows* and *diagnostics* to perform device specific functions such as *get bus access*, *bus memory address*, *DMA operations*, and *so forth*.

## Rendering Context Manager

The Rendering Context Manager (RCM) is a loadable module.

**Note:** Previously, the high functional terminal interface provided *AIXwindows* with the ***gsc\_handle***. This handle is used in all of the ***aixgsc*** system calls. The RCM provides this service for the *lft* interface.

To ensure that *lft* can recover the display in case *AIXwindows* should terminate abnormally, *AIXwindows* issues the *ioctl* to RCM after opening the pseudo-device. RCM passes on the *ioctl* to the *lft*. This way, the *close* function in RCM is invoked (Because *AIXwindows* is the only application that has opened RCM), and RCM notifies the *lft* interface to start reusing the display. To support this communication, the RCM provides the required *ioctl* support.

### The RCM to lft Interface Initialization

1. RCM performs the *open /dev/lft*.
2. Upon receiving a list of displays from X, RCM passes the information to the *lft* through an *ioctl*.
3. RCM resets the adapter.

### If AIXwindows Terminates Abnormally

1. RCM receives notification from X about the displays it was using.
2. RCM resets the adapter.
3. RCM passes the information to the *lft* by way of an *ioctl*.

### AIXwindows to lft Initialization

The *AIXwindows* to *lft* initialization includes the following:

1. *AIXwindows* opens */dev/rcm*.
2. *AIXwindows* gets the ***gsc\_handle*** from RCM via an *ioctl*.
3. *AIXwindows* becomes a graphics process *aixgsc* (*MAKE\_GP*, ...)
4. *AIXwindows*, through an *ioctl*, informs RCM about the displays it wishes to use.
5. *AIXwindows* opens all of the input devices it needs and passes the same input ring to each of them.

### Upon Normal Termination

1. X issues a *close* to all of the input devices it opened.
2. X informs RCM, through an *ioctl*, about the displays it was using.

## Diagnostics

*Diagnostics* and other applications that require access to the graphics adapter use the *AIXwindows* to *lft* interface.

---

## Accented Characters

Here are the valid sets of characters for each of the diacritics that the Low Function Terminal (LFT) subsystem uses to validate the two-key nonspacing character sequence.

### List of Diacritics Supported by the HFT LFT Subsystem

There are seven diacritic characters for which sets of characters are provided:

- Acute
- Grave
- Circumflex
- Umlaut
- Tilde
- Overcircle
- Cedilla

### Valid Sets of Characters (Categorized by Diacritics)

The following are acute function code values:

<b>Acute Function</b>	<b>Code Value</b>
Acute accent	0xef
Apostrophe (acute)	0x27
e Acute small	0x82
e Acute capital	0x90
a Acute small	0xa0
i Acute small	0xa1
o Acute small	0xa2
u Acute small	0xa3
a Acute capital	0xb5
i Acute capital	0xd6
y Acute small	0xec
y Acute capital	0xed
o Acute capital	0xe0
u Acute capital	0xe9

The following are grave function code values:

<b>Grave Function</b>	<b>Code Value</b>
Grave accent	0x60
a Grave small	0x85
e Grave small	0x8a
i Grave small	0x8d
o Grave small	0x95
u Grave small	0x97
a Grave capital	0xb7
e Grave capital	0xd4
i Grave capital	0xde
o Grave capital	0xe3
u Grave capital	0xeb

The following are circumflex function code values:

<b>Circumflex Function</b>	<b>Code Value</b>
----------------------------	-------------------



^ Circumflex accent	0x5e
a Circumflex small	0x83
e Circumflex small	0x88
i Circumflex small	0x8c
o Circumflex small	0x93
u Circumflex small	0x96
a Circumflex capital	0xb6
e Circumflex capital	0xd2
i Circumflex capital	0xd7
o Circumflex capital	0xe2
u Circumflex capital	0xea

The following are umlaut function code values:

<b>Umlaut Function</b>	<b>Code Value</b>
Umlaut accent	0xf9
u Umlaut small	0x81
a Umlaut small	0x84
e Umlaut small	0x89
i Umlaut small	0x8b
a Umlaut capital	0x8e
O Umlaut capital	0x99
u Umlaut capital	0x9a
e Umlaut capital	0xd3
i Umlaut capital	0xd8

The following are tilde function code values:

<b>Tilde Function</b>	<b>Code Value</b>
Tilde accent	0x7e
n Tilde small	0xa4
n Tilde capital	0xa5
a Tilde small	0xc6
a Tilde capital	0xc7
o Tilde small	0xe4
o Tilde capital	0xe5
Overcircle Function	Code Value
Overcircle accent	0x7d
a Overcircle small	0x86
a Overcircle capital	0x8f
Cedilla Function	Code Value
Cedilla accent	0xf7
c Cedilla capital	0x80
c Cedilla small	0x87

---

## Related Information

National Language Support Overview, Setting National Language Support for Devices, Locales in *Operating system and device management*

Keyboard Overview in *Keyboard Technical Reference*

Understanding the Japanese Input Method (JIM), Understanding the Korean Input Method (KIM), Understanding the Traditional Chinese Input Method (TIM) in *AIX Version 6.1 General Programming Concepts: Writing and Debugging Programs*.

## **Commands References**

The **iconv** command in *AIX Version 6.1 Commands Reference, Volume 3*.

---

## Chapter 10. Logical Volume Subsystem

A logical volume subsystem provides flexible access and control for complex physical storage systems.

The following topics describe how the logical volume device driver (LVDD) interacts with physical volumes:

- “Direct Access Storage Devices (DASDs)”
- “Physical Volumes”
- “Understanding the Logical Volume Device Driver” on page 212
- “Understanding Logical Volumes and Bad Blocks” on page 215

---

### Direct Access Storage Devices (DASDs)

Direct access storage devices (DASDs) are *fixed* or *removable* storage devices. Typically, these devices are hard disks. A fixed storage device is any storage device defined during system configuration to be an integral part of the system DASD. The operating system detects an error if a fixed storage device is not available at some time during normal operation.

A removable storage device is any storage device defined by the person who administers your system during system configuration to be an optional part of the system DASD. The removable storage device can be removed from the system at any time during normal operation. As long as the device is logically unmounted first, the operating system does not detect an error.

The following types of devices are not considered DASD and are not supported by the logical volume manager (LVM):

- Diskettes
- CD-ROM (compact disk read-only memory)
- DVD-ROM (DVD read-only memory)
- WORM (write-once read-many)

For a description of the block level, see “DASD Device Block Level Description” on page 313.

---

### Physical Volumes

A logical volume is a portion of a physical volume viewed by the system as a volume. Logical records are records defined in terms of the information they contain rather than physical attributes.

A physical volume is a DASD structured for requests at the physical level, that is, the level at which a processing unit can request device-independent operations on a physical block address basis. A physical volume is composed of the following:

- A device-dependent reserved area
- A variable number of physical blocks that serve as DASD descriptors
- An integral number of partitions, each containing a fixed number of physical blocks

When performing I/O at a physical level, no bad-block relocation is supported. Bad blocks are not hidden at this level as they are at the logical level. Typical operations at the physical level are `read-physical-block` and `write-physical-block`. For more information on bad blocks, see “Understanding Logical Volumes and Bad Blocks” on page 215.

The following are terms used when discussing DASD volumes:

**block**                    A contiguous, 512-byte region of a physical volume that corresponds in size to a DASD sector

**partition** A set of blocks (with sequential cylinder, head, and sector numbers) contained within a single physical volume

The number of blocks in a partition, as well as the number of partitions in a given physical volume, are fixed when the physical volume is installed in a volume group. Every physical volume in a volume group has exactly the same partition size. There is no restriction on the types of DASDs (for example, Small Computer Systems Interface (SCSI), Enhanced Small Device Interface (ESDI), or Intelligent Peripheral Interface (IPI)) that can be placed in a given volume group.

**Note:** A given physical volume must be assigned to a volume group before that physical volume can be used by the LVM.

## Physical Volume Implementation Limitations

When composing a physical volume from a DASD, the following implementation restrictions apply to DASD characteristics:

- 1 to 32 physical volumes per volume group
- 1 to 128 physical volumes in a big volume group
- The partition size is restricted to  $2^{**n}$  bytes, for  $20 \leq n \leq 30$
- The physical block size is restricted to 512 bytes

## Physical Volume Layout

A physical volume consists of a logically contiguous string of physical sectors. Sectors are numbered 0 through the last physical sector number (LPSN) on the physical volume. The total number of physical sectors on a physical volume is  $LPSN + 1$ . The actual physical location and physical order of the sectors are transparent to the sector numbering scheme.

**Note:** Sector numbering applies to user-accessible data sectors only. Spare sectors and Customer-Engineer (CE) sectors are not included. CE sectors are reserved for use by diagnostic test routines or microcode.

## Reserved Sectors on a Physical Volume

A physical volume reserves the first 128 sectors to store various types of DASD configuration and operation information. The `/usr/include/sys/hd_psn.h` file describes the information stored on the reserved sectors. The locations of the items in the reserved area are expressed as physical sector numbers in this file, and the lengths of those items are in number of sectors.

The 128-sector reserved area of a physical volume includes a boot record, the bad-block directory, the LVM record, and the mirror write consistency (MWC) record. The boot record consists of one sector containing information that allows the read-only system (ROS) to boot the system. A description of the boot record can be found in the `/usr/include/sys/bootrecord.h` file.

The boot record also contains the `pv_id` field. This field is a 64-bit number uniquely identifying a physical volume. This identifier might be assigned by the manufacturer of the physical volume. However, if a physical volume is part of a volume group, the `pv_id` field will be assigned by the LVM.

The bad-block directory records the blocks on the physical volume that have been diagnosed as unusable. The structure of the bad-block directory and its entries can be found in the `/usr/include/sys/bbdir.h` file.

The LVM record consists of one sector and contains information used by the LVM when the physical volume is a member of the volume group. The LVM record is described in the `/usr/include/lvmrec.h` file.

The MWC record consists of one sector. It identifies which logical partitions might be inconsistent if the system is not shut down properly. When the volume group is varied back online for use, this information is used to make logical partitions consistent again.

## Sectors Reserved for the Logical Volume Manager (LVM)

If a physical volume is part of a volume group, the physical volume is used by the LVM and contains two additional reserved areas. One area contains the volume group descriptor area/volume group status area and follows the first 128 reserved sectors. The other area is at the end of the physical volume reserved as a relocation pool for bad blocks that must be software-relocated. Both of these areas are described by the LVM record. The space between these last two reserved areas is divided into equal-sized partitions.

The volume group descriptor area (VGDA) is divided into the following:

- The volume group header. This header contains general information about the volume group and a time stamp used to verify the consistency of the VGDA.
- A list of logical volume entries. The logical volume entries describe the states and policies of logical volumes. This list defines the maximum number of logical volumes allowed in the volume group. The maximum is specified when a volume group is created.
- A list of physical volume entries. The size of the physical volume list is variable because the number of entries in the partition map can vary for each physical volume. For example, a 200 MB physical volume with a partition size of 1 MB has 200 partition map entries.
- A name list. This list contains the special file names of each logical volume in the volume group.
- A volume group trailer. This trailer contains an ending time stamp for the volume group descriptor area.

When a volume group is varied online, a majority (also called a quorum) of VGDA's must be present to perform recovery operations unless you have specified the **force** flag. (The vary-on operation, performed by using the **varyonvg** command, makes a volume group available to the system.) See Logical volume storage in *Operating system and device management* for introductory information about the vary-on process and quorums.

**Attention:** Use of the **force** flag can result in data inconsistency.

A volume group with only one physical volume must contain two copies of the physical volume group descriptor area. For any volume group containing more than one physical volume, there are at least three on-disk copies of the volume group descriptor area. The default placement of these areas on the physical volume is as follows:

- For the first physical volume installed in a volume group, two copies of the volume group descriptor area are placed on the physical volume.
- For the second physical volume installed in a volume group, one copy of the volume group descriptor area is placed on the physical volume.
- For the third physical volume installed in a volume group, one copy of the volume group descriptor area is placed on the physical volume. The second copy is removed from the first volume.
- For additional physical volumes installed in a volume group, one copy of the volume group descriptor area is placed on the physical volume.

When a vary-on operation is performed, a majority of copies of the volume group descriptor area must be able to come online before the vary-on operation is considered successful. A quorum ensures that at least one copy of the volume group descriptor areas available to perform recovery was also one of the volume group descriptor areas that were online during the previous vary-off operation. If not, the consistency of the volume group descriptor area cannot be ensured.

The volume group status area (VGSA) contains the status of all physical volumes in the volume group. This status is limited to active or missing. The VGSA also contains the state of all allocated physical

partitions (PP) on all physical volumes in the volume group. This state is limited to active or stale. A PP with a stale state is not used to satisfy a read request and is not updated on a write request.

A PP changes from stale to active after a successful resynchronization of the logical partition (LP) that has multiple copies, or mirrors, and is no longer consistent with its peers in the LP. This inconsistency can be caused by a write error or by not having a physical volume available when the LP is written to or updated.

A PP changes from stale to active after a successful resynchronization of the LP. A resynchronization operation issues resynchronization requests starting at the beginning of the LP and proceeding sequentially through its end. The LVDD reads from an active partition in the LP and then writes that data to any stale partition in the LP. When the entire LP has been traversed, the partition state is changed from stale to active.

Normal I/O can occur concurrently in an LP that is being resynchronized.

**Note:** If a write error occurs in a stale partition while a resynchronization is in progress, that partition remains stale.

If all stale partitions in an LP encounter write errors, the resynchronization operation is ended for this LP and must be restarted from the beginning.

The vary-on operation uses the information in the VGSA to initialize the LVDD data structures when the volume group is brought online.

---

## Understanding the Logical Volume Device Driver

The Logical Volume Device Driver (LVDD) is a pseudo-device driver that operates on logical volumes through the `/dev/lvn` special file. Like the physical disk device driver, this pseudo-device driver provides character and block entry points with compatible arguments. Each volume group has an entry in the kernel device switch table. Each entry contains entry points for the device driver and a pointer to the volume group data structure. The logical volumes of a volume group are distinguished by their minor numbers.

**Attention:** Each logical volume has a control block located in the first 512 bytes. Data begins in the second 512-byte block. Care must be taken when reading and writing directly to the logical volume, such as when using applications that write to raw logical volumes, because the control block is not protected from such writes. If the control block is overwritten, commands that use the control block will use default information instead.

Character I/O requests are performed by issuing a read or write request on a `/dev/rlvn` character special file for a logical volume. The read or write is processed by the file system SVC handler, which calls the LVDD `ddread` or `ddwrite` entry point. The `ddread` or `ddwrite` entry point transforms the character request into a block request. This is done by building a buffer for the request and calling the LVDD `ddstrategy` entry point.

Block I/O requests are performed by issuing a read or write on a block special file `/dev/lvn` for a logical volume. These requests go through the SVC handler to the `bread` or `bwrite` block I/O kernel services. These services build buffers for the request and call the LVDD `ddstrategy` entry point. The LVDD `ddstrategy` entry point then translates the logical address to a physical address (handling bad block relocation and mirroring) and calls the appropriate physical disk device driver.

On completion of the I/O, the physical disk device driver calls the `iodone` kernel service on the device interrupt level. This service then calls the LVDD I/O completion-handling routine. Once this is completed, the LVDD calls the `iodone` service again to notify the requester that the I/O is completed.

The LVDD is logically split into top and bottom halves. The top half contains the `ddopen`, `ddclose`, `ddread`, `ddwrite`, `ddioctl`, and `ddconfig` entry points. The bottom half contains the `ddstrategy` entry point,

which contains block read and write code. This is done to isolate the code that must run fully pinned and has no access to user process context. The bottom half of the device driver runs on interrupt levels and is not permitted to page fault. The top half runs in the context of a process address space and can page fault.

## Data Structures

The interface to the **ddstrategy** entry point is one or more logical **buf** structures in a list. The logical **buf** structure is defined in the `/usr/include/sys/buf.h` file and contains all needed information about an I/O request, including a pointer to the data buffer. The **ddstrategy** entry point associates one or more (if mirrored) physical **buf** structures (or **pbufs**) with each logical **buf** structure and passes them to the appropriate physical device driver.

The **pbuf** structure is a standard **buf** structure with some additional fields. The LVDD uses these fields to track the status of the physical requests that correspond to each logical I/O request. A pool of pinned **pbuf** structures is allocated and managed by the LVDD.

There is one device switch entry for each volume group defined on the system. Each volume group entry contains a pointer to the volume group data structure describing it.

## Top Half of LVDD

The top half of the LVDD contains the code that runs in the context of a process address space and can page fault. It contains the following entry points:

<b>ddopen</b>	Called by the file system when a logical volume is mounted, to open the logical volume specified.
<b>ddclose</b>	Called by the file system when a logical volume is unmounted, to close the logical volume specified.
<b>ddconfig</b>	Initializes data structures for the LVDD.
<b>ddread</b>	Called by the <b>read</b> subroutine to translate character I/O requests to block I/O requests. This entry point verifies that the request is on a 512-byte boundary and is a multiple of 512 bytes in length.

Most of the time a request will be sent down as a single request to the LVDD **ddstrategy** entry point which handles logical block I/O requests. However, the **ddread** routine might divide very large requests into multiple requests that are passed to the LVDD **ddstrategy** entry point.

If the `ext` parameter is set (called by the **readx** subroutine), the **ddread** entry point passes this parameter to the LVDD **ddstrategy** routine in the `b_options` field of the buffer header.

<b>ddwrite</b>	Called by the <b>write</b> subroutine to translate character I/O requests to block I/O requests. The LVDD <b>ddwrite</b> routine performs the same processing for a write request as the LVDD <b>ddread</b> routine does for read requests.
----------------	---

**ddioctl** Supports the following operations:

### CACLNUP

Causes the mirror write consistency (MWC) cache to be written to all physical volumes (PVs) in a volume group.

### IOCINFO, XLATE

Return LVM configuration information and PP status information.

### LV\_INFO

Provides information about a logical volume.

### PBUFCNT

Increases the number of physical buffer headers (pbufs) in the LVM pbuf pool.

## Bottom Half of the LVDD

The bottom half of the device driver supports the **ddstrategy** entry point. This entry point processes all logical block requests and performs the following functions:

- Validates I/O requests.

- Checks requests for conflicts (such as overlapping block ranges) with requests currently in progress.
- Translates logical addresses to physical addresses.
- Handles mirroring and bad-block relocation.

The bottom half of the LVDD runs on interrupt levels and, as a result, is not permitted to page fault. The bottom half of the LVDD is divided into the following three layers:

- Strategy layer
- Scheduler layer
- Physical layer

Each logical I/O request passes down through the bottom three layers before reaching the physical disk device driver. Once the I/O is complete, the request returns back up through the layers to handle the I/O completion processing at each layer. Finally, control returns to the original requestor.

### Strategy Layer

The strategy layer deals only with logical requests. The **ddstrategy** entry point is called with one or more logical **buf** structures. A list of **buf** structures for requests that are not blocked are passed to the second layer, the scheduler.

### Scheduler Layer

The scheduler layer schedules physical requests for logical operations and handles mirroring and the MWC cache. For each logical request the scheduler layer schedules one or more physical requests. These requests involve translating logical addresses to physical addresses, handling mirroring, and calling the LVDD physical layer with a list of physical requests.

When a physical I/O operation is complete for one phase or mirror of a logical request, the scheduler initiates the next phase (if there is one). If no more I/O operations are required for the request, the scheduler calls the strategy termination routine. This routine notifies the originator that the request has been completed.

The scheduler also handles the MWC cache for the volume group. If a logical volume is using mirror write consistency, then requests for this logical volume are held within the scheduling layer until the MWC cache blocks can be updated on the target physical volumes. When the MWC cache blocks have been updated, the request proceeds with the physical data write operations.

When MWC is being used, system performance can be adversely affected. This is caused by the overhead of logging or journalling that a write request is active in one or more logical track groups (LTGs) (128K, 256K, 512K or 1024K). This overhead is for mirrored writes only. It is necessary to guarantee data consistency between mirrors particularly if the system crashes before the write to all mirrors has been completed.

Mirror write consistency can be turned off for an entire logical volume. It can also be inhibited on a request basis by turning on the **NO\_MWC** flag as defined in the `/usr/include/sys/lvdd.h` file.

### Physical Layer

The physical layer of the LVDD handles startup and termination of the physical request. The physical layer calls a physical disk device driver's **ddstrategy** entry point with a list of **buf** structures linked together. In turn, the physical layer is called by the **iodone** kernel service when each physical request is completed.

This layer also performs bad-block relocation and detection/correction of bad blocks, when necessary. These details are hidden from the other two layers.



## Interface to Physical Disk Device Drivers

Physical disk device drivers adhere to the following criteria if they are to be accessed by the LVDD:

- Disk block size must be 512 bytes.
- The physical disk device driver needs to accept a list of requests defined by **buf** structures, which are linked together by the `av_forw` field in each **buf** structure.
- For unrecoverable media errors, physical disk device drivers need to set the following:
  - The **B\_ERROR** flag must be set to on (defined in the `/usr/include/sys/buf.h` file) in the `b_flags` field.
  - The `b_error` field must be set to **E\_MEDIA** (defined in the `/usr/include/sys/errno.h` file).
  - The `b_resid` field must be set to the number of bytes in the request that were not read or written successfully. The `b_resid` field is used to determine the block in error.

**Note:** For write requests, the LVDD attempts to hardware-relocate the bad block. If this is unsuccessful, then the block is software-relocated. For read requests, the information is recorded and the block is relocated on the next write request to that block.

- For a successful request that generated an excessive number of retries, the device driver can return good data. To indicate this situation it must set the following:
  - The `b_error` field is set to **ESOFT**; this is defined in the `/usr/include/sys/errno.h` file.
  - The `b_flags` field has the **B\_ERROR** flag set to on.
  - The `b_resid` field is set to a count indicating the first block in the request that had excessive retries. This block is then relocated.
- The physical disk device driver needs to accept a request of one block with **HWRELOC** (defined in the `/usr/include/sys/lvdd.h` file) set to on in the `b_options` field. This indicates that the device driver is to perform a hardware relocation on this request. If the device driver does not support hardware relocation the following should be set:
  - The `b_error` field is set to **EIO**; this is defined in the `/usr/include/sys/errno.h` file.
  - The `b_flags` field has the **B\_ERROR** flag set on.
  - The `b_resid` field is set to a count indicating the first block in the request that has excessive retries.
- The physical disk device driver should support the system dump interface as defined.
- The physical disk device driver must support write verification on an I/O request. Requests for write verification are made by setting the `b_options` field to **WRITEV**. This value is defined in the `/usr/include/sys/lvdd.h` file.

---

## Understanding Logical Volumes and Bad Blocks

The physical layer of the logical volume device driver (LVDD) initiates all bad-block processing and isolates all of the decision making from the physical disk device driver. This happens so the physical disk device driver does not need to handle mirroring, which is the duplication of data transparent to the user.

### Relocating Bad Blocks

The physical layer of the LVDD checks each physical request to see if there are any known software-relocated bad blocks in the request. The LVDD determines if a request contains known software-relocated bad blocks by hashing the physical address. Then a hash chain of the LVDD defects directory is searched to see if any bad-block entries are in the address range of the request.

If bad blocks exist in a physical request, the request is split into pieces. The first piece contains any blocks up to the relocated block. The second piece contains the relocated block (the relocated address is specified in the bad-block entry) of the defects directory. The third piece contains any blocks after the relocated block to the end of the request or to the next relocated block. These separate pieces are processed sequentially until the entire request has been satisfied.

Once the I/O for the first of the separate pieces has completed, the **iodone** kernel service calls the LVDD physical layer's termination routine (specified in the **b\_done** field of the **buf** structure). The termination routine initiates I/O for the second piece of the original request (containing the relocated block), and then for the third piece. When the entire physical operation is completed, the appropriate scheduler's policy routine (in the second layer of the LVDD) is called to start the next phase of the logical operation.

## Detecting and Correcting Bad Blocks

If a logical volume is mirrored, a newly detected bad block is fixed by relocating that block. A good mirror is read and then the block is relocated using data from the good mirror. With mirroring, the user does not need to know when bad blocks are found. However, the physical disk device driver does log permanent I/O errors so the user can determine the rate of media surface errors.

When a bad block is detected during I/O, the physical disk device driver sets the error fields in the **buf** structure to indicate that there was a media surface error. The physical layer of the LVDD then initiates any bad-block processing that must be done.

If the operation was a nonmirrored read, the block is not relocated because the data in the relocated block is not initialized until a write is performed to the block. To support this delayed relocation, an entry for the bad block is put into the LVDD defects directory and into the bad-block directory on disk. These entries contain no relocated block address and the status for the block is set to indicate that relocation is desired.

On each I/O request, the physical layer checks whether there are any bad blocks in the request. If the request is a write and contains a block that is in a **relocation-desired** state, the request is sent to the physical disk device driver with safe hardware relocation requested. If the request is a read, a read of the known defective block is attempted.

If the operation was a read operation in a mirrored LP, a request to read one of the other mirrors is initiated. If the second read is successful, then the read is turned into a write request and the physical disk device driver is called with safe hardware relocation specified to fix the bad mirror.

If the hardware relocation fails or the device does not support safe hardware relocation, the physical layer of the LVDD attempts software relocation. At the end of each volume is a reserved area used by the LVDD as a pool of relocation blocks. When a bad block is detected and the disk device driver is unable to relocate the block, the LVDD picks the next unused block in the relocation pool and writes to this new location. A new entry is added to the LVDD defects directory in memory (and to the bad-block directory on disk) that maps the bad-block address to the new relocation block address. Any subsequent I/O requests to the bad-block address are routed to the relocation address.

**Attention:** Formatting a fixed disk deletes any data on the disk. Format a fixed disk only when absolutely necessary and preferably after backing up all data on the disk.

If you need to format a fixed disk completely (including reinitializing any bad blocks), use the formatting function supplied by the **diag** command. (The **diag** command typically, but not necessarily, writes over all data on a fixed disk. Refer to the documentation that comes with the fixed disk to determine the effect of formatting with the **diag** command.)

---

## Related Information

### Subroutine References

The **readx** subroutine, **write** subroutine in *AIX Version 6.1 Technical Reference: Base Operating System and Extensions Volume 2*.

### Files Reference

The **lvdd** special file in *AIX Version 6.1 Files Reference*.

## Technical References

The **buf** structure in *AIX Version 6.1 Technical Reference: Kernel and Subsystems Volume 1*.

The **bread** kernel service, **bwrite** kernel service, **iodone** kernel service in *AIX Version 6.1 Technical Reference: Kernel and Subsystems Volume 1*.



---

## Chapter 11. Printer Addition Management Subsystem

If you are configuring a printer for your system, there are basically two types of printers: printers already supported by the operating system and new printer types. Printer Support in *Printers and printing* lists supported printers.

---

### Printer Types Currently Supported

To configure a supported type of printer, you need only to run the **mkvirprt** command to create a customized printer file for your printer. This customized printer file, which is in the **/var/spool/lpd/piop/@local/custom** directory, describes the specific parameters for your printer. For more information see Configuring a Printer without Adding a Queue in *Printers and printing*.

---

### Printer Types Currently Unsupported

To configure a currently unsupported type of printer, you must develop and add a predefined printer definition for your printer. This new option is then entered in the list of available choices when the user selects a printer to configure for the system. The actual data used by the printer subsystem comes from the Customized printer definition created by the **mkvirprt** command.

“Adding a New Printer Type to Your System” provides general instructions for adding an undefined printer. To add an undefined printer, you modify an existing printer definition. Undefined printers fall into two categories:

- Printers that closely emulate a supported printer. You can use SMIT or the virtual printer commands to make the changes you need.
- Printers that do not emulate a supported printer or that emulate several data streams. It is simpler to make the necessary changes for these printers by editing the printer colon file. See Adding a Printer Using the Printer Colon File in *Printers and printing*.

“Adding an Unsupported Device to the System” on page 100 offers an overview of the major steps required to add an unsupported device of any type to your system.

---

### Adding a New Printer Type to Your System

To add an unsupported printer to your system, you must add a new Printer definition to the printer directories. For more complicated scenarios, you might also need to add a new printer-specific formatter to the printer backend.

Example of Print Formatter in *Printers and printing* shows how the print formatter interacts with the printer formatter subroutines.

### Additional Steps for Adding a New Printer Type

However, if you want the new Printer definition to carry the name of the new printer, you must develop a new Predefined definition to carry the new printer information besides adding a new Printer definition. Use the **piopredef** command to do this.

Steps for adding a new printer-specific formatter to the printer backend are discussed in Adding a Printer Formatter to the Printer Backend . Example of Print Formatter in *Printers and printing* shows how print formatters can interact with the printer formatter subroutines.

**Note:** These instructions apply to the addition of a new printer definition to the system, not to the addition of a physical printer device itself. For information on adding a new printer device, refer to device

configuration and management. If your new printer requires an interface other than the parallel or serial interface provided by the operating system, you must also provide a new device driver.

If the printer being added does not emulate a supported printer or if it emulates several data streams, you need to make more changes to the Printer definition. It is simpler to make the necessary changes for these printers by editing the printer colon file. See *Adding a Printer Using the Printer Colon File* in *Printers and printing*.

## Modifying Printer Attributes

Edit the customized file (`/var/spool/lpd/pio/custom /var/spool/lpd/pio/@local/custom`

`QueueName:QueueDeviceName`), adding or changing the printer attributes to match the new printer.

For example, assume that you created a new file based on the existing 4201-3 printer. The customized file for the 4201-3 printer contains the following template that the printer formatter uses to initialize the printer:

```
%I[ez,em,eA,cv,eC,e0,cp,cc, . . .
```

The formatter fills in the string as directed by this template and sends the resulting sequence of commands to the 4201-3 printer. Specifically, this generates a string of escape sequences that initialize the printer and set such parameters as vertical and horizontal spacing and page length. You would construct a similar command string to properly initialize the new printer and put it into 4201-emulation mode. Although many of the escape sequences might be the same, at least one will be different: the escape sequence that is the command to put the printer into the specific printer-emulation mode. Assume that you added an **ep** attribute that specifies the string to initialize the printer to 4201-3 emulation mode, as follows:

```
\033\012\013
```

The Printer Initialization field will then be:

```
%I[ep,ez,em,eA,cv,eC,e0,cp,cc, . . .
```

You must create a virtual printer for each printer-emulation mode you want to use. See *Real and Virtual Printers* in *Printers and printing*.

---

## Adding a Printer Definition

To add a new printer to the system, you must first create a description of the printer by adding a new printer definition to the printer definition directories.

Typically, to add a new printer definition to the database, you first modify an existing printer definition and then create a customized printer definition in the Customized Printer Directory.

Once you have added the new customized printer definition to the directory, the **mkvirprt** command uses it to present the new printer as a choice for printer addition and selection. Because the new printer definition is a customized printer definition, it appears in the list of printers under the name of the original printer from which it was customized.

A totally new printer must be added as a predefined printer definition in the `/usr/lib/lpd/pio/predef` directory. If the user chooses to work with printers once this new predefined printer definition is added to the Predefined Printer Directory, the **mkvirprt** command can then list all the printers in that directory. The added printer appears on the list of printers given to the user as if it had been supported all along. Specific information about this printer can then be extended, added, modified, or deleted, as necessary.

Printer Support in *Printers and printing* lists the supported printer types and names of representative printers.

---

## Adding a Printer Formatter to the Printer Backend

If your new printer's data stream differs significantly from one of the numerous printer data streams currently handled by the operating system, you must define a new backend formatter. Adding a new formatter does not require the addition of a new backend. Instead, all you typically need are modifications to the formatter commands associated with that printer under the supervision of the existing printer backend. If a new backend is required, see Printer Backend Overview for Programming in *Printers and printing*.

---

## Understanding Embedded References in Printer Attribute Strings

The attribute string retrieved by the **piocmdout**, **piogetstr**, and **piogetvals** subroutines can contain embedded references to other attribute strings or integers. The attribute string can also contain embedded logic that dynamically determines the content of the constructed string. This allows the constructed string to reflect the state of the formatter environment when one of these subroutines is called.

Embedded references and logic are defined with escape sequences that are placed at appropriate locations in the attribute string. The first character of each escape sequence is always the % character. This character indicates the beginning of an escape sequence. The second character (and sometimes subsequent characters) define the operation to be performed. The remainder of the characters (if any) in the escape sequence are operands to be used in performing the specified operation.

The escape sequences that can be specified in an attribute string are based on the **terminfo** parameterized string escape sequences for terminals. These escape sequences have been modified and extended for printers.

The attribute names that can be referenced by attribute strings are:

- The names of all attribute variables (which can be integer or string variables) defined to the **piogetvals** subroutine. When references are made to these variables, the **piogetvals**-defined versions are the values used.
- All other attributes names in the database. These attributes are considered string constants.

Any attribute value (integer variable, string variable, or string constant) can be referenced by any attribute string. Consequently, it is important that the formatter ensures that the values for all the integer variables and string variables defined to the **piogetvals** subroutine are kept current.

The formatter must not assume that the particular attribute string whose name it specifies to the **piogetstr** or **piocmdout** subroutine does not reference certain variables. The attribute string is retrieved from the database that is external to the formatter. The values in the database represented by the string can be changed to reference additional variables without the formatter's knowledge.

---

## Related Information

*Printers and printing*

### Subroutine References

The **piocmdout** subroutine, **piogetstr** subroutine, **piogetvals** subroutine in *AIX Version 6.1 Technical Reference: Base Operating System and Extensions Volume 1*.

### Commands References

The **mkvirprt** command in *AIX Version 6.1 Commands Reference, Volume 3*.

The **piopredef** command in *AIX Version 6.1 Commands Reference, Volume 4*.





---

## Chapter 12. Small Computer System Interface Subsystem (Parallel SCSI)

This overview describes the interface between a small computer system interface (SCSI) device driver and a SCSI adapter device driver. The information in the article is specific for the parallel SCSI implementation. Parallel SCSI is the traditional physical transport that was originally described in the SCSI-1 and SCSI-2 standards, culminating in support of Ultra 320. In the original implementation, SCSI was both a physical transport (parallel bus) and a logical command and response protocol. In the following topics, the term SCSI is used only for the traditional, parallel bus implementation of SCSI. The following topics can help you design and write a parallel SCSI device driver that can interface with an existing parallel SCSI adapter device driver. You can also use the following topics to design and write a parallel SCSI adapter device driver that interfaces with existing parallel SCSI device drivers.

For information about the implementation of SCSI on Fibre Channel (FCP), iSCSI and SAS transport types, see SCSI Architectural Model Subsystem.

---

### SCSI Subsystem Overview

The main topics covered in this overview are:

- Responsibilities of the SCSI Adapter Device Driver
- Responsibilities of the SCSI Device Driver
- Initiator-Mode Support
- Target-Mode Support

This section frequently refers to both a *SCSI device driver* and a *SCSI adapter device driver*. These two distinct device drivers work together in a layered approach to support attachment of a range of SCSI devices. The SCSI adapter device driver is the *lower* device driver of the pair, and the SCSI device driver is the *upper* device driver.

### Responsibilities of the SCSI Adapter Device Driver

The SCSI adapter device driver (the lower layer) is the software interface to the system hardware. This hardware includes the SCSI bus hardware plus any other system I/O hardware required to run an I/O request. The SCSI adapter device driver hides the details of the I/O hardware from the SCSI device driver. The design of the software interface allows a user with limited knowledge of the system hardware to write the upper device driver.

The SCSI adapter device driver manages the SCSI bus but not the SCSI devices. It can send and receive SCSI commands, but it cannot interpret the contents of the commands. The lower driver also provides recovery and logging for errors related to the SCSI bus and system I/O hardware. Management of the device specifics is left to the SCSI device driver. The interface of the two drivers allows the upper driver to communicate with different SCSI bus adapters without requiring special code paths for each adapter.

### Responsibilities of the SCSI Device Driver

The SCSI device driver (the upper layer) provides the rest of the operating system with the software interface to a given SCSI device or device class. The upper layer recognizes which SCSI commands are required to control a particular SCSI device or device class. The SCSI device driver builds I/O requests containing device SCSI commands and sends them to the SCSI adapter device driver in the sequence needed to operate the device successfully. The SCSI device driver cannot manage adapter resources or give the SCSI command to the adapter. Specifics about the adapter and system hardware are left to the lower layer.

The SCSI device driver also provides recovery and logging for errors related to the SCSI device it controls.

The operating system provides several kernel services allowing the SCSI device driver to communicate with SCSI adapter device driver entry points without having the actual name or address of those entry points. The description contained in “Logical File System Kernel Services” on page 63 can provide more information.

## Communication between SCSI Devices

When two SCSI devices communicate, one assumes the initiator-mode role, and the other assumes the target-mode role. The initiator-mode device generates the SCSI command, which requests an operation, and the target-mode device receives the SCSI command and acts. It is possible for a SCSI device to perform both roles simultaneously.

When writing a new SCSI adapter device driver, the writer must know which mode or modes must be supported to meet the requirements of the SCSI adapter and any interfaced SCSI device drivers. When a SCSI adapter device driver is added so that a new SCSI adapter works with all existing SCSI device drivers, both initiator-mode and target-mode must be supported in the SCSI adapter device driver.

### Initiator-Mode Support

The interface between the SCSI device driver and the SCSI adapter device driver for initiator-mode support (that is, the attached device acts as a target) is accessed through calls to the SCSI adapter device driver **open**, **close**, **ioctl**, and **strategy** routines. I/O requests are queued to the SCSI adapter device driver through calls to its strategy entry point.

Communication between the SCSI device driver and the SCSI adapter device driver for a particular initiator I/O request is made through the **sc\_buf** structure, which is passed to and from the strategy routine in the same way a standard driver uses a **struct buf** structure.

### Target-Mode Support

The interface between the SCSI device driver and the SCSI adapter device driver for target-mode support (that is, the attached device acts as an initiator) is accessed through calls to the SCSI adapter device driver **open**, **close**, and **ioctl** subroutines. Buffers that contain data received from an attached initiator device are passed from the SCSI adapter device driver to the SCSI device driver, and back again, in **tm\_buf** structures.

Communication between the SCSI adapter device driver and the SCSI device driver for a particular data transfer is made by passing the **tm\_buf** structures by pointer directly to routines whose entry points have been previously registered. This registration occurs as part of the sequence of commands the SCSI device driver executes using calls to the SCSI adapter device driver when the device driver opens a target-mode device instance.

---

## Understanding SCSI Asynchronous Event Handling

**Note:** This operation is not supported by all SCSI I/O controllers.

A SCSI device driver can register a particular device instance for receiving asynchronous event status by calling the **SCIOEVENT** **ioctl** operation for the SCSI-adapter device driver. When an event covered by the **SCIOEVENT** **ioctl** operation is detected by the SCSI adapter device driver, it builds an **sc\_event\_info** structure and passes a pointer to the structure and to the asynchronous event-handler routine entry point, which was previously registered. The fields in the structure are filled in by the SCSI adapter device driver as follows:

<b>id</b>	For initiator mode, this is set to the SCSI ID of the attached SCSI target device. For target mode, this is set to the SCSI ID of the attached SCSI initiator device.
<b>lun</b>	For initiator mode, this is set to the SCSI LUN of the attached SCSI target device. For target mode, this is set to 0).

<b>mode</b>	Identifies whether the initiator or target mode device is being reported. The following values are possible:  <b>SC_IM_MODE</b> An initiator mode device is being reported.  <b>SC_TM_MODE</b> A target mode device is being reported.
<b>events</b>	This field is set to indicate what event or events are being reported. The following values are possible, as defined in the <code>/usr/include/sys/scsi.h</code> file:  <b>SC_FATAL_HDW_ERR</b> A fatal adapter hardware error occurred.  <b>SC_ADAP_CMD_FAILED</b> An unrecoverable adapter command failure occurred.  <b>SC_SCSI_RESET_EVENT</b> A SCSI bus reset was detected.  <b>SC_BUFS_EXHAUSTED</b> In target-mode, a maximum buffer usage event has occurred.
<b>adap_devno</b>	This field is set to indicate the device major and minor numbers of the adapter on which the device is located.
<b>async_correlator</b>	This field is set to the value passed to the SCSI adapter device driver in the <b>sc_event_struct</b> structure. The SCSI device driver may optionally use this field to provide an efficient means of associating event status with the device instance it goes with. Alternatively, the SCSI device driver uses the combination of the <code>id</code> , <code>lun</code> , <code>mode</code> , and <code>adap_devno</code> fields to identify the device instance.

**Note:** Reserved fields should be set to 0 by the SCSI adapter device driver.

The information reported in the `sc_event_info.events` field does not queue to the SCSI device driver, but is instead reported as one or more flags as they occur. Because the data does not queue, the SCSI adapter device driver writer can use a single **sc\_event\_info** structure and pass it one at a time, by pointer, to each asynchronous event handler routine for the appropriate device instance. After determining for which device the events are being reported, the SCSI device driver must copy the `sc_event_info.events` field into local space and must not modify the contents of the rest of the **sc\_event\_info** structure.

Because the event status is optional, the SCSI device driver writer determines what action is necessary to take upon receiving event status. The writer may decide to save the status and report it back to the calling application, or the SCSI device driver or application level program can take error recovery actions.

## Defined Events and Recovery Actions

The adapter fatal hardware failure event is intended to indicate that no further commands to or from this SCSI device are likely to succeed, because the adapter it is attached to has failed. It is recommended that the application end the session with the device.

The unrecoverable adapter command failure event is not necessarily a fatal condition, but it can indicate that the adapter is not functioning properly. Possible actions by the application program include:

- Ending of the session with the device in the near future
- Ending of the session after multiple (two or more) such events
- Attempting to continue the session indefinitely

The SCSI Bus Reset detection event is mainly intended as information only, but may be used by the application to perform further actions, if necessary.

The maximum buffer usage detected event applies only to a given target-mode device; it will not be reported for an initiator-mode device. This event indicates to the application that this particular target-mode device instance has filled its maximum allotted buffer space. The application should perform **read** system calls fast enough to prevent this condition. If this event occurs, data is not lost, but it is delayed to prevent further buffer usage. Data reception will be restored when the application empties enough buffers to continue reasonable operations. The **num\_bufs** attribute may need to be increased to help minimize this problem. Also, it is possible that regardless of the number of buffers, the application simply is not processing received data fast enough. This may require some fine tuning of the application's data processing routines.

## Asynchronous Event-Handling Routine

The SCSI-device driver asynchronous event-handling routine is typically called directly from the hardware interrupt-handling routine for the SCSI adapter device driver. The SCSI device driver writer must be aware of how this affects the design of the SCSI device driver.

Because the event handling routine is running on the hardware interrupt level, the SCSI device driver must be careful to limit operations in that routine. Processing should be kept to a minimum. In particular, if any error recovery actions are performed, it is recommended that the event-handling routine set state or status flags only and allow a process level routine to perform the actual operations.

The SCSI device driver must be careful to disable interrupts at the correct level in places where the SCSI device driver's lower execution priority routines manipulate variables that are also modified by the event-handling routine. To allow the SCSI device driver to disable at the correct level, the SCSI adapter device driver writer must provide a configuration database attribute that defines the interrupt class, or priority, it runs on. This attribute must be named **intr\_priority** so that the SCSI device driver configuration method knows which attribute of the parent adapter to query. The SCSI device driver configuration method should then pass this interrupt priority value to the SCSI device driver along with other configuration data for the device instance.

The SCSI device driver writer must follow any other general system rules for writing a routine that must execute in an interrupt environment. For example, the routine must not attempt to sleep or wait on I/O operations. It can perform wakeups to allow the process level to handle those operations.

Because the SCSI device driver copies the information from the **sc\_event\_info.events** field on each call to its asynchronous event-handling routine, there is no resource to free or any information which must be passed back later to the SCSI adapter device driver.

---

## SCSI Error Recovery

The SCSI error-recovery process handles different issues depending on whether the SCSI device is in initiator mode or target mode. If the device is in initiator mode, the error-recovery process varies depending on whether or not the device is supporting command queuing.

### SCSI Initiator-Mode Recovery When Not Command Tag Queuing

If an error such as a check condition or hardware failure occurs, transactions queued within the SCSI adapter device driver are terminated abnormally with **iodone** calls. The transaction active during the error is returned with the **sc\_buf.bufstruct.b\_error** field set to **EIO**. Other transactions in the queue are returned with the **sc\_buf.bufstruct.b\_error** field set to **ENXIO**. The SCSI device driver should process or recover the condition, rerunning any mode selects or device reservations to recover from this condition properly. After this recovery, it should reschedule the transaction that had the error. In many cases, the SCSI device driver only needs to retry the unsuccessful operation.

The SCSI adapter device driver should never retry a SCSI command on error after the command has successfully been given to the adapter. The consequences for retrying a SCSI command at this point range from minimal to catastrophic, depending upon the type of device. Commands for certain devices

cannot be retried immediately after a failure (for example, tapes and other sequential access devices). If such an error occurs, the failed command returns an appropriate error status with an **iodone** call to the SCSI device driver for error recovery. Only the SCSI device driver that originally issued the command knows if the command can be retried on the device. The SCSI adapter device driver must only retry commands that were never successfully transferred to the adapter. In this case, if retries are successful, the **sc\_buf** status should not reflect an error. However, the SCSI adapter device driver should perform error logging on the retried condition.

The first transaction passed to the SCSI adapter device driver during error recovery must include a special flag. This **SC\_RESUME** flag in the `sc_buf.flags` field must be set to inform the SCSI adapter device driver that the SCSI device driver has recognized the fatal error and is beginning recovery operations. Any transactions passed to the SCSI adapter device driver, after the fatal error occurs and before the **SC\_RESUME** transaction is issued, should be flushed; that is, returned with an error type of **ENXIO** through an **iodone** call.

**Note:** If a SCSI device driver continues to pass transactions to the SCSI adapter device driver after the SCSI adapter device driver has flushed the queue, these transactions are also flushed with an error return of **ENXIO** through the **iodone** service. This gives the SCSI device driver a positive indication of all transactions flushed.

If the SCSI device driver is executing a gathered write operation, the error-recovery information mentioned previously is still valid, but the caller must restore the contents of the `sc_buf.resvdw1` field and the **uio** struct that the field pointed to before attempting the retry. The retry must occur from the SCSI device driver's process level; it cannot be performed from the caller's **iodone** subroutine. Also, additional return codes of **EFAULT** and **ENOMEM** are possible in the `sc_buf.bufstruct.b_error` field for a gathered write operation.

## SCSI Initiator-Mode Recovery During Command Tag Queuing

If the SCSI device driver is queuing multiple transactions to the device and either a check condition error or a command terminated error occurs, the SCSI adapter driver does not clear all transactions in its queues for the device. It returns the failed transaction to the SCSI device driver with an indication that the queue for this device is not cleared by setting the **SC\_DID\_NOT\_CLEAR\_Q** flag in the `sc_buf.adap_q_status` field. The SCSI adapter driver halts the queue for this device awaiting error recovery notification from the SCSI device driver. The SCSI device driver then has three options to recover from this error:

- Send one error recovery command (request sense) to the device.
- Clear the SCSI adapter driver's queue for this device.
- Resume the SCSI adapter driver's queue for this device.

When the SCSI adapter driver's queue is halted, the SCSI device driver can get sense data from a device by setting the **SC\_RESUME** flag in the `sc_buf.flags` field and the **SC\_NO\_Q** flag in `sc_buf.q_tag_msg` field of the request-sense **sc\_buf**. This action notifies the SCSI adapter driver that this is an error-recovery transaction and should be sent to the device while the remainder of the queue for the device remains halted. When the request sense completes, the SCSI device driver needs to either clear or resume the SCSI adapter driver's queue for this device.

The SCSI device driver can notify the SCSI adapter driver to clear its halted queue by sending a transaction with the **SC\_Q\_CLR** flag in the `sc_buf.flags` field. This transaction must not contain a SCSI command because it is cleared from the SCSI adapter driver's queue without being sent to the adapter. However, this transaction must have the SCSI ID field (`sc_buf.scsi_command.scsi_id`) and the LUN fields (`sc_buf.scsi_command.scsi_cmd.lun` and `sc_buf.lun`) filled in with the device's SCSI ID and logical unit number (LUN). If addressing LUNs 8 - 31, the `sc_buf.lun` field should be set to the logical unit number and the `sc_buf.scsi_command.scsi_cmd.lun` field should be zeroed out. See the descriptions of these fields for further explanation. Upon receiving an **SC\_Q\_CLR** transaction, the SCSI adapter driver flushes all transactions for this device and sets their `sc_buf.bufstruct.b_error` fields to **ENXIO**. The SCSI device

driver must wait until the `sc_buf` with the `SC_Q_CLR` flag set is returned before it resumes issuing transactions. The first transaction sent by the SCSI device driver after it receives the returned `SC_Q_CLR` transaction must have the `SC_RESUME` flag set in the `sc_buf.flags` fields.

If the SCSI device driver wants the SCSI adapter driver to resume its halted queue, it must send a transaction with the `SC_Q_RESUME` flag set in the `sc_buf.flags` field. This transaction can contain an actual SCSI command, but it is not required. However, this transaction must have the `sc_buf.scsi_command.scsi_id`, `sc_buf.scsi_command.scsi_cmd.lun`, and the `sc_buf.lun` fields filled in with the device's SCSI ID and logical unit number. See the description of these fields for further details. If this is the first transaction issued by the SCSI device driver after receiving the error (indicating that the adapter driver's queue is halted), then the `SC_RESUME` flag must be set as well as the `SC_Q_RESUME` flag.

## Analyzing Returned Status

The following order of precedence should be followed by SCSI device drivers when analyzing the returned status:

1. If the `sc_buf.bufstruct.b_flags` field has the `B_ERROR` flag set, then an error has occurred and the `sc_buf.bufstruct.b_error` field contains a valid `errno` value.

If the `b_error` field contains the `ENXIO` value, either the command needs to be restarted or it was canceled at the request of the SCSI device driver.

If the `b_error` field contains the `EIO` value, then either one or no flag is set in the `sc_buf.status_validity` field. If a flag is set, an error in either the `scsi_status` or `general_card_status` field is the cause.

If the `status_validity` field is 0, then the `sc_buf.bufstruct.b_resid` field should be examined to see if the SCSI command issued was in error. The `b_resid` field can have a value without an error having occurred. To decide whether an error has occurred, the SCSI device driver must evaluate this field with regard to the SCSI command being sent and the SCSI device being driven.

If the SCSI device driver is queuing multiple transactions to the device and if either `SC_CHECK_CONDITION` or `SC_COMMAND_TERMINATED` is set in `scsi_status`, then the value of `sc_buf.adap_q_status` must be analyzed to determine if the adapter driver has cleared its queue for this device. If the SCSI adapter driver has not cleared its queue after an error, then it holds that queue in a halted state.

If `sc_buf.adap_q_status` is set to 0, the SCSI adapter driver has cleared its queue for this device and any transactions outstanding are flushed back to the SCSI device driver with an error of `ENXIO`.

If the `SC_DID_NOT_CLEAR_Q` flag is set in the `sc_buf.adap_q_status` field, the adapter driver has not cleared its queue for this device. When this condition occurs, the SCSI adapter driver allows the SCSI device driver to send one error recovery transaction (request sense) that has the field `sc_buf.q_tag_msg` set to `SC_NO_Q` and the field `sc_buf.flags` set to `SC_RESUME`. The SCSI device driver can then notify the SCSI adapter driver to clear or resume its queue for the device by sending a `SC_Q CLR` or `SC_Q RESUME` transaction.

If the SCSI device driver does not queue multiple transactions to the device (that is, the `SC_NO_Q` is set in `sc_buf.q_tag_msg`), then the SCSI adapter clears its queue on error and sets `sc_buf.adap_q_status` to 0.

2. If the `sc_buf.bufstruct.b_flags` field does not have the `B_ERROR` flag set, then no error is being reported. However, the SCSI device driver should examine the `b_resid` field to check for cases where less data was transferred than expected. For some SCSI commands, this occurrence might not represent an error. The SCSI device driver must determine if an error has occurred.

If a nonzero `b_resid` field does represent an error condition, then the device queue is not halted by the SCSI adapter device driver. It is possible for one or more succeeding queued commands to be sent to the adapter (and possibly the device). Recovering from this situation is the responsibility of the SCSI device driver.

3. In any of the above cases, if `sc_buf.bufstruct.b_flags` field has the `B_ERROR` flag set, then the queue of the device in question has been halted. The first `sc_buf` structure sent to recover the error (or continue operations) must have the `SC_RESUME` bit set in the `sc_buf.flags` field.

## Target-Mode Error Recovery

If an error occurs during the reception of **send** command data, the SCSI adapter device driver sets the **TM\_ERROR** flag in the `tm_buf.user_flag` field. The SCSI adapter device driver also sets the **SC\_ADAPTER\_ERROR** bit in the `tm_buf.status_validity` field and sets a single flag in the `tm_buf.general_card_status` field to indicate the error that occurred.

In the SCSI subsystem, an error during a **send** command does not affect future target-mode data reception. Future **send** commands continue to be processed by the SCSI adapter device driver and queue up, as necessary, after the data with the error. The SCSI device driver continues processing the **send** command data, satisfying user read requests as usual except that the error status is returned for the appropriate user request. Any error recovery or synchronization procedures the user requires for a target-mode received-data error must be implemented in user-supplied software.

---

## A Typical Initiator-Mode SCSI Driver Transaction Sequence

A simplified sequence of events for a transaction between a SCSI device driver and a SCSI adapter device driver follows. In this sequence, routine names preceded by a **dd\_** are part of the SCSI device driver, where as those preceded by a **sc\_** are part of the SCSI adapter device driver.

1. The SCSI device driver receives a call to its **dd\_strategy** routine; any required internal queuing occurs in this routine. The **dd\_strategy** entry point then triggers the operation by calling the **dd\_start** entry point. The **dd\_start** routine invokes the **sc\_strategy** entry point by calling the **devstrategy** kernel service with the relevant **sc\_buf** structure as a parameter.
2. The **sc\_strategy** entry point initially checks the **sc\_buf** structure for validity. These checks include validating the `devno` field, matching the SCSI ID/LUN to internal tables for configuration purposes, and validating the request size.
3. Although the SCSI adapter device driver cannot reorder transactions, it does perform queue chaining. If no other transactions are pending for the requested device, the **sc\_strategy** routine immediately calls the **sc\_start** routine with the new transaction. If there are other transactions pending, the new transaction is added to the tail of the device chain.
4. At each interrupt, the **sc\_intr** interrupt handler verifies the current status. The SCSI adapter device driver fills in the `sc_buf.status_validity` field, updating the `scsi_status` and `general_card_status` fields as required.
5. The SCSI adapter device driver also fills in the `bufstruct.b_resid` field with the number of bytes not transferred from the request. If all the data was transferred, the `b_resid` field is set to a value of 0. When a transaction completes, the **sc\_intr** routine causes the **sc\_buf** entry to be removed from the device queue and calls the **iodone** kernel service, passing the just dequeued **sc\_buf** structure for the device as the parameter.

The **sc\_start** routine is then called again to process the next transaction on the device queue. The **iodone** kernel service calls the SCSI device driver **dd\_iodone** entry point, signaling the SCSI device driver that the particular transaction has completed.

6. The SCSI device driver **dd\_iodone** routine investigates the I/O completion codes in the **sc\_buf** status entries and performs error recovery, if required. If the operation completed correctly, the SCSI device driver dequeues the original buffer structures. It calls the **iodone** kernel service with the original buffer pointers to notify the originator of the request.

---

## Understanding SCSI Device Driver Internal Commands

During initialization, error recovery, and open or close operations, SCSI device drivers initiate some transactions not directly related to an operating system request. These transactions are called *internal commands* and are relatively simple to handle.

Internal commands differ from operating system-initiated transactions in several ways. The primary difference is that the SCSI device driver is required to generate a **struct buf** that is not related to a specific request. Also, the actual SCSI commands are typically more control-oriented than data transfer-related.

The only special requirement for commands with short data-phase transfers (less than or equal to 256 bytes) is that the SCSI device driver must have pinned the memory being transferred into or out of system memory pages. However, due to system hardware considerations, additional precautions must be taken for data transfers into system memory pages when the transfers are larger than 256 bytes. The problem is that any system memory area with a DMA data operation in progress causes the entire memory page that contains it to become inaccessible.

As a result, a SCSI device driver that initiates an internal command with more than 256 bytes must have preallocated and pinned an area of some multiple whose size is the system page size. The driver must not place in this area any other data areas that it may need to access while I/O is being performed into or out of that page. Memory pages so allocated must be avoided by the device driver from the moment the transaction is passed to the adapter device driver until the device driver **iodone** routine is called for the transaction (and for any other transactions to those pages).

---

## Understanding the Execution of Initiator I/O Requests

During normal processing, many transactions are queued in the SCSI device driver. As the SCSI device driver processes these transactions and passes them to the SCSI adapter device driver, the SCSI device driver moves them to the in-process queue. When the SCSI adapter device driver returns through the **iodone** service with one of these transactions, the SCSI device driver either recovers any errors on the transaction or returns using the **iodone** kernel service to the calling level.

The SCSI device driver can send only one **sc\_buf** structure per call to the SCSI adapter device driver. Thus, the **sc\_buf.bufstruct.av\_forw** pointer should be null when given to the SCSI adapter device driver, which indicates that this is the only request. The SCSI device driver can queue multiple **sc\_buf** requests by making multiple calls to the SCSI adapter device driver strategy routine.

## Spanned (Consolidated) Commands

Some kernel operations might be composed of sequential operations to a device. For example, if consecutive blocks are written to disk, blocks might or might not be in physically consecutive buffer pool blocks.

To enhance SCSI bus performance, the SCSI device driver should consolidate multiple queued requests when possible into a single SCSI command. To allow the SCSI adapter device driver the ability to handle the scatter and gather operations required, the **sc\_buf.bp** should always point to the first **buf** structure entry for the spanned transaction. A null-terminated list of additional **struct buf** entries should be chained from the first field through the **buf.av\_forw** field to give the SCSI adapter device driver enough information to perform the DMA scatter and gather operations required. This information must include at least the buffer's starting address, length, and cross-memory descriptor.

The spanned requests should always be for requests in either the read or write direction but not both, because the SCSI adapter device driver must be given a single SCSI command to handle the requests. The spanned request should always consist of complete I/O requests (including the additional **struct buf** entries). The SCSI device driver should not attempt to use partial requests to reach the maximum transfer size.

The maximum transfer size is actually adapter-dependent. The **IOCINFO** ioctl operation can be used to discover the SCSI adapter device driver's maximum allowable transfer size. To ease the design, implementation, and testing of components that might need to interact with multiple SCSI-adapter device



drivers, a required minimum size has been established that all SCSI adapter device drivers must be capable of supporting. The value of this minimum/maximum transfer size is defined as the following value in the `/usr/include/sys/scsi.h` file:

```
SC_MAXREQUEST      /* maximum transfer request for a single */  
                   /* SCSI command (in bytes) */
```

If a transfer size larger than the supported maximum is attempted, the SCSI adapter device driver returns a value of **EINVAL** in the `sc_buf.bufstruct.b_error` field.

Due to system hardware requirements, the SCSI device driver must consolidate only commands that are memory page-aligned at both their starting and ending addresses. Specifically, this applies to the consolidation of *inner* memory buffers. The ending address of the first buffer and the starting address of all subsequent buffers should be memory page-aligned. However, the starting address of the first memory buffer and the ending address of the last do not need to be aligned so.

The purpose of consolidating transactions is to decrease the number of SCSI commands and bus phases required to perform the required operation. The time required to maintain the simple chain of **buf** structure entries is significantly less than the overhead of multiple (even two) SCSI bus transactions.

## Fragmented Commands

Single I/O requests larger than the maximum transfer size must be divided into smaller requests by the SCSI device driver. For calls to a SCSI device driver's character I/O (read/write) entry points, the **uphysio** kernel service can be used to break up these requests. For a *fragmented command* such as this, the `sc_buf.bp` field should be null so that the SCSI adapter device driver uses only the information in the **sc\_buf** structure to prepare for the DMA operation.

## Gathered Write Commands

The gathered write commands facilitate communication applications that are required to send header and trailer messages with data buffers. These headers and trailers are typically the same or similar for each transfer. Therefore, there might be a single copy of these messages but multiple data buffers.

The gathered write commands, accessed through the `sc_buf.resvd1` field, differ from the spanned commands, accessed through the `sc_buf.bp` field, in several ways:

- Gathered write commands can transfer data regardless of address alignment, where as spanned commands must be memory page-aligned in address and length, making small transfers difficult.
- Gathered write commands can be implemented either in software (which requires the extra step of copying the data to temporary buffers) or hardware. Spanned commands can be implemented in system hardware due to address-alignment requirements. As a result, spanned commands are potentially faster to run.
- Gathered write commands are not able to handle read requests. Spanned commands can handle both read and write requests.
- Gathered write commands can be initiated only on the process level, but spanned commands can be initiated on either the process or interrupt level.

To execute a gathered write command, the SCSI device driver must:

- Fill in the `resvd1` field with a pointer to the **uio** struct
- Call the SCSI adapter device driver on the same process level with the **sc\_buf** structure in question
- Be attempting a write
- Not have put a non-null value in the `sc_buf.bp` field

If any of these conditions are not met, the gathered write commands do not succeed and the `sc_buf.bufstruct.b_error` is set to **EINVAL**.

This interface allows the SCSI adapter device driver to perform the gathered write commands in both software or and hardware as long as the adapter supports this capability. Because the gathered write commands can be performed in software (by using such kernel services as **uiomove**), the contents of the `resvd1` field and the **uio** struct can be altered. Therefore, the caller must restore the contents of both the `resvd1` field and the **uio** struct before attempting a retry. Also, the retry must occur from the process level; it must not be performed from the caller's **iodone** subroutine.

To support SCSI adapter device drivers that perform the gathered write commands in software, additional return values in the `sc_buf.bufstruct.b_error` field are possible when gathered write commands are unsuccessful.

**ENOMEM** Error due to lack of system memory to perform copy.  
**EFAULT** Error due to memory copy problem.

**Note:** The gathered write command facility is optional for both the SCSI device driver and the SCSI adapter device driver. Attempting a gathered write command to a SCSI adapter device driver that does not support gathered write can cause a system crash. Therefore, any SCSI device driver must issue a **SCIOGTHW** ioctl operation to the SCSI adapter device driver before using gathered writes. A SCSI adapter device driver that supports gathered writes must support the **SCIOGTHW** ioctl as well. The ioctl returns a successful return code if gathered writes are supported. If the ioctl fails, the SCSI device driver must not attempt a gathered write. Typically, a SCSI device driver places the **SCIOGTHW** call in its open routine for device instances that it will send gathered writes to.

---

## SCSI Command Tag Queuing

**Note:** This operation is not supported by all SCSI I/O controllers.

SCSI command tag queuing refers to queuing multiple commands to a SCSI device. Queuing to the SCSI device can improve performance because the device itself determines the most efficient way to order and process commands. SCSI devices that support command tag queuing can be divided into two classes: those that clear their queues on error and those that do not. Devices that do not clear their queues on error resume processing of queued commands when the error condition is cleared typically by receiving the next command. Devices that do clear their queues flush all commands currently outstanding.

Command tag queuing requires the SCSI adapter, the SCSI device, the SCSI device driver, and the SCSI adapter driver to support this capability. For a SCSI device driver to queue multiple commands to a SCSI device (that supports command tag queuing), it must be able to provide at least one of the following values in the `sc_buf.q_tag_msg`: **SC\_SIMPLE\_Q**, **SC\_HEAD\_OF\_Q**, or **SC\_ORDERED\_Q**. The SCSI disk device driver and SCSI adapter driver do support this capability. This implementation provides some queuing-specific changeable attributes for disks that can queue commands. With this information, the disk device driver attempts to queue to the disk, first by queuing commands to the adapter driver. The SCSI adapter driver then queues these commands to the adapter, providing that the adapter supports command tag queuing. If the SCSI adapter does not support command tag queuing, then the SCSI adapter driver sends only one command at a time to the SCSI adapter and so multiple commands are not queued to the SCSI disk.

---

## Understanding the `sc_buf` Structure

The **sc\_buf** structure is used for communication between the SCSI device driver and the SCSI adapter device driver during an initiator I/O request. This structure is passed to and from the strategy routine in the same way a standard driver uses a **struct buf** structure.

## Fields in the `sc_buf` Structure

The `sc_buf` structure contains certain fields used to pass a SCSI command and associated parameters to the SCSI adapter device driver. Other fields within this structure are used to pass returned status back to the SCSI device driver. The `sc_buf` structure is defined in the `/usr/include/sys/scsi.h` file.

Fields in the `sc_buf` structure are used as follows:

1. Reserved fields should be set to a value of 0, except where noted.
2. The `bufstruct` field contains a copy of the standard `buf` buffer structure that documents the I/O request. Included in this structure, for example, are the buffer address, byte count, and transfer direction. The `b_work` field in the `buf` structure is reserved for use by the SCSI adapter device driver. The current definition of the `buf` structure is in the `/usr/include/sys/buf.h` include file.
3. The `bp` field points to the original buffer structure received by the SCSI Device Driver from the caller, if any. This can be a chain of entries in the case of spanned transfers (SCSI commands that transfer data from or to more than one system-memory buffer). A null pointer indicates a nonspanned transfer. The null value specifically tells the SCSI adapter device driver that all the information needed to perform the DMA data transfer is contained in the `bufstruct` fields of the `sc_buf` structure. If the `bp` field is set to a non-null value, the `sc_buf.resvd1` field must have a value of null, or else the operation is not allowed.
4. The `scsi_command` field, defined as a `scsi` structure, contains, for example, the SCSI ID, SCSI command length, SCSI command, and a flag variable:
  - a. The `scsi_length` field is the number of bytes in the actual SCSI command. This is normally 6, 10, or 12 (decimal).
  - b. The `scsi_id` field is the SCSI physical unit ID.
  - c. The `scsi_flags` field contains the following bit flags:

### **SC\_NODISC**

Do not allow the target to disconnect during this command.

### **SC\_ASYNC**

Do not allow the adapter to negotiate for synchronous transfer to the SCSI device.

During normal use, the `SC_NODISC` bit should not be set. Setting this bit allows a device executing commands to monopolize the SCSI bus. Sometimes it is desirable for a particular device to maintain control of the bus once it has successfully arbitrated for it; for instance, when this is the only device on the SCSI bus or the only device that will be in use. For performance reasons, it might not be desirable to go through SCSI selections again to save SCSI bus overhead on each command.

Also during normal use, the `SC_ASYNC` bit must not be set. It should be set only in cases where a previous command to the device ended in an unexpected SCSI bus free condition. This condition is noted as `SC_SCSI_BUS_FAULT` in the `general_card_status` field of the `sc_cmd` structure. Because other errors might also result in the `SC_SCSI_BUS_FAULT` flag being set, the `SC_ASYNC` bit should only be set on the last retry of the failed command.

- d. The `sc_cmd` structure contains the physical SCSI command block. The 6 to 12 bytes of a single SCSI command are stored in consecutive bytes, with the op code and logical unit identified individually. The `sc_cmd` structure contains the following fields:
  - The `scsi_op_code` field specifies the standard SCSI op code for this command.
  - The `lun` field specifies the standard SCSI logical unit for the physical SCSI device controller. Typically, there will be one LUN per controller (LUN=0, for example) for devices with imbedded controllers. Only the upper 3 bits of this field contain the actual LUN ID. If addressing LUN's 0 - 7, this `lun` field should always be filled in with the LUN value. When addressing LUN's 8 - 31, this `lun` field should be set to 0 and the LUN value should be placed into the `sc_buf.lun` field described in this section.
  - The `scsi_bytes` field contains the remaining command-unique bytes of the SCSI command block. The actual number of bytes depends on the value in the `scsi_op_code` field.

- The `resvd1` field is set to a non-null value to indicate a request for a gathered write. A gathered write means the SCSI command conducts a system-to-device data transfer where multiple, noncontiguous system buffers contain the write data. This data is transferred in order as a single data transfer for the SCSI command in this `sc_buf` structure.

The contents of the `resvd1` field, if non-null, must be a pointer to the `uio` structure that is passed to the SCSI device driver. The SCSI adapter device driver treats the `resvd1` field as a pointer to a `uio` structure that accesses the `iovec` structures containing pointers to the data. There are no address-alignment restrictions on the data in the `iovec` structures. The only restriction is that the total transfer length of all the data must not exceed the maximum transfer length for the adapter device driver.

The `sc_buf.bufstruct.b_un.b_addr` field, which normally contains the starting system-buffer address, is ignored and can be altered by the SCSI adapter device driver when the `sc_buf` is returned. The `sc_buf.bufstruct.b_bcount` field should be set by the caller to the total transfer length for the data.

5. The `timeout_value` field specifies the time-out limit (in seconds) to be used for completion of this command. A time-out value of 0 means no time-out is applied to this I/O request.
6. The `status_validity` field contains an output parameter that can have one of the following bit flags as a value:

#### **SC\_SCSI\_ERROR**

The `scsi_status` field is valid.

#### **SC\_ADAPTER\_ERROR**

The `general_card_status` field is valid.

7. The `scsi_status` field in the `sc_buf` structure is an output parameter that provides valid SCSI command completion status when its `status_validity` bit is nonzero. The `sc_buf.bufstruct.b_error` field should be set to **EIO** anytime the `scsi_status` field is valid. Typical status values include:

#### **SC\_GOOD\_STATUS**

The target successfully completed the command.

#### **SC\_CHECK\_CONDITION**

The target is reporting an error, exception, or other conditions.

#### **SC\_BUSY\_STATUS**

The target is currently busy and cannot accept a command now.

#### **SC\_RESERVATION\_CONFLICT**

The target is reserved by another initiator and cannot be accessed.

#### **SC\_COMMAND\_TERMINATED**

The target terminated this command after receiving a terminate I/O process message from the SCSI adapter.

#### **SC\_QUEUE\_FULL**

The target's command queue is full, so this command is returned.

8. The `general_card_status` field is an output parameter that is valid when its `status_validity` bit is nonzero. The `sc_buf.bufstruct.b_error` field should be set to **EIO** anytime the `general_card_status` field is valid. This field contains generic SCSI adapter card status. It is intentionally general in coverage so that it can report error status from any typical SCSI adapter.

If an error is detected during execution of a SCSI command, and the error prevented the SCSI command from actually being sent to the SCSI bus by the adapter, then the error should be processed or recovered, or both, by the SCSI adapter device driver.

If it is recovered successfully by the SCSI adapter device driver, the error is logged, as appropriate, but is not reflected in the `general_card_status` byte. If the error cannot be recovered by the SCSI adapter device driver, the appropriate `general_card_status` bit is set and the `sc_buf` structure is returned to the SCSI device driver for further processing.

If an error is detected after the command was actually sent to the SCSI device, then it should be processed or recovered, or both, by the SCSI device driver.

For error logging, the SCSI adapter device driver logs SCSI bus- and adapter-related conditions, where as the SCSI device driver logs SCSI device-related errors. In the following description, a capital letter "A" after the error name indicates that the SCSI adapter device driver handles error logging. A capital letter "H" indicates that the SCSI device driver handles error logging.

Some of the following error conditions indicate a SCSI device failure. Others are SCSI bus- or adapter-related.

**SC\_HOST\_IO\_BUS\_ERR (A)**

The system I/O bus generated or detected an error during a DMA or Programmed I/O (PIO) transfer.

**SC\_SCSI\_BUS\_FAULT (H)**

The SCSI bus protocol or hardware was unsuccessful.

**SC\_CMD\_TIMEOUT (H)**

The command timed out before completion.

**SC\_NO\_DEVICE\_RESPONSE (H)**

The target device did not respond to selection phase.

**SC\_ADAPTER\_HDW\_FAILURE (A)**

The adapter indicated an onboard hardware failure.

**SC\_ADAPTER\_SFW\_FAILURE (A)**

The adapter indicated microcode failure.

**SC\_FUSE\_OR\_TERMINAL\_PWR (A)**

The adapter indicated a blown terminator fuse or bad termination.

**SC\_SCSI\_BUS\_RESET (A)**

The adapter indicated the SCSI bus has been reset.

9. When the SCSI device driver queues multiple transactions to a device, the `adap_q_status` field indicates whether or not the SCSI adapter driver has cleared its queue for this device after an error has occurred. The flag of **SC\_DID\_NOT\_CLEAR\_Q** indicates that the SCSI adapter driver has not cleared its queue for this device and that it is in a halted state (so none of the pending queued transactions are sent to the device).
10. The `lun` field provides addressability of up to 32 logical units (LUNs). This field specifies the standard SCSI LUN for the physical SCSI device controller. If addressing LUN's 0 - 7, both this `lun` field (`sc_buf.lun`) and the `lun` field located in the `scsi_command` structure (`sc_buf.scsi_command.scsi_cmd.lun`) should be set to the LUN value. If addressing LUN's 8 - 31, this `lun` field (`sc_buf.lun`) should be set to the LUN value and the `lun` field located in the `scsi_command` structure (`sc_buf.scsi_command.scsi_cmd.lun`) should be set to 0.

Logical Unit Numbers (LUNs)		
lun Fields	LUN 0 - 7	LUN 8 - 31
<code>sc_buf.lun</code>	<i>LUN Value</i>	<i>LUN Value</i>
<code>sc_buf.scsi_command.scsi_cmd.lun</code>	<i>LUN Value</i>	0

**Note:** *LUN value* is the current value of LUN.

11. The `q_tag_msg` field indicates if the SCSI adapter can attempt to queue this transaction to the device. This information causes the SCSI adapter to fill in the Queue Tag Message Code of the queue tag message for a SCSI command. The following values are valid for this field:

**SC\_NO\_Q**

Specifies that the SCSI adapter does not send a queue tag message for this command, and

so the device does not allow more than one SCSI command on its command queue. This value must be used for all commands sent to SCSI devices that do not support command tag queuing.

#### **SC\_SIMPLE\_Q**

Specifies placing this command in the device's command queue. The device determines the order that it executes commands in its queue. The SCSI-2 specification calls this value the "Simple Queue Tag Message."

#### **SC\_HEAD\_OF\_Q**

Specifies placing this command first in the device's command queue. This command does not preempt an active command at the device, but it is executed before all other commands in the command queue. The SCSI-2 specification calls this value the "Head of Queue Tag Message."

#### **SC\_ORDERED\_Q**

Specifies placing this command in the device's command queue. The device processes these commands in the order that they are received. The SCSI-2 specification calls this value the "Ordered Queue Tag Message."

**Note:** Commands with the value of **SC\_NO\_Q** for the `q_tag_msg` field (except for request sense commands) should not be queued to a device whose queue contains a command with another value for `q_tag_msg`. If commands with the **SC\_NO\_Q** value (except for request sense) are sent to the device, then the SCSI device driver must make sure that no active commands are using different values for `q_tag_msg`. Similarly, the SCSI device driver must also make sure that a command with a `q_tag_msg` value of **SC\_ORDERED\_Q**, **SC\_HEAD\_OF\_Q**, or **SC\_SIMPLE\_Q** is not sent to a device that has a command with the `q_tag_msg` field of **SC\_NO\_Q**.

12. The `flags` field contains bit flags sent from the SCSI device driver to the SCSI adapter device driver. The following flags are defined:

#### **SC\_RESUME**

When set, means the SCSI adapter device driver should resume transaction queuing for this ID/LUN. Error recovery is complete after a **SCIOHALT** operation, check condition, or severe SCSI bus error. This flag is used to restart the SCSI adapter device driver following a reported error.

#### **SC\_DELAY\_CMD**

When set, means the SCSI adapter device driver should delay sending this command (following a SCSI reset or BDR to this device) by at least the number of seconds specified to the SCSI adapter device driver in its configuration information. For SCSI devices that do not require this function, this flag should not be set.

#### **SC\_Q\_CLR**

When set, means the SCSI adapter driver should clear its transaction queue for this ID/LUN. The transaction containing this flag setting does not require an actual SCSI command in the `sc_buf` because it is flushed back to the SCSI device driver with the rest of the transactions for this ID/LUN. However, this transaction must have the SCSI ID field (`sc_buf.scsi_command.scsi_id`) and the LUN fields (`sc_buf.scsi_command.scsi_cmd.lun` and `sc_buf.lun`) filled in with the device's SCSI ID and logical unit number (LUN). This flag is valid only during error recovery of a check condition or command terminated at a command tag queuing device when the **SC\_DID\_NOT\_CLR\_Q** flag is set in the `sc_buf.adap_q_status` field.

**Note:** When addressing LUN's 8 - 31, be sure to see the description of the `sc_buf.lun` field within the `sc_buf` structure.

#### **SC\_Q\_RESUME**

When set, means that the SCSI adapter driver should resume its halted transaction queue for this ID/LUN. The transaction containing this flag setting does not require an actual SCSI

command to be sent to the SCSI adapter driver. However, this transaction must have the `sc_buf.scsi_command.scsi_id` and `sc_buf.scsi_command.scsi_cmd.lun` fields filled in with the device's SCSI ID and logical unit number. If the transaction containing this flag setting is the first issued by the SCSI device driver after it receives an error (indicating that the adapter driver's queue is halted), then the **SC\_RESUME** flag must be set also.

**Note:** When addressing LUN's 8 - 31, be sure to see the description of the `sc_buf.lun` field within the `sc_buf` structure.

---

## Other SCSI Design Considerations

The following topics cover design considerations of SCSI device and adapter device drivers:

- Responsibilities of the SCSI Device Driver
- SCSI Options to the `openx` Subroutine
- Using the `SC_FORCED_OPEN` Option
- Using the `SC_RETAIN_RESERVATION` Option
- Using the `SC_DIAGNOSTIC` Option
- Using the `SC_NO_RESERVE` Option
- Using the `SC_SINGLE` Option
- Closing the SCSI Device
- SCSI Error Processing
- Device Driver and Adapter Device Driver Interfaces
- Performing SCSI Dumps

## Responsibilities of the SCSI Device Driver

SCSI device drivers are responsible for the following actions:

- Interfacing with block I/O and logical-volume device-driver code in the operating system.
- Translating I/O requests from the operating system into SCSI commands suitable for the particular SCSI device. These commands are then given to the SCSI adapter device driver for execution.
- Issuing any and all SCSI commands to the attached device. The SCSI adapter device driver sends no SCSI commands except those it is directed to send by the calling SCSI device driver.
- Managing SCSI device reservations and releases. In the operating system, it is assumed that other SCSI initiators might be active on the SCSI bus. Usually, the SCSI device driver reserves the SCSI device at open time and releases it at close time (except when told to do otherwise through parameters in the SCSI device driver interface). Once the device is reserved, the SCSI device driver must be prepared to reserve the SCSI device again whenever a Unit Attention condition is reported through the SCSI request-sense data.

## SCSI Options to the `openx` Subroutine

SCSI device drivers in the operating system must support eight defined extended options in their `open` routine (that is, an **openx** subroutine). Additional extended options to the `open` are also allowed, but they must not conflict with predefined `open` options. The defined extended options are bit flags in the `ext` `open` parameter. These options can be specified singly or in combination with each other. The required `ext` options are defined in the `/usr/include/sys/scsi.h` header file and can have one of the following values:

<b>SC_FORCED_OPEN</b>	Do not honor device reservation-conflict status.
<b>SC_RETAIN_RESERVATION</b>	Do not release SCSI device on close.
<b>SC_DIAGNOSTIC</b>	Enter diagnostic mode for this device.
<b>SC_NO_RESERVE</b>	Prevents the reservation of the device during an <b>openx</b> subroutine call to that device. Allows multiple hosts to share a device.
<b>SC_SINGLE</b>	Places the selected device in Exclusive Access mode.

<b>SC_RESV_05</b>	Reserved for future expansion.
<b>SC_RESV_07</b>	Reserved for future expansion.
<b>SC_RESV_08</b>	Reserved for future expansion.

## Using the **SC\_FORCED\_OPEN** Option

The **SC\_FORCED\_OPEN** option causes the SCSI device driver to call the SCSI adapter device driver's Bus Device Reset ioctl (**SCIORESET**) operation on the first open. This forces the device to release another initiator's reservation. After the **SCIORESET** command is completed, other SCSI commands are sent as in a normal open. If any of the SCSI commands fail due to a reservation conflict, the open registers the failure as an **EBUSY** status. This is also the result if a reservation conflict occurs during a normal open. The SCSI device driver should require the caller to have appropriate authority to request the **SC\_FORCED\_OPEN** option because this request can force a device to drop a SCSI reservation. If the caller attempts to initiate this system call without the proper authority, the SCSI device driver should return a value of -1, with the **errno** global variable set to a value of **EPERM**.

## Using the **SC\_RETAIN\_RESERVATION** Option

The **SC\_RETAIN\_RESERVATION** option causes the SCSI device driver not to issue the SCSI release command during the close of the device. This guarantees a calling program control of the device (using SCSI reservation) through open and close cycles. For shared devices (for example, disk or CD-ROM), the SCSI device driver must OR together this option for all opens to a given device. If any caller requests this option, the **close** routine does not issue the release even if other opens to the device do not set **SC\_RETAIN\_RESERVATION**. The SCSI device driver should require the caller to have appropriate authority to request the **SC\_RETAIN\_RESERVATION** option because this request can allow a program to monopolize a device (for example, if this is a nonshared device). If the caller attempts to initiate this system call without the proper authority, the SCSI device driver should return a value of -1, with the **errno** global variable set to a value of **EPERM**.

## Using the **SC\_DIAGNOSTIC** Option

The **SC\_DIAGNOSTIC** option causes the SCSI device driver to enter Diagnostic mode for the given device. This option directs the SCSI device driver to perform only minimal operations to open a logical path to the device. No SCSI commands should be sent to the device in the **open** or **close** routine when the device is in Diagnostic mode. One or more ioctl operations should be provided by the SCSI device driver to allow the caller to issue SCSI commands to the attached device for diagnostic purposes.

The **SC\_DIAGNOSTIC** option gives the caller an exclusive open to the selected device. This option requires appropriate authority to run. If the caller attempts to initiate this system call without the proper authority, the SCSI device driver should return a value of -1, with the **errno** global variable set to a value of **EPERM**. The **SC\_DIAGNOSTIC** option may be run only if the device is not already opened for normal operation. If this ioctl operation is attempted when the device is already opened, or if an **openx** call with the **SC\_DIAGNOSTIC** option is already in progress, a return value of -1 should be passed, with the **errno** global variable set to a value of **EACCES**. Once successfully opened with the **SC\_DIAGNOSTIC** flag, the SCSI device driver is placed in Diagnostic mode for the selected device.

## Using the **SC\_NO\_RESERVE** Option

The **SC\_NO\_RESERVE** option causes the SCSI device driver not to issue the SCSI reserve command during the opening of the device and not to issue the SCSI release command during the close of the device. This allows multiple hosts to share the device. The SCSI device driver should require the caller to have appropriate authority to request the **SC\_NO\_RESERVE** option, because this request allows other hosts to modify data on the device. If a caller does this kind of request then the caller must ensure data integrity between multiple hosts. If the caller attempts to initiate this system call without the proper authority, the SCSI device driver should return a value of -1, with the **errno** global variable set to a value of **EPERM**.



## Using the SC\_SINGLE Option

The **SC\_SINGLE** option causes the SCSI device driver to issue a normal open, but does not allow another caller to issue another open until the first caller has closed the device. This request gives the caller an exclusive open to the selected device. If this **openx** is attempted when the device is already open, a return value of -1 is passed, with the **errno** global variable set to a value of **EBUSY**.

Once successfully opened, the device is placed in Exclusive Access mode. If another caller tries to do any type of **open**, a return value of -1 is passed, with the **errno** global variable set to a value of **EACCES**.

The remaining options for the *ext* parameter are reserved for future requirements.

**Implementation note:** The following table shows how the various combinations of *ext* options should be handled in the SCSI device driver.

EXT OPTIONS <i>openx ext option</i>	Device Driver Action
none	Open: normal. Close: normal.
diag	Open: no SCSI commands. Close: no SCSI commands.
diag + force	Open: issue SCIORESET otherwise, no SCSI commands issued. Close: no SCSI commands.
diag + force + no_reserve	Open: issue SCIORESET; otherwise, no SCSI commands issued. Close: no SCSI commands.
diag + force + no_reserve + single	Open: issue SCIORESET; otherwise, no SCSI commands issued. Close: no SCSI commands.
diag + force +retain	Open: issue SCIORESET; otherwise, no SCSI commands issued. Close: no SCSI commands.
diag + force +retain + no_reserve	Open: issue SCIORESET; otherwise, no SCSI commands issued. Close: no SCSI commands.
diag + force +retain + no_reserve + single	Open: issue SCIORESET; otherwise, no SCSI commands issued. Close: no SCSI commands.
diag + force +retain + single	Open: issue SCIORESET; otherwise, no SCSI commands issued. Close: no SCSI commands.
diag + force + single	Open: issue SCIORESET; otherwise, no SCSI commands issued. Close: no SCSI commands.
diag+no_reserve	Open: no SCSI commands. Close: no SCSI commands.
diag + retain	Open: no SCSI commands. Close: no SCSI commands.
diag + retain + no_reserve	Open: no SCSI commands. Close: no SCSI commands.
diag + retain + no_reserve + single	Open: no SCSI commands. Close: no SCSI commands.
diag + retain + single	Open: no SCSI commands. Close: no SCSI commands.
diag + single	Open: no SCSI commands. Close: no SCSI commands.
diag + single + no_reserve	Open: no SCSI commands. Close: no SCSI commands.
force	Open: normal, except SCIORESET issued prior to any SCSI commands. Close: normal.
force + no_reserve	Open: normal except SCIORESET issued prior to any SCSI commands. No RESERVE command issued. Close: normal except no RELEASE.
force + retain	Open: normal, except SCIORESET issued prior to any SCSI commands. Close: no RELEASE.

EXT OPTIONS <i>openx ext option</i>	Device Driver Action
force + retain + no_reserve	Open: normal except SCIORESET issued prior to any SCSI commands. No RESERVE command issued. Close: no RELEASE.
force + retain + no_reserve + single	Open: normal except SCIORESET issued prior to any SCSI commands. No RESERVE command issued. Close: no RELEASE.
force + retain + single	Open: normal except SCIORESET issued prior to any SCSI commands. Close: no RELEASE.
force + single	Open: normal except SCIORESET issued prior to any SCSI commands. Close: normal.
force + single + no_reserve	Open: normal except SCIORESET issued prior to any SCSI commands. No RESERVE command issued. Close: no RELEASE.
no_reserve	Open: no RESERVE. Close: no RELEASE.
retain	Open: normal. Close: no RELEASE.
retain + no_reserve	Open: no RESERVE. Close: no RELEASE.
retain + single	Open: normal. Close: no RELEASE.
retain + single + no_reserve	Open: normal except no RESERVE command issued. Close: no RELEASE.
single	Open: normal. Close: normal.
single + no_reserve	Open: no RESERVE. Close: no RELEASE.

## Closing the SCSI Device

When a SCSI device driver is preparing to close a device through the SCSI adapter device driver, it must ensure that all transactions are complete. When the SCSI adapter device driver receives a **SCIOSTOP** ioctl operation and there are pending I/O requests, the ioctl operation does not return until all have completed. New requests received during this time are rejected from the adapter device driver's **ddstrategy** routine.

When the SCSI adapter device driver receives an **SCIOSTOPTGT** ioctl operation, it must forcibly free any receive data buffers that have been queued to the SCSI device driver for this device and have not been returned to the SCSI adapter device driver through the buffer free routine. The SCSI device driver is responsible for making sure all the receive data buffers are freed before calling the **SCIOSTOPTGT** ioctl operation. However, the SCSI adapter device driver must check that this is done, and, if necessary, forcibly free the buffers. The buffers must be freed because those not freed result in memory areas being permanently lost to the system (until the next boot).

To allow the SCSI adapter device driver to free buffers that are sent to the SCSI device driver but never returned, it must track which **tm\_bufs** are currently queued to the SCSI device driver. Tracking **tm\_bufs** requires the SCSI adapter device driver to violate the general SCSI rule, which states the SCSI adapter device driver should not modify the **tm\_bufs** structure while it is queued to the SCSI device driver. This exception to the rule is necessary because it is never acceptable not to free memory allocated from the system.

## SCSI Error Processing

It is the responsibility of the SCSI device driver to process SCSI check conditions and other returned errors properly. The SCSI adapter device driver only passes SCSI commands without otherwise processing them and is not responsible for device error recovery.

## Device Driver and Adapter Device Driver Interfaces

The SCSI device drivers can have both character (raw) and block special files in the `/dev` directory. The SCSI adapter device driver has only character (raw) special files in the `/dev` directory and has only the **ddconfig**, **ddopen**, **ddclose**, **dddump**, and **ddioctl** entry points available to operating system programs. The **ddread** and **ddwrite** entry points are not implemented.

Internally, the **devsw** table has entry points for the **ddconfig**, **ddopen**, **ddclose**, **dddump**, **ddioctl**, and **ddstrategy** routines. The SCSI device drivers pass their SCSI commands to the SCSI adapter device driver by calling the SCSI adapter device driver **ddstrategy** routine. (This routine is unavailable to other operating system programs due to the lack of a block-device special file.)

Access to the SCSI adapter device driver's **ddconfig**, **ddopen**, **ddclose**, **dddump**, **ddioctl**, and **ddstrategy** entry points by the SCSI device drivers is performed through the kernel services provided. These include such services as **fp\_opendev**, **fp\_close**, **fp\_ioctl**, **devdump**, and **devstrategy**.

## Performing SCSI Dumps

A SCSI adapter device driver must have a **dddump** entry point if it is used to access a system dump device. A SCSI device driver must have a **dddump** entry point if it drives a dump device. Examples of dump devices are disks and tapes.

**Note:** SCSI adapter-device-driver writers should be aware that system services providing interrupt and timer services are unavailable for use in the dump routine. Kernel DMA services are assumed to be available for use by the dump routine. The SCSI adapter device driver should be designed to ignore extra **DUMPINIT** and **DUMPSTART** commands to the **dddump** entry point.

The **DUMPQUERY** option should return a minimum transfer size of 0 bytes, and a maximum transfer size equal to the maximum transfer size supported by the SCSI adapter device driver.

Calls to the SCSI adapter device driver **DUMPWRITE** option should use the *arg* parameter as a pointer to the **sc\_buf** structure to be processed. Using this interface, a SCSI **write** command can be run on a previously started (opened) target device. The *uiop* parameter is ignored by the SCSI adapter device driver during the **DUMPWRITE** command. Spanned, or consolidated, commands are not supported using the **DUMPWRITE** option. Gathered **write** commands are also not supported using the **DUMPWRITE** option. No queuing of **sc\_buf** structures is supported during dump processing because the dump routine runs essentially as a subroutine call from the caller's dump routine. Control is returned when the entire **sc\_buf** structure has been processed.

**Attention:** Also, both adapter-device-driver and device-driver writers should be aware that any error occurring during the **DUMPWRITE** option is considered unsuccessful. Therefore, no error recovery is employed during the **DUMPWRITE**. Return values from the call to the **dddump** routine indicate the specific nature of the failure.

Successful completion of the selected operation is indicated by a 0 return value to the subroutine. Unsuccessful completion is indicated by a return code set to one of the following values for the **errno** global variable. The various **sc\_buf** status fields, including the **b\_error** field, are not set by the SCSI adapter device driver at completion of the **DUMPWRITE** command. Error logging is, of necessity, not supported during the dump.

- An **errno** value of **EINVAL** indicates that a request that was not valid passed to the SCSI adapter device driver, such as to attempt a **DUMPSTART** command before successfully executing a **DUMPINIT** command.
- An **errno** value of **EIO** indicates that the SCSI adapter device driver was unable to complete the command due to a lack of required resources or an I/O error.
- An **errno** value of **ETIMEDOUT** indicates that the adapter did not respond with completion status before the passed command time-out value expired.

---

## SCSI Target-Mode Overview

**Note:** This operation is not supported by all SCSI I/O controllers.

The SCSI target-mode interface is intended to be used with the SCSI initiator-mode interface to provide the equivalent of a full-duplex communications path between processor type devices. Both communicating devices must support target-mode and initiator-mode. To work with the SCSI subsystem in this manner, an attached device's target-mode and initiator-mode interfaces must meet certain minimum requirements:

- The device's target-mode interface must be capable of receiving and processing at least the following SCSI commands:
  - **send**
  - **request sense**
  - **inquiry**

The data returned by the **inquiry** command must set the peripheral device type field to processor device. The device should support the vendor and product identification fields. Additional functional SCSI requirements, such as SCSI message support, must be addressed by examining the detailed functional specification of the SCSI initiator that the target-mode device is attached to.

- The attached device's initiator mode interface must be capable of sending the following SCSI commands:
  - **send**
  - **request sense**

In addition, the **inquiry** command should be supported by the attached initiator if it needs to identify SCSI target devices. Additional functional SCSI requirements, such as SCSI message support, must be addressed by examining the detailed functional specification of the SCSI target that the initiator-mode device is attached to.

## Configuring and Using SCSI Target Mode

The adapter, acting as either a target or initiator device, requires its own SCSI ID. This ID, as well as the IDs of all attached devices on this SCSI bus, must be unique and between 0 and 7, inclusive. Because each device on the bus must be at a unique ID, the user must complete any installation and configuration of the SCSI devices required to set the correct IDs before physically cabling the devices together. Failure to do so will produce unpredictable results.

SCSI target mode in the SCSI subsystem does not attempt to implement any receive-data protocol, with the exception of actions taken to prevent an application from excessive receive-data-buffer usage. Any protocol required to maintain or otherwise manage the communications of data must be implemented in user-supplied programs. The only delays in receiving data are those inherent in the SCSI subsystem and the hardware environment in which it operates.

The SCSI target mode is capable of simultaneously receiving data from all attached SCSI IDs using SCSI **send** commands. In target-mode, the host adapter is assumed to act as a single SCSI Logical Unit Number (LUN) at its assigned SCSI ID. Therefore, only one logical connection is possible between each attached SCSI initiator on the SCSI Bus and the host adapter. The SCSI subsystem is designed to be fully capable of simultaneously sending SCSI commands in initiator-mode while receiving data in target-mode.

## Managing Receive-Data Buffers

In the SCSI subsystem target-mode interface, the SCSI adapter device driver is responsible for managing the receive-data buffers versus the SCSI device driver because the buffering is dependent upon how the adapter works. It is not possible for the SCSI device driver to run a single approach that is capable of making full use of the performance advantages of various adapters' buffering schemes. With the SCSI adapter device driver layer performing the buffer management, the SCSI device driver can be interfaced to a variety of adapter types and can potentially get the best possible performance out of each adapter. This

approach also allows multiple SCSI target-mode device drivers to be run on top of adapters that use a shared-pool buffer management scheme. This would not be possible if the target-mode device drivers managed the buffers.

## Understanding Target-Mode Data Pacing

Because it is possible for the attached initiator device to send data faster than the host operating system and associated application can process it, eventually the situation arises in which all buffers for this device instance are in use at the same time. There are two possible scenarios:

- The previous **send** command has been received by the adapter, but there is no space for the next **send** command.
- The **send** command is not yet completed, and there is no space for the remaining data.

In both cases, the combination of the SCSI adapter device driver and the SCSI adapter must be capable of stopping the flow of data from the initiator device.

### SCSI Adapter Device Driver

The adapter can handle both cases described previously by simply accepting the **send** command (if newly received) and then disconnecting during the data phase. When buffer space becomes available, the SCSI adapter reconnects and continues the data transfer. As an alternative, when handling a newly received command, a check condition can be given back to the initiator to indicate a lack of resources. The implementation of this alternative is adapter-dependent. The technique of accepting the command and then disconnecting until buffer space is available should result in better throughput, as it avoids both a **request sense** command and the retry of the **send** command.

For adapters allowing a shared pool of buffers to be used for all attached initiators' data transfers, an additional problem can result. If any single initiator instance is allowed to transfer data continually, the entire shared pool of buffers can fill up. These filled-up buffers prevent other initiator instances from transferring data. To solve this problem, the combination of the SCSI adapter device driver and the host SCSI adapter must stop the flow of data from a particular initiator ID on the bus. This could include disconnecting during the data phase for a particular ID but allowing other IDs to continue data transfer. This could begin when the number of **tm\_buf** structures on a target-mode instance's **tm\_buf** queue equals the number of buffers allocated for this device. When a threshold percentage of the number of buffers is processed and returned to the SCSI adapter device driver's buffer-free routine, the ID can be enabled again for the continuation of data transfer.

### SCSI Device Driver

The SCSI device driver can optionally be informed by the SCSI adapter device driver whenever all buffers for this device are in use. This is known as a maximum-buffer-usage event. To pass this information, the SCSI device driver must be registered for notification of asynchronous event status from the SCSI adapter device driver. Registration is done by calling the SCSI adapter device-driver **ioctl** entry point with the **SCIOEVENT** operation. If registering for event notification, the SCSI device driver receives notification of all asynchronous events, not just the maximum buffer usage event.

## Understanding the SCSI Target Mode Device Driver Receive Buffer Routine

The SCSI target-mode device-driver **receive buffer** routine must be a pinned routine that the SCSI adapter device driver can directly address. This routine is called directly from the SCSI adapter device driver hardware interrupt handling routine. The SCSI device driver writer must be aware of how this routine affects the design of the SCSI device driver.

First, because the **receive buffer** routine is running on the hardware interrupt level, the SCSI device driver must limit operations in order to limit routine processing time. In particular, the data copy, which occurs because the data is queued ahead of the user read request, must not occur in the **receive buffer** routine. Data copying in this routine will adversely affect system response time. Data copy is best performed in a

process level SCSI device-driver routine. This routine sleeps, waiting for data, and is awakened by the **receive buffer** routine. Typically, this process level routine is the SCSI device driver's **read** routine.

Second, the **receive buffer** routine is called at the SCSI adapter device driver hardware interrupt level, so care must be taken when disabling interrupts. They must be disabled to the correct level in places in the SCSI device driver's lower run priority routines, which manipulate variables also modified in the **receive buffer** routine. To allow the SCSI device driver to disable to the correct level, the SCSI adapter device-driver writer must provide a configuration database attribute, named **intr\_priority**, that defines the interrupt class, or priority, that the adapter runs on. The SCSI device-driver configuration method should pass this attribute to the SCSI device driver along with other configuration data for the device instance.

Third, the SCSI device-driver writer must follow any other general system rules for writing a routine that must run in an interrupt environment. For example, the routine must not attempt to sleep or wait on I/O operations. It can perform wake-up calls to allow the process level to handle those operations.

Duties of the SCSI device driver **receive buffer** routine include:

- Matching the data with the appropriate target-mode instance.
- Queuing the **tm\_buf** structures to the appropriate target-mode instance.
- Waking up the process-level routine for further processing of the received data.

After the **tm\_buf** structure has been passed to the SCSI device driver **receive buffer** routine, the SCSI device driver is considered to be responsible for it. Responsibilities include processing the data and any error conditions and also maintaining the next pointer for chained **tm\_buf** structures. The SCSI device driver's responsibilities for the **tm\_buf** structures end when it passes the structure back to the SCSI adapter device driver.

Until the **tm\_buf** structure is again passed to the SCSI device driver **receive buffer** routine, the SCSI adapter device driver is considered responsible for it. The SCSI adapter device-driver writer must be aware that during the time the SCSI device driver is responsible for the **tm\_buf** structure, it is still possible for the SCSI adapter device driver to access the structure's contents. Access is possible because only one copy of the structure is in memory, and only a pointer to the structure is passed to the SCSI device driver.

**Note:** Under no circumstances should the SCSI adapter device driver access the structure or modify its contents while the SCSI device driver is responsible for it, or the other way around.

It is recommended that the SCSI device-driver writer implement a threshold level to wake up the process level with available **tm\_buf** structures. This way, processing for some of the buffers, including copying the data to the user buffer, can be overlapped with time spent waiting for more data. It is also recommended the writer implement a threshold level for these buffers to handle cases where the **send** command data length exceeds the aggregate receive-data buffer space. A suggested threshold level is 25% of the device's total buffers. That is, when 25% or more of the number of buffers allocated for this device is queued and no end to the **send** command is encountered, the SCSI device driver receive buffer routine should wake the process level to process these buffers.

## Understanding the **tm\_buf** Structure

The **tm\_buf** structure is used for communication between the SCSI device driver and the SCSI adapter device driver for a target-mode received-data buffer. The **tm\_buf** structure is passed by pointer directly to routines whose entry points have been registered through the **SCIOSTARTTGT** ioctl operation of the SCSI adapter device driver. The SCSI device driver is required to call this ioctl operation when opening a target-mode device instance.

### Fields in the **tm\_buf** Structure

The **tm\_buf** structure contains certain fields used to pass a received data buffer from the SCSI adapter device driver to the SCSI device driver. Other fields are used to pass returned status back to the SCSI device driver. After processing the data, the **tm\_buf** structure is passed back from the SCSI device driver

to the SCSI adapter device driver to allow the buffer to be reused. The **tm\_buf** structure is defined in the `/usr/include/sys/scsi.h` file and contains the following fields:

**Note:** Reserved fields must not be modified by the SCSI device driver, unless noted otherwise. Nonreserved fields can be modified, except where noted otherwise.

1. The `tm_correlator` field is an optional field for the SCSI device driver. This field is a copy of the field with the same name that was passed by the SCSI device driver in the **SCIOSTARTTGT** ioctl. The SCSI device driver should use this field to speed the search for the target-mode device instance the **tm\_buf** structure is associated with. Alternatively, the SCSI device driver can combine the `tm_buf.user_id` and `tm_buf.adap_devno` fields to find the associated device.
2. The `adap_devno` field is the device major and minor numbers of the adapter instance on which this target mode device is defined. This field can be used to find the particular target-mode instance the **tm\_buf** structure is associated with.

**Note:** The SCSI device driver must not modify this field.

3. The `data_addr` field is the kernel space address where the data begins for this buffer.
4. The `data_len` field is the length of valid data in the buffer starting at the **tm\_buf.data\_addr** location in memory.
5. The `user_flag` field is a set of bit flags that can be set to communicate information about this data buffer to the SCSI device driver. Except where noted, one or more of the following flags can be set:

#### **TM\_HASDATA**

Set to indicate a valid **tm\_buf** structure

#### **TM\_MORE\_DATA**

Set if more data is coming (that is, more **tm\_buf** structures) for a particular **send** command. This is only possible for adapters that support spanning the **send** command data across multiple receive buffers. This flag cannot be used with the **TM\_ERROR** flag.

#### **TM\_ERROR**

Set if any error occurred on a particular **send** command. This flag cannot be used with the **TM\_MORE\_DATA** flag.

6. The `user_id` field is set to the SCSI ID of the initiator that sent the data to this target mode instance. If more than one adapter is used for target mode in this system, this ID might not be unique. Therefore, this field must be used in combination with the `tm_buf.adap_devno` field to find the target-mode instance this ID is associated with.

**Note:** The SCSI device driver must not modify this field.

7. The `status_validity` field contains the following bit flag:

#### **SC\_ADAPTER\_ERROR**

Indicates the `tm_buf.general_card_status` is valid.

8. The `general_card_status` field is a returned status field that gives a broad indication of the class of error encountered by the adapter. This field is valid when its status-validity bit is set in the `tm_buf.status_validity` field. The definition of this field is the same as that found in the **sc\_buf** structure definition, except the **SC\_CMD\_TIMEOUT** value is not possible and is never returned for a target-mode transfer.
9. The next field is a **tm\_buf** pointer that is either null, meaning this is the only or last **tm\_buf** structure, or else contains a non-null pointer to the next **tm\_buf** structure.

## Understanding the Running of SCSI Target-Mode Requests

The target-mode interface provided by the SCSI subsystem is designed to handle data reception from SCSI **send** commands. The host SCSI adapter acts as a secondary device that waits for an attached initiator device to issue a SCSI **send** command. The SCSI **send** command data is received by buffers managed by the SCSI adapter device driver. The **tm\_buf** structure is used to manage individual buffers.

For each buffer of data received from an attached initiator, the SCSI adapter device driver passes a **tm\_buf** structure to the SCSI device driver for processing. Multiple **tm\_buf** structures can be linked together and passed to the SCSI device driver at one time. When the SCSI device driver has processed one or more **tm\_buf** structures, it passes the **tm\_buf** structures back to the SCSI adapter device driver so they can be reused.

### Detailed Running of Target-Mode Requests

When a **send** command is received by the host SCSI adapter, data is placed in one or more receive-data buffers. These buffers are made available to the adapter by the SCSI adapter device driver. The procedure by which the data gets from the SCSI bus to the system-memory buffer is adapter-dependent. The SCSI adapter device driver takes the received data and updates the information in one or more **tm\_buf** structures in order to identify the data to the SCSI device driver. This process includes filling the **tm\_correlator**, **adap\_devno**, **data\_addr**, **data\_len**, **user\_flag**, and **user\_id** fields. Error status information is put in the **status\_validity** and **general\_card\_status** fields. The next field is set to null to indicate this is the only element, or set to non-null to link multiple **tm\_buf** structures. If there are multiple **tm\_buf** structures, the final **tm\_buf.next** field is set to null to end the chain. If there are multiple **tm\_buf** structures and they are linked, they must all be from the same initiator SCSI ID. The **tm\_buf.tm\_correlator** field, in this case, has the same value as it does in the **SCIOSTARTTGT** ioctl operation to the SCSI adapter device driver. The SCSI device driver should use this field to speed the search for the target-mode instance designated by this **tm\_buf** structure. For example, when using the value of **tm\_buf.tm\_correlator** as a pointer to the device-information structure associated with this target-mode instance.

Each **send** command, no matter how short its data length, requires its own **tm\_buf** structure. For host SCSI adapters capable of spanning multiple receive-data buffers with data from a single **send** command, the SCSI adapter device driver must set the **TM\_MORE\_DATA** flag in the **tm\_buf.user\_flag** fields of all but the final **tm\_buf** structure holding data for the **send** command. The SCSI device driver must be designed to support the **TM\_MORE\_DATA** flag. Using this flag, the target-mode SCSI device driver can associate multiple buffers with the single transfer they represent. The end of a **send** command will be the boundary used by the SCSI device driver to satisfy a user read request.

The SCSI adapter device driver is responsible for sending the **tm\_buf** structures for a particular initiator SCSI ID to the SCSI device driver in the order they were received. The SCSI device driver is responsible for processing these **tm\_buf** structures in the order they were received. There is no particular ordering implied in the processing of simultaneous **send** commands from different SCSI IDs, as long as the data from an individual SCSI ID's **send** command is processed in the order it was received.

The pointer to the **tm\_buf** structure chain is passed by the SCSI adapter device driver to the SCSI device driver's receive buffer routine. The address of this routine is registered with the SCSI adapter device driver by the SCSI device driver using the **SCIOSTARTTGT** ioctl. The duties of the receive buffer routine include queuing the **tm\_buf** structures and waking up a process-level routine (typically the SCSI device driver's **read** routine) to process the received data.

When the process-level SCSI device driver routine finishes processing one or more **tm\_buf** structures, it passes them to the SCSI adapter device driver's buffer-free routine. The address of this routine is registered with the SCSI device driver in an output field in the structure passed to the SCSI adapter device driver **SCIOSTARTTGT** ioctl operation. The buffer-free routine must be a pinned routine the SCSI device driver can directly access. The buffer-free routine is typically called directly from the SCSI device driver buffer-handling routine. The SCSI device driver chains one or more **tm\_buf** structures by using the next field (a null value for the last **tm\_buf** next field ends the chain). It then passes a pointer, which points to the head of the chain, to the SCSI adapter device driver buffer-free routine. These **tm\_buf** structures must all be for the same target-mode instance. Also, the SCSI device driver must not modify the **tm\_buf.user\_id** or **tm\_buf.adap\_devno** field.

The SCSI adapter device driver takes the **tm\_buf** structures passed to its buffer-free routine and attempts to make the described receive buffers available to the adapter for future data transfers. Because it is desirable to keep as many buffers as possible available to the adapter, the SCSI device driver should pass



processed **tm\_buf** structures to the SCSI-adapter device driver's buffer-free routine as quickly as possible. The writer of a SCSI device driver should avoid requiring the last buffer of a **send** command to be received before processing buffers, as this could cause a situation where all buffers are in use and the **send** command has not completed. It is recommended that the writer therefore place a threshold of 25% on the free buffers. That is, when 25% or more of the number of buffers allocated for this device have been processed and the **send** command is not completed, the SCSI device driver should free the processed buffers by passing them to the SCSI adapter device driver's buffer-free routine.

---

## Required SCSI Adapter Device Driver ioctl Commands

Various ioctl operations must be performed for proper operation of the SCSI adapter device driver. The ioctl operations described here are the minimum set of commands the SCSI adapter device driver must implement to support SCSI device drivers. Other operations might be required in the SCSI adapter device driver to support, for example, system management facilities and diagnostics. SCSI device driver writers also need to understand these ioctl operations.

Every SCSI adapter device driver must support the **IOCINFO** ioctl operation. The structure to be returned to the caller is the **devinfo** structure, including the **scsi** union definition for the SCSI adapter, which can be found in the `/usr/include/sys/devinfo.h` file. The SCSI device driver should request the **IOCINFO** ioctl operation (probably during its open routine) to get the maximum transfer size of the adapter.

**Note:** The SCSI adapter device driver ioctl operations can only be called from the process level. They cannot be run from a call on any more favored priority levels. Attempting to call them from a more favored priority level can result in a system crash.

## Initiator-Mode ioctl Commands

The following **SCIOSTART** and **SCIOSTOP** operations must be sent by the SCSI device driver (for the open and close routines, respectively) for each device. They cause the SCSI adapter device driver to allocate and initialize internal resources. The **SCIOHALT** ioctl operation is used to abort pending or running commands, usually after signal processing by the SCSI device driver. This might be used by a SCSI device driver to end an operation instead of waiting for completion or a time out. The **SCIORESET** operation is provided for clearing device hard errors and competing initiator reservations during open processing by the SCSI device driver. The **SCIOGTHW** operation is supported by SCSI adapter device drivers that support gathered write commands to target devices.

Except where noted otherwise, the *arg* parameter for each of the ioctl operations described here must contain a long integer. In this field, the least significant byte is the SCSI LUN and the next least significant byte is the SCSI ID value. (The upper two bytes are reserved and should be set to 0.) This provides the information required to allocate or deallocate resources and perform SCSI bus operations for the ioctl operation requested.

The following information is provided on the various ioctl operations:

### SCIOSTART

This operation allocates and initializes SCSI device-dependent information local to the SCSI adapter device driver. Run this operation only on the first open of an ID/LUN device. Subsequent **SCIOSTART** commands to the same ID/LUN fail unless an intervening **SCIOSTOP** command is issued.

The following values for the **errno** global variable are supported:

**0** Indicates successful completion.

**EIO** Indicates lack of resources or other error-preventing device allocation.

### EINVAL

Indicates that the selected SCSI ID and LUN are already in use, or the SCSI ID matches the adapter ID.

## ETIMEDOUT

Indicates that the command did not complete.

## SCIOSTOP

This operation deallocates resources local to the SCSI adapter device driver for this SCSI device. This should be run on the last close of an ID/LUN device. If an **SCIOSTART** operation has not been previously issued, this command is unsuccessful.

The following values for the **errno** global variable should be supported:

**0** Indicates successful completion.

**EIO** Indicates error preventing device deallocation.

## EINVAL

Indicates that the selected SCSI ID and LUN have not been started.

## ETIMEDOUT

Indicates that the command did not complete.

## SCIOCMD

The SCIOCMD operation provides the means for issuing any SCSI command to the specified device after the SCSI device has been successfully started (SCIOSTART). The SCSI adapter driver performs no error recovery other than issuing a request sense for a SCSI check condition error. If the caller allocated an autosense buffer, then the request sense data is returned in that buffer. The SCSI adapter driver will not log any errors in the system error log for failures on a SCIOCMD operation. The following is a typical call:

```
rc = ioctl(adapter, SCIOCMD, &iocmd);
```

where *adapter* is a file descriptor and *iocmd* is an **sc\_passthru** structure as defined in the **/usr/include/sys/scsi.h** header file. The SCSI ID and LUN should be placed in the **sc\_passthru** parameter block.

The SCSI status byte and the adapter status bytes are returned through the **sc\_passthru** structure. If the SCIOCMD operation returns a value of -1 and the **errno** global variable is set to a nonzero value, the requested operation has failed. In this case, the caller should evaluate the returned status bytes to determine why the operation failed and what recovery actions should be taken.

If a SCIOCMD operation fails because a field in the **sc\_passthru** structure has an invalid value, then the subroutine will return a value of -1 and set the **errno** global variable to **EINVAL**. In addition the **einval\_arg** field will be set to the field number (starting with 1 for the version field) of the field that had an invalid value. A value of 0 for the **einval\_arg** field indicates no additional information on the failure is available.

The **devinfo** structure defines the maximum transfer size for the command. If an attempt is made to transfer more than the maximum, a value of -1 is returned and the **errno** global variable set to a value of **EINVAL**. Refer to the Small Computer System Interface (SCSI) Specification for the applicable device to get request sense information.

Possible **errno** values are:

**EIO** A system error has occurred. Consider retrying the operation several (three or more) times, because another attempt might be successful. If an **EIO** error occurs and the **status\_validity** field is set to **SC SCSI\_ERROR**, then the **scsi\_status** field has a valid value and should be inspected.

If the **status\_validity** field is zero and remains so on successive retries, then an unrecoverable error has occurred with the device.

If the **status\_validity** field is `SC_SCSI_ERROR` and the **scsi\_status** field contains a Check Condition status, then a SCSI request sense should be issued using the `SCIOCMD` ioctl to recover the the sense data.

#### **EFAULT**

A user process copy has failed.

#### **EINVAL**

The device is not opened or the caller has set a field in the **sc\_passthru** structure to an invalid value.

#### **EACCES**

The adapter is in diagnostics mode.

#### **ENOMEM**

A memory request has failed.

#### **ETIMEDOUT**

The command has timed out, which indicates the operation did not complete before the time-out value was exceeded. Consider retrying the operation.

#### **ENODEV**

The device is not responding.

**Note:** This operation requires the **SCIOSTART** operation to be run first.

If the FCP **SCIOCMD** ioctl operation completes successfully, then the **adap\_set\_flags** field might have the **SC\_RET\_ID** flag set. This field is set only if the **world\_wide\_name** and **node\_name** fields were provided in the ioctl call and the FC adapter driver detects that the **scsi\_id** field of this device has changed. The **scsi\_id** field will contain the new **scsi\_id** value.

The **version** field of the **scsi\_passthru** structure can be set to the value of `SC_VERSION_2` in `/usr/include/sys/scsi.h` or `SCSI_VERSION_2` in `/usr/include/sys/scsi_buf.h`, and the user can provide the following fields:

- **variable\_cdb\_ptr** - pointer to a buffer that contains the SCSI *cdb* variable.
- **variable\_cdb\_length** - the length of the variable *cdb* to which the *variable\_cdb\_ptr* points.

When the **SCIOCMD** ioctl request with the **version** field set to `SCSI_VERSION_2` completes and the device did not fully satisfy the request, the **residual** field indicates left over data. If the request completes successfully, the **residual** field indicates the device does not have all the requested data. If the request did not complete successfully, check the **status\_validity** to see whether a valid SCSI bus problem exists. If a valid SCSI bus problem exists, the **residual** field indicates the number of bytes by which the device failed to complete the request.

For more information, see *SCIOCMD SCSI Adapter Device Driver ioctl Operation in AIX Version 6.1 Technical Reference: Kernel and Subsystems Volume 2*.

#### **SCIOHALT**

This operation halts outstanding transactions to this ID/LUN device and causes the SCSI adapter device driver to stop accepting transactions for this device. This situation remains in effect until the SCSI device driver sends another transaction with the **SC\_RESUME** flag set (in the **sc\_buf.flags** field) for this ID/LUN combination. The **SCIOHALT** ioctl operation causes the SCSI adapter device driver to fail the command in progress, if any, as well as all queued commands for the device with a return value of **ENXIO** in the `sc_buf.bufstruct.b_error` field. If an **SCIOSTART** operation has not been previously issued, this command fails.

The following values for the **errno** global variable are supported:

- 0** Indicates successful completion.
- EIO** Indicates an unrecovered I/O error occurred.

**EINVAL**

Indicates that the selected SCSI ID and LUN have not been started.

**ETIMEDOUT**

Indicates that the command did not complete.

**SCIORESET**

This operation causes the SCSI adapter device driver to send a SCSI Bus Device Reset (BDR) message to the selected SCSI ID. For this operation, the SCSI device driver should set the LUN in the *arg* parameter to the LUN ID of a LUN on this SCSI ID, which has been successfully started using the **SCIOSTART** operation.

The SCSI device driver should use this command only when directed to do a *forced open*. This occurs in two possible situations: one, when it is desirable to force the device to drop a SCSI reservation; two, when the device needs to be reset to clear an error condition (for example, when running diagnostics on this device).

**Note:** In normal system operation, this command should not be issued, as it would force the device to drop a SCSI reservation another initiator (and, hence, another system) might have. If an **SCIOSTART** operation has not been previously issued, this command is unsuccessful.

The following values for the **errno** global variable are supported:

**0** Indicates successful completion.

**EIO** Indicates an unrecovered I/O error occurred.

**EINVAL**

Indicates that the selected SCSI ID and LUN have not been started.

**ETIMEDOUT**

Indicates that the command did not complete.

**SCIOGTHW**

This operation is only supported by SCSI adapter device drivers that support gathered write commands. The purpose of the operation is to indicate support for gathered writes to SCSI device drivers that intend to use this facility. If the SCSI adapter device driver does not support gathered write commands, it must fail the operation. The SCSI device driver should call this operation from its open routine for a particular device instance. If the operation is unsuccessful, the SCSI device driver should not attempt to run a gathered write command.

The *arg* parameter to the **SCIOGTHW** is set to null by the caller to indicate that no input parameter is passed:

The following values for the **errno** global variable are supported:

**0** Indicates successful completion and in particular that the adapter driver supports gathered writes.

**EINVAL**

Indicates that the SCSI adapter device driver does not support gathered writes.

**Target-Mode ioctl Commands**

The following **SCIOSTARTTGT** and **SCIOSTOPTGT** operations must be sent by the SCSI device driver (for the open and close routines, respectively) for each target-mode device instance. This causes the SCSI adapter device driver to allocate and initialize internal resources, and, if necessary, prepare the hardware for operation.

Target-mode support in the SCSI device driver and SCSI adapter device driver is optional. A failing return code from these commands, in the absence of any programming error, indicates target mode is not supported. If the SCSI device driver requires target mode, it must check the return code to verify the SCSI adapter device driver supports it.

Only a kernel process or device driver can call these ioctls. If attempted by a user process, the ioctl will fail, and the **errno** global variable will be set to **EPERM**.

The following information is provided on the various target-mode ioctl operations:

### **SCIOSTARTTGT**

This operation opens a logical path to a SCSI initiator device. It allocates and initializes SCSI device-dependent information local to the SCSI adapter device driver. This is run by the SCSI device driver in its open routine. Subsequent **SCIOSTARTTGT** commands to the same ID (LUN is always 0) are unsuccessful unless an intervening **SCIOSTOPTGT** is issued. This command also causes the SCSI adapter device driver to allocate system buffer areas to hold data received from the initiator, and makes the adapter ready to receive data from the selected initiator.

The *arg* parameter to the **SCIOSTARTTGT** should be set to the address of an **sc\_strt\_tgt** structure, which is defined in the **/usr/include/sys/scsi.h** file. The following parameters are supported:

*id*        The caller fills in the SCSI ID of the attached SCSI initiator.

*lun*        The caller sets the LUN to 0, as the initiator LUN is ignored for received data.

*buf\_size*

The caller specifies size in bytes to be used for each receive buffer allocated for this host target instance.

*num\_bufs*

The caller specifies how many buffers to allocate for this target instance.

*tm\_correlator*

The caller optionally places a value in this field to be passed back in each **tm\_buf** for this target instance.

*recv\_func*

The caller places in this field the address of a pinned routine the SCSI adapter device driver should call to pass **tm\_bufs** received for this target instance.

*free\_func*

This is an output parameter the SCSI adapter device driver fills with the address of a pinned routine that the SCSI device driver calls to pass **tm\_bufs** after they have been processed. The SCSI adapter device driver ignores the value passed as input.

**Note:** All reserved fields should be set to 0 by the caller.

The following values for the **errno** global variable are supported:

**0**        Indicates successful completion.

### **EINVAL**

An **SCIOSTARTTGT** command has already been issued to this SCSI ID.

The passed SCSI ID is the same as that of the adapter.

The LUN ID field is not set to zero.

The *buf\_size* is not valid. This is an adapter dependent value.

The *Num\_bufs* is not valid. This is an adapter dependent value.

The *recv\_func* value, which cannot be null, is not valid.

**EPERM**

Indicates the caller is not running in kernel mode, which is the only mode allowed to run this operation.

**ENOMEM**

Indicates that a memory allocation failure has occurred.

**EIO** Indicates an I/O error occurred, preventing the device driver from completing **SCIOSTARTTGT** processing.

**SCIOSTOPTGT**

This operation closes a logical path to a SCSI initiator device. It causes the SCSI adapter device driver to deallocate device dependent information areas allocated in response to a **SCIOSTARTTGT** operation. It also causes the SCSI adapter device driver to deallocate system buffer areas used to hold data received from the initiator, and to disable the host adapter's ability to receive data from the selected initiator.

The *arg* parameter to the **SCIOSTOPTGT** ioctl should be set to the address of an **sc\_stop\_tgt** structure, which is defined in the **/usr/include/sys/scsi.h** file. The caller fills in the *id* field with the SCSI ID of the SCSI initiator, and sets the *lun* field to 0 as the initiator LUN is ignored for received data. Reserved fields should be set to 0 by the caller.

The following values for the **errno** global variable should be supported:

**0** Indicates successful completion.

**EINVAL**

An **SCIOSTARTTGT** command has not been previously issued to this SCSI ID.

**EPERM**

Indicates the caller is not running in kernel mode, which is the only mode allowed to run this operation.

**Target- and Initiator-Mode ioctl Commands**

For either target or initiator mode, the SCSI device driver can issue an **SCIOEVENT** ioctl operation to register for receiving asynchronous event status from the SCSI adapter device driver for a particular device instance. This is an optional call for the SCSI device driver, and is optionally supported for the SCSI adapter device driver. A failing return code from this command, in the absence of any programming error, indicates it is not supported. If the SCSI device driver requires this function, it must check the return code to verify the SCSI adapter device driver supports it.

Only a kernel process or device driver can invoke these ioctls. If attempted by a user process, the ioctl will fail, and the **errno** global variable will be set to **EPERM**.

The event registration performed by this ioctl operation is allowed once per device session. Only the first **SCIOEVENT** ioctl operation is accepted after the device session is opened. Succeeding **SCIOEVENT** ioctl operations will fail, and the **errno** global variable will be set to **EINVAL**. The event registration is canceled automatically when the device session is closed.

The *arg* parameter to the **SCIOEVENT** ioctl operation should be set to the address of an **sc\_event\_struct** structure, which is defined in the **/usr/include/sys/scsi.h** file. The following parameters are supported:

<i>id</i>	The caller sets <i>id</i> to the SCSI ID of the attached SCSI target device for initiator-mode. For target-mode, the caller sets the <i>id</i> to the SCSI ID of the attached SCSI initiator device.
<i>lun</i>	The caller sets the <i>lun</i> field to the SCSI LUN of the attached SCSI target device for initiator-mode. For target-mode, the caller sets the <i>lun</i> field to 0.

<i>mode</i>	Identifies whether the initiator- or target-mode device is being registered. These values are possible:  <b>SC_IM_MODE</b> This is an initiator mode device.  <b>SC_TM_MODE</b> This is a target mode device.
<i>async_correlator</i>	The caller places a value in this optional field, which is saved by the SCSI adapter device driver and returned when an event occurs in this field in the <b>sc_event_info</b> structure. This structure is defined in the <b>/user/include/sys/scsi.h</b> file.
<i>async_func</i>	The caller fills in the address of a pinned routine that the SCSI adapter device driver calls whenever asynchronous event status is available. The SCSI adapter device driver passes a pointer to a <b>sc_event_info</b> structure to the caller's <b>async_func</b> routine.

**Note:** All reserved fields should be set to 0 by the caller.

The following values for the **errno** global variable are supported:

<b>0</b>	Indicates successful completion.
<b>EINVAL</b>	Either an <b>SCIOSTART</b> or <b>SCIOSTARTTGT</b> has not been issued to this device instance, or this device is already registered for async events.
<b>EPERM</b>	Indicates the caller is not running in kernel mode, which is the only mode allowed to run this operation.

---

## Related Information

Logical File System Kernel Services

## Technical References

The following reference articles can be found in *AIX Version 6.1 Technical Reference: Kernel and Subsystems Volume 2*:

- scdisk SCSI Device Driver
- scsidisk SCSI Device Driver
- SCSI Adapter Device Driver
- SCIOCMD SCSI Adapter Device Driver ioctl Operation
- SCIODIAG (Diagnostic) SCSI Adapter Device Driver ioctl Operation
- SCIODNLD (Download) SCSI Adapter Device Driver ioctl Operation
- SCIOEVENT (Event) SCSI Adapter Device Driver ioctl Operation
- SCIOGTHW (Gathered Write) SCSI Adapter Device Driver ioctl Operation
- SCIOHALT (HALT) SCSI Adapter Device Driver ioctl Operation
- SCIOINQU (Inquiry) SCSI Adapter Device Driver ioctl Operation
- SCIOREAD (Read) SCSI Adapter Device Driver ioctl Operation
- SCIORESET (Reset) SCSI Adapter Device Driver ioctl Operation
- SCIOSTART (Start SCSI) SCSI Adapter Device Driver ioctl Operation
- SCIOSTARTTGT (Start Target) SCSI Adapter Device Driver ioctl Operation
- SCIOSTOP (Stop Device) SCSI Adapter Device Driver ioctl Operation
- SCIOSTOPTGT (Stop Target) SCSI Adapter Device Driver ioctl Operation
- SCIOSTUNIT (Start Unit) SCSI Adapter Device Driver ioctl Operation
- SCIOTRAM (Diagnostic) SCSI Adapter Device Driver ioctl Operation
- SCIOTUR (Test Unit Ready) SCSI Adapter Device Driver ioctl Operation





---

## Chapter 13. SCSI Architectural Model Subsystem

This overview describes the interface between a SCSI Architectural Model (SAM) device driver and a SAM adapter device driver. *SAM* is a set of multiple physical transport types, all of which make use of the SCSI command set. You can use the following physical transport types in SAM:

- Fibre Channel Protocol for SCSI (FCP)
- iSCSI
- Serial Attached SCSI (SAS)

For information about the traditional parallel bus implementation of SCSI, see Small Computer System Interface Subsystem (Parallel SCSI).

The SAM subsystem is directed toward those wishing to design and write a SAM storage device driver that interfaces with an existing SAM adapter device driver. It is also meant for those wishing to design and write a SAM adapter device driver that interfaces with existing SAM storage device drivers.

---

### Programming SAM Device Drivers

The SAM subsystem has two parts:

- Device Driver
- Adapter Device Driver

The adapter device driver is designed to shield you from having to communicate directly with the system I/O hardware. This gives you the ability to successfully write a device driver without having a detailed knowledge of the system hardware. You can look at the subsystem as a two-tiered structure in which the adapter device driver is the bottom or supporting layer. As a programmer, you need only worry about the upper layer. This chapter only discusses writing a device driver, because the adapter device driver is already provided.

The adapter device driver, or lower layer, is responsible only for the communications to and from the bus, and any error logging and recovery. The upper layer is responsible for all of the device-specific commands. The device driver should handle all commands directed towards its specific device by building the necessary sequence of I/O requests to the adapter device driver in order to properly communicate with the device.

These I/O requests contain the commands that are needed by the device. One important aspect to note is that the device driver cannot access any of the adapter resources and should never try to pass the commands directly to the adapter, since it has absolutely no knowledge of the meaning of those commands.

### FCP, iSCSI, and Virtual SCSI Client Device Drivers

The role of the device driver is to pass information between the operating system and the adapter device driver by accepting I/O requests and passing these requests to the adapter device driver. The device driver should accept either character or block I/O requests, build the necessary commands, and then issue these commands to the device through the adapter device driver.

The device driver should also process the various required reservations and releases needed for the device. The device driver is notified through the **iodone** kernel service once the adapter has completed the processing of the command. The device driver should then notify its calling process that the request has completed processing through the **iodone** kernel service.

## FCP, iSCSI, and Virtual SCSI Client Adapter Device Driver

Unlike most other device drivers, the adapter device driver does *not* support the **read** and **write** subroutines. It only supports the **open**, **close**, **ioctl**, **config**, and **strategy** subroutines. Included with the **open** subroutine call is the **openx** subroutine that allows adapter diagnostics.

A device driver does not need to access the diagnostic commands. Commands received from the device driver through the **strategy** routine of the adapter are processed from a queue. Once the command has completed, the device driver is notified through the **iodone** kernel service.

## FCP, iSCSI, and Virtual SCSI Client Adapter and Device Interface

**Note:** Virtual SCSI is available only on IBM eServer i5 and IBM eServer p5 models.

The adapter device driver does not contain the **ddread** and **ddwrite** entry points, but does contain the **ddconfig**, **ddopen**, **ddclose**, **dddump**, and **ddioctl** entry points.

Therefore, the adapter device driver's entry in the kernel devsw table contains only those entries plus an additional **ddstrategy** entry point. This **ddstrategy** routine is the path that the device driver uses to pass commands to the device driver. Access to these entry points is possible through the following kernel services:

- **fp\_open**
- **fp\_close**
- **devdump**
- **fp\_ioctl**
- **devstrat**

The adapter is accessed by the device driver through the **/dev/fscsi#** special files, where **#** indicates ascending numbers 0, 1, 2, and so on. The adapter is designed so that multiple devices on the same adapter can be accessed at the same time.

The iSCSI adapter is accessed by the device driver through the **/dev/iscsin** special files, where **n** indicates ascending numbers 0, 1, 2, and so on. The adapter is designed so that multiple devices on the same adapter can be accessed at the same time.

The Virtual SCSI Client adapter is accessed by the device driver through the **/dev/vscsiX** special files, where **X** indicates ascending numbers 0, 1, 2, and so on. The adapter is designed such that multiple devices on the same adapter can be accessed at the same time.

For additional information on spanned and gathered write commands, see "Understanding the Execution of SAM Initiator I/O Requests" on page 288.

### scsi\_buf Structure

The I/O requests made from the device driver to the adapter device driver are completed through the use of the **scsi\_buf** structure, which is defined in the **/usr/include/sys/scsi\_buf.h** header file. This structure, which is similar to the **buf** structure in other drivers, is passed between the two subsystem drivers through the **strategy** routine. The following is a brief description of the fields contained in the **scsi\_buf** structure:

- Reserved fields should be set to a value of 0, except where noted.
- The **bufstruct** field contains a copy of the standard **buf** buffer structure that documents the I/O request. Included in this structure, for example, are the buffer address, byte count, and transfer direction. The **b\_work** field in the **buf** structure is reserved for use by the adapter device driver. The current definition of the **buf** structure is in the **/usr/include/sys/buf.h** include file.
- The **bp** field points to the original buffer structure received by the Device Driver from the caller, if any. This can be a chain of entries in the case of spanned transfers (commands that transfer data from or to

more than one system-memory buffer). A null pointer indicates a nonspanned transfer. The null value specifically tells the adapter device driver that all the information needed to perform the DMA data transfer is contained in the **bufstruct** fields of the **scsi\_buf** structure.

- The **scsi\_command** field, defined as a **scsi\_cmd structure**, contains, for example, the SCSI command length, SCSI command, and a flag variable:
  - The **scsi\_length** field is the number of bytes in the actual SCSI command. This is normally 6,10,12, or 16 (decimal).
  - The **FCP\_flags** field contains the following bit flags:

#### **SC\_NODISC**

Do not allow the target to disconnect during this command.

#### **SC\_ASYNC**

Do not allow the adapter to negotiate for synchronous transfer to the device.

During normal use, the **SC\_NODISC** bit should not be set. Setting this bit allows a device executing commands to monopolize the transport layer. Sometimes it is desirable for a particular device to maintain control of the transport layer once it has successfully arbitrated for it; for instance, when this is the only device on the transport layer or the only device that will be in use. For performance reasons, it might not be desirable to go through selections again to save transport layer overhead on each command.

Also during normal use, the **SC\_ASYNC** bit must not be set. It should be set only in cases where a previous command to the device ended in an unexpected transport free condition. This condition is noted as **SCSI\_TRANSPORT\_FAULT** in the **adapter\_status** field of the **scsi\_cmd** structure. Because other errors might also result in the **SCSI\_TRANSPORT\_FAULT** flag being set, the **SC\_ASYNC** bit should only be set on the last retry of the failed command.

- The **FCP\_flags** field is not used by the Virtual SCSI client driver.
- The **scsi\_cdb** structure contains the physical SCSI command block. The 6 to 16 bytes of a single SCSI command are stored in consecutive bytes, with the op code identified individually. The **scsi\_cdb** structure contains the following fields:
  1. The **scsi\_op\_code** field specifies the standard op code for this command.
  2. The **scsi\_bytes** field contains the remaining command-unique bytes of the command block. The actual number of bytes depends on the value in the **scsi\_op\_code** field.
- The **timeout\_value** field specifies the time-out limit (in seconds) to be used for completion of this command. A time-out value of 0 means no time-out is applied to this I/O request.
- The **status\_validity** field contains an output parameter that can have one of the following bit flags as a value:

#### **SC\_SCSI\_ERROR**

The **scsi\_status** field is valid.

#### **SC\_ADAPTER\_ERROR**

The **adapter\_status** field is valid.

- The **scsi\_status** field in the **scsi\_buf** structure is an output parameter that provides valid command completion status when its **status\_validity** bit is nonzero. The **scsi\_buf.bufstruct.b\_error** field should be set to EIO anytime the **scsi\_status** field is valid. Typical status values include:

#### **SC\_GOOD\_STATUS**

The target successfully completed the command.

#### **SC\_CHECK\_CONDITION**

The target is reporting an error, exception, or other conditions.

#### **SC\_BUSY\_STATUS**

The target is currently transporting and cannot accept a command now.

#### **SC\_RESERVATION\_CONFLICT**

The target is reserved by another initiator and cannot be accessed.

#### **SC\_COMMAND\_TERMINATED**

The target terminated this command after receiving a terminate I/O process message from the adapter.

#### **SC\_QUEUE\_FULL**

The target's command queue is full, so this command is returned.

#### **SC\_ACA\_ACTIVE**

The device has an ACA (auto contingent allegiance) condition that requires a Clear ACA to request to clear it.

- The **adapter\_status** field is an output parameter that is valid when its **status\_validity** bit is nonzero. The **scsi\_buf.bufstruct.b\_error** field should be set to EIO anytime the **adapter\_status** field is valid. This field contains generic adapter card status. It is intentionally general in coverage so that it can report error status from any typical adapter.

If an error is detected during execution of a command, and the error prevented the command from actually being sent to the transport layer by the adapter, then the error should be processed or recovered, or both, by the adapter device driver.

If it is recovered successfully by the adapter device driver, the error is logged, as appropriate, but is not reflected in the **adapter\_status** byte. If the error cannot be recovered by the adapter device driver, the appropriate **adapter\_status** bit is set and the **scsi\_buf** structure is returned to the device driver for further processing.

If an error is detected after the command was actually sent to the device, then it should be processed or recovered, or both, by the device driver.

For error logging, the adapter device driver logs transport layer and adapter-related conditions, and the device driver logs device-related errors. In the following description, a capital letter (A) after the error name indicates that the adapter device driver handles error logging. A capital letter (H) indicates that the device driver handles error logging.

Some of the following error conditions indicate a device failure. Others are transport layer or adapter-related.

#### **SCSI\_HOST\_IO\_BUS\_ERR (A)**

The system I/O transport layer generated or detected an error during a DMA or Programmed I/O (PIO) transfer.

#### **SCSI\_TRANSPORT\_FAULT (H)**

The transport protocol or hardware was unsuccessful.

#### **SCSI\_CMD\_TIMEOUT (H)**

The command timed out before completion.

#### **SCSI\_NO\_DEVICE\_RESPONSE (H)**

The target device did not respond to selection phase.

#### **SCSI\_ADAPTER\_HDW\_FAILURE (A)**

The adapter indicated an onboard hardware failure.

#### **SCSI\_ADAPTER\_SFW\_FAILURE (A)**

The adapter indicated microcode failure.

#### **SCSI\_FUSE\_OR\_TERMINAL\_PWR (A)**

The adapter indicated a blown terminator fuse or bad termination.

#### **SCSI\_TRANSPORT\_RESET (A)**

The adapter indicated the transport layer has been reset.

#### **SCSI\_WW\_NAME\_CHANGE (A)**

The adapter indicated the device at this SCSI ID has a new world wide name.

#### **SCSI\_TRANSPORT\_BUSY (A)**

The adapter indicated the transport layer is busy.

### SCSI\_TRANSPORT\_DEAD (A)

The adapter indicated the transport layer currently inoperative and is likely to remain this way for an extended time.

- The **add\_status** field contains additional device status. For devices, this field contains the Response code returned.
- When the device driver queues multiple transactions to a device, the **adap\_q\_status** field indicates whether or not the adapter driver has cleared its queue for this device after an error has occurred. The flag of **SC\_DID\_NOT\_CLEAR\_Q** indicates that the adapter driver has not cleared its queue for this device and that it is in a halted state (so none of the pending queued transactions are sent to the device).
- The **q\_tag\_msg** field indicates if the adapter can attempt to queue this transaction to the device. This information causes the adapter to fill in the Queue Tag Message Code of the queue tag message for a command. The following values are valid for this field:

#### SC\_NO\_Q

Specifies that the adapter does not send a queue tag message for this command, and so the device does not allow more than one command on its command queue. This value must be used for all commands sent to devices that do not support command tag queuing.

#### SC\_SIMPLE\_Q

Specifies placing this command in the device's command queue. The device determines the order that it executes commands in its queue. The SCSI-2 specification calls this value the Simple Queue Tag Message.

#### SC\_HEAD\_OF\_Q

Specifies placing this command first in the device's command queue. This command does not preempt an active command at the device, but it is executed before all other commands in the command queue. The SCSI-2 specification calls this value the Head of Queue Tag Message.

#### SC\_ORDERED\_Q

Specifies placing this command in the device's command queue. The device processes these commands in the order that they are received. The SCSI-2 specification calls this value the Ordered Queue Tag Message.

#### SC\_ACA\_Q

Specifies placing this command in the device's command queue, when the device has an ACA (auto contingent allegiance) condition. The SCSI-3 Architecture Model calls this value the ACA task attribute.

**Note:** Commands with the value of SC\_NO\_Q for the **q\_tag\_msg** field (except for request sense commands) should not be queued to a device whose queue contains a command with another value for **q\_tag\_msg**. If commands with the SC\_NO\_Q value (except for request sense) are sent to the device, then the device driver must make sure that no active commands are using different values for **q\_tag\_ms**. Similarly, the device driver must also make sure that a command with a **q\_tag\_msg** value of SC\_ORDERED\_Q, SC\_HEAD\_Q, or SC\_SIMPLE\_Q is not sent to a device that has a command with the **q\_tag\_msg** field of SC\_NO\_Q.

- The flags field contains bit flags sent from the device driver to the adapter device driver. The following flags are defined:

#### SC\_RESUME

When set, means the adapter device driver should resume transaction queuing for this ID/LUN. Error recovery is complete after a **SCIOHALT** operation, check condition, or severe transport error. This flag is used to restart the adapter device driver following a reported error.

#### SC\_DELAY\_CMD

When set, means the adapter device driver should delay sending this command (following a

reset or BDR to this device) by at least the number of seconds specified to the adapter device driver in its configuration information. For devices that do not require this function, this flag should not be set.

#### **SC\_Q\_CLR**

When set, means the adapter driver should clear its transaction queue for this ID/LUN. The transaction containing this flag setting does not require an actual command in the **scsi\_buf** because it is flushed back to the device driver with the rest of the transactions for this ID/LUN. However, this transaction must have the SCSI ID field (**scsi\_buf.scsi\_id**) and the LUN field (**scsi\_buf.lun\_id**) filled in with the device's SCSI ID and LUN. This flag is valid only during error recovery of a check condition or command terminated at a command tag queuing device when the **SC\_DID\_NOT\_CLR\_Q** flag is set in the **scsi\_buf.adap\_q\_status** field.

#### **SC\_Q\_RESUME**

When set, means that the adapter driver should resume its halted transaction queue for this ID/LUN. The transaction containing this flag setting does not require an actual command to be sent to the adapter driver. However, this transaction must have the SCSI ID field (**scsi\_buf.scsi\_id**) and the LUN field (**scsi\_buf.lun\_id**) filled in with the device's SCSI ID and logical unit number (LUN). If the transaction containing this flag setting is the first issued by the device driver after it receives an error (indicating that the adapter driver's queue is halted), then the **SC\_RESUME** flag must be set also.

#### **SC\_CLEAR\_ACA**

When set, means the SCSI adapter driver should issue a Clear ACA task management request for this ID/LUN. This flag should be used in conjunction with either the **SC\_Q\_CLR** or **SC\_Q\_RESUME** flags to clear or resume the SCSI adapter driver's queue for this device. If neither of these flags is used, then this transaction is treated as if the **SC\_Q\_RESUME** flag is also set. The transaction containing the **SC\_CLEAR\_ACA** flag setting does not require an actual SCSI command in the **sc\_buf**. If this transaction contains a SCSI command then it will be processed depending on whether **SC\_Q\_CLR** or **SC\_Q\_RESUME** is set. This transaction must have the SCSI ID field (**scsi\_buf.scsi\_id**) and the LUN field (**scsi\_buf.lun\_id**) filled in with the device's SCSI ID and LUN. This flag is valid only during error recovery of a check condition or command terminated at a command tag queuing.

#### **SC\_TARGET\_RESET**

When set, means the SCSI adapter driver should issue a Target Reset task management request for this ID/LUN. This flag should be used in conjunction with the **SC\_Q\_CLR** flag. The transaction containing this flag setting does allow an actual command to be sent to the adapter driver. However, this transaction must have the SCSI ID field (**scsi\_buf.scsi\_id**) filled in with the device's SCSI ID. If the transaction containing this flag setting is the first issued by the device driver after it receives an error (indicating that the adapter driver's queue is halted), then the **SC\_RESUME** flag must be set also.

#### **SC\_LUN\_RESET**

When set, means the SCSI adapter driver should issue a Lun Reset task management request for this ID/LUN. This flag should be used in conjunction with the **SC\_Q\_CLR** flag. The transaction containing this flag setting does allow an actual command to be sent to the adapter driver. However, this transaction must have the the SCSI ID field (**scsi\_buf.scsi\_id**) and the LUN field (**scsi\_buf.lun\_id**) filled in with the device's SCSI ID and logical unit number (LUN). If the transaction containing this flag setting is the first issued by the device driver after it receives an error (indicating that the adapter driver's queue is halted), then the **SC\_RESUME** flag must be set also.

- The **dev\_flags** field contains additional values sent from the FCP device driver to the FCP adapter device driver. This field is not used for iSCSI or Virtual SCSI device drivers. The following values are defined:

#### **FC\_CLASS1**

When set, this tells the SCSI adapter driver that it should issue this request as a Fibre Channel Class 1 request. If the SCSI adapter driver does not support this class, then it will fail the

**scsi\_buf** with an error of EINVAL. If no Fibre Channel Class is specified in the **scsi\_buf** then the SCSI adapter will default to a Fibre Channel Class.

#### FC\_CLASS2

When set, this tells the SCSI adapter driver that it should issue this request as a Fibre Channel Class 2 request. If the SCSI adapter driver does not support this class, then it will fail the **scsi\_buf** with an error of EINVAL. If no Fibre Channel Class is specified in the **scsi\_buf** then the SCSI adapter will default to a Fibre Channel Class.

#### FC\_CLASS3

When set, this tells the SCSI adapter driver that it should issue this request as a Fibre Channel Class 3 request. If the SCSI adapter driver does not support this class, then it will fail the **scsi\_buf** with an error of EINVAL. If no Fibre Channel Class is specified in the **scsi\_buf** then the SCSI adapter will default to a Fibre Channel Class.

#### FC\_CLASS4

When set, this tells the SCSI adapter driver that it should issue this request as a Fibre Channel Class 4 request. If the SCSI adapter driver does not support this class, then it will fail the **scsi\_buf** with an error of EINVAL. If no Fibre Channel Class is specified in the **scsi\_buf** then the SCSI adapter will default to a Fibre Channel Class.

- The **add\_work** field is reserved for use by the adapter device driver.
- The **adap\_set\_flags** field contains an output parameter that can have one of the following bit flags as a value:

#### SC\_AUTOSENSE\_DATA\_VALID

Autosense data was placed in the autosense buffer referenced by the **autosense\_buffer\_ptr** field.

- The **autosense\_length** field contains the length in bytes of the SCSI device driver's sense buffer, which is referenced via the **autosense\_buffer\_ptr** field. For devices this field must be non-zero, otherwise the autosense data will be lost.
- The **autosense\_buffer\_ptr** field contains the address of the SCSI devices driver's autosense buffer for this command. For devices this field must be non-NULL, otherwise the autosense data will be lost.
- The **dev\_burst\_len** field contains the burst size if this write operation in bytes. This should only be set by the device driver if it has negotiated with the device and it allows burst of write data without transfer ready's. For most operations, this should be set to 0.
- The **scsi\_id** field contains the 64-bit SCSI ID for this device. This field must be set for FCP devices.
- The **lun\_id** field contains the 64-bit lun ID for this device. This field must be set for devices.
- The **kernext\_handle** field contains the pointer returned from the **kernext\_handle** field of the **scsi\_sciolst** argument for the SCIOSTART ioctl.

## Adapter and Device Driver Intercommunication

In a typical request to the device driver, a call is first made to the device driver's **strategy** routine, which takes care of any necessary queuing. The device driver's **strategy** routine then calls the device driver's **start** routine, which fills in the **scsi\_buf** structure and calls the adapter driver's **strategy** routine through the **devstrat** kernel service.

The adapter driver's **strategy** routine validates all of the information contained in the **scsi\_buf** structure and also performs any necessary queuing of the transaction request. If no queuing is necessary, the adapter driver's **start** subroutine is called.

When an interrupt occurs, adapter driver **interrupt** routine fills in the **status\_validity** field and the appropriate **scsi\_status** or **adapter\_status** field of the **scsi\_buf** structure. The **bufstruct.b\_resid** field is also filled in with the value of nontransferred bytes. The adapter driver's **interrupt** routine then passes this newly filled in **scsi\_buf** structure to the **iodone** kernel service, which then signals the device driver's **iodone** subroutine. The adapter driver's **start** routine is also called from the **interrupt** routine to process any additional transactions on the queue.

The device driver's **iodone** routine should then process all of the applicable fields in the queued **scsi\_buf** structure for any errors and attempt error recovery if necessary. The device driver should then dequeue the **scsi\_buf** structure and then pass a pointer to the structure back to the **iodone** kernel service so that it can notify the originator of the request.

## FCP, iSCSI, and Virtual SCSI Client Adapter Device Driver Routines

This section describes the following routines:

- **config**
- **open**
- **close**
- **openx**
- **strategy**
- **ioctl**
- **start**
- **interrupt**

### config Routine

The **config** routine performs all of the processing needed to configure, unconfigure, and read Vital Product Data (VPD) for the adapter. When this routine is called to configure an adapter, it performs the required checks and building of data structures needed to prepare the adapter for the processing of requests.

When asked to unconfigure or terminate an adapter, this routine deallocates any structures defined for the adapter and marks the adapter as unconfigured. This routine can also be called to return the Vital Product Data for the adapter, which contains information that is used to identify the serial number, change level, or part number of the adapter.

### open Routine

The **open** routine establishes a connection between a special file and a file descriptor. This file descriptor is the link to the special file that is the access point to a device and is used by all subsequent calls to perform I/O requests to the device. Interrupts are enabled and any data structures needed by the adapter driver are also initialized.

### close Routine

The **close** routine marks the adapter as closed and disables all future interrupts, which causes the driver to reject all future requests to this adapter.

### openx Routine

The **openx** routine allows a process with the proper authority to open the adapter in diagnostic mode. If the adapter is already open in either normal or diagnostic mode, the **openx** subroutine has a return value of -1. Improper authority results in an **errno** value of EPERM, while an already open error results in an **errno** value of EACCES. If the adapter is in diagnostic mode, only the **close** and **ioctl** routines are allowed. All other routines return a value of -1 and an **errno** value of EACCES.

While in diagnostics mode, the adapter can run diagnostics, run wrap tests, and download microcode. The **openx** routine is called with an **ext** parameter that contains the adapter mode and the SC\_DIAGNOSTIC value, both of which are defined in the **sys/scsi.h** header file.

### strategy Routine

The **strategy** routine is the link between the device driver and the adapter device driver for all normal I/O requests. Whenever the device driver receives a call, it builds an **scsi\_buf** structure with the correct parameters and then passes it to this routine, which in turn queues up the request if necessary. Each request on the pending queue is then processed by building the necessary commands required to carry out the request. When the command has completed, the device driver is notified through the **iodone** kernel service.





- **SCIOLPAYLD**
- **SCIO LCHBA**
- **SCIO LPASSTHRUHBA**

## IOCINFO for FCP Adapters

This operation lets an FCP device driver obtain important information about a FCP adapter, including the adapter's SCSI ID, the maximum data transfer size in bytes, and the FC topology to which the adapter is connected. By knowing the maximum data transfer size, a FCP device driver can control several different devices on several different adapters. This operation returns a **devinfo** structure as defined in the **sys/devinfo.h** header file with the device type **DD\_BUS** and subtype **DS\_FCP**. The following is an example of a call to obtain the information:

```
rc = fp_ioctl(fp, IOCINFO, &infostruct, NULL);
```

where *fp* is a pointer to a file structure and *infostruct* is a **devinfo** structure. A non-zero rc value indicates an error. Note that the **devinfo** structure is a union of several structures and that **fcp** is the structure that applies to the adapter. For example, the maximum transfer size value is contained in the **infostruct.un.fcp.max\_transfer** variable and the card ID is contained in **infostruct.un.fcp.scsi\_id**.

## IOCINFO for iSCSI Adapters

This operation lets an iSCSI device driver obtain important information about an iSCSI adapter, including the adapter's maximum data transfer size in bytes. By knowing the maximum data transfer size, an iSCSI device driver can control several different devices on several different adapters. This operation returns a **devinfo** structure as defined in the **sys/devinfo.h** header file with the device type **DD\_BUS** and subtype **DS\_ISCSI**. The following is an example of a call to obtain the information:

```
rc = fp_ioctl(fp, IOCINFO, &infostruct, NULL);
```

where *fp* is a pointer to a file structure and *infostruct* is a **devinfo** structure. A non-zero rc value indicates an error. Note that the **devinfo** structure is a union of several structures and that **iscsi** is the structure that applies to the adapter. For example, the maximum transfer size value is contained in the **infostruct.un.iscsi.max\_transfer** variable.

## IOCINFO for SAS Adapters

This operation allows a SAS device driver to obtain important information about a SAS adapter, such as the maximum data-transfer size in bytes of the adapter. By knowing the maximum data transfer size, a SAS device driver can control several different devices on several different adapters. This operation returns a **devinfo** structure as defined in the **sys/devinfo.h** header file with the **DD\_BUS** device type and the **DS\_SAS** subtype. The following example shows a call to obtain the information:

```
rc = fp_ioctl(fp, IOCINFO, &infostruct, NULL);
```

where the *fp* parameter points to a file structure and the *infostruct* parameter is a **devinfo** structure. A nonzero rc value indicates an error. The **devinfo** structure is a union of several structures and the **sas** structure applies to the adapter. For example, the maximum transfer-size value is contained in the **infostruct.un.sas.max\_transfer** variable.

## IOCINFO for Virtual SCSI Adapters

This operation lets a Virtual SCSI device driver obtain important information about a Virtual SCSI adapter, including the adapter's maximum data transfer size in bytes. This information is determined by the Virtual SCSI server driver. By knowing the maximum data transfer size, a Virtual SCSI device driver can control several different devices on several different adapter instances. This operation returns a **devinfo** structure as defined in the **sys/devinfo.h** header file with the device type **DD\_BUS** and subtype **DS\_CVSCSI**. The following is an example of a call to obtain the information:

```
rc = fp_ioctl(fp, IOCINFO, &infostruct, NULL);
```

where *fp* is a pointer to a file structure and *infostruct* is a **devinfo** structure. A non-zero *rc* value indicates an error. The **devinfo** structure is a union of several structures and Virtual SCSI is the structure that applies to the adapter. For example, the maximum transfer size value is contained in the *infostruct.un.vscsi.max\_transfer* variable.

## SCIOLSTART

This operation opens a logical path to the target device and causes the SAM adapter device driver to allocate and initialize all of the data areas needed for the target device. The **SCIOLSTOP** operation should be issued when those data areas are no longer needed. The **SCIOLSTART** operation must be issued before any nondiagnostic operation except for **IOCINFO**. The following is a typical call:

```
rc = fp_ioctl(fp, SCIOLSTART, &sciolst);
```

where *fp* is a pointer to a file structure and *sciolst* is a **scsi\_sciolst** structure (defined in **/usr/include/sys/scsi\_buf.h**) that contains the SCSI and Logical Unit Number (LUN) ID values of the device to be started. In addition, the **scsi\_sciolst** structure can be used to specify an explicit login for this operation.

For FCP adapters, the **version** field of the **scsi\_sciolst** structure must be set to a minimum value of **SCSI\_VERSION\_1**, which is defined in the **/usr/include/sys/scsi\_buf.h** file. In addition, the following fields can be set:

- **world\_wide\_name** - The caller can set the **world\_wide\_name** field to the World Wide Name of the attached target device. If the **world\_wide\_name** field is set and the **version** field is set to **SCSI\_VERSION\_1**, the World Wide Name can be used to address the target instead of the **scsi\_id** field. If **Dynamic Tracking of FC devices** is enabled, the **world\_wide\_name** field must be set to ensure communication with the device because the **scsi\_id** field of a device can change after dynamic tracking events.
- **node\_name** - The caller can set the **node\_name** field to the Node Name of the attached target device. For AIX 5.2 through AIX 5.2.0.9, if the **world\_wide\_name** field and the **version** field are set to **SCSI\_VERSION\_1** but the **node\_name** field is not set, the **scsi\_id** field is used for device lookup instead of the **world\_wide\_name**.

If a World Wide Name or Node Name is provided and it does not match the World Wide Name or Node Name that was detected for the target, an error log is generated and the **SCIOLSTART** operation fails with an **errno** of **ENXIO**.

Upon successfully return from an **SCIOLSTART** operation, both the **world\_wide\_name** field and the **node\_name** field are set to the World Wide Name and Node Name of this device. These values are inspected to ensure that the **SCIOLSTART** operation was delivered to the intended device.

If **Dynamic Tracking of FC devices** is enabled, the **node\_name** field must be set to ensure communication with the device because the **scsi\_id** field of a device can change after dynamic tracking events.

For iSCSI adapters, this **version** field of the **scsi\_sciolst** must be set to a minimum value of **SCSI\_VERSION\_1** (defined in the **/usr/include/sys/scsi\_buf.h** file). In addition, iSCSI adapters require the caller to set the following fields:

- **lun\_id** of the device's LUN ID
- **parms.iscsi.name** to the device's iSCSI target name
- **parms.iscsi.iscsi\_ip\_addr** to the device's IP V4 or IP V6 address
- **parms.iscsi.port\_num** to the device's TCP port number

If the iSCSI **SCIOLSTART** ioctl operation completes successfully, then the **adap\_set\_flags** field should have the **SCIOL\_RET\_ID\_ALIAS** flag and the **scsi\_id** field set to a SCSI ID alias that can be used for subsequent ioctl calls to this device other than **SCIOLSTART**.

For Virtual SCSI adapters, the version field of the **scsi\_sciolst** structure must be set to the value of **SCSI\_VERSION\_1** (defined in the `/usr/include/sys/scsi_buf.h` file). In addition, Virtual SCSI adapters require the caller to set the **lun\_id** field to the Logical Unit Id (LUN) of the device being started.

For SAS adapters, the **version** field of the **scsi\_sciolst** structure must be set to a minimum value of **SCSI\_VERSION\_1** (defined in the `/usr/include/sys/scsi_buf.h` file). In addition, SAS adapter device drivers require the caller to set the following fields:

- **scsi\_id** (The field can also be referred to as the SAS address.)
- **lun\_id**

For AIX 5.2 with 5200-01 and later, if the FCP **SCIOSTART** ioctl operation completes successfully, and the **adap\_set\_flags** field has the **SCIOL\_DYNTRK\_ENABLED** flag set, then **Dynamic Tracking of FC Devices** has been enabled for this device.

All FC adapter ioctl calls for AIX 5.2 with 5200-01 and later, should set the **version** field to **SCSI\_VERSION\_1** if indicated in the ioctl structure comments in the header files. The **world\_wide\_name** and **node\_name** fields of all **SCSI\_VERSION\_1** ioctl structures should also be set. This is especially important if dynamic tracking has been enabled on this adapter. With dynamic tracking, the FC adapter driver can recover from **scsi\_id** changes of FC devices while devices are online. Because the **scsi\_id** can change, use of the **world\_wide\_name** and **node\_name** fields is necessary to ensure communication with the intended device.

Failure to use a **SCSI\_VERSION\_1** ioctl structure for **SCIOSTART** when dynamic tracking is enabled can produce undesired results, and temporarily disable dynamic tracking for a given device. If a target has at least one lun activated by **SCIOSTART** with the version field set to **SCSI\_VERSION\_1**, then a **SCSI\_VERSION\_0 SCIOSTART** fails. If this is the first lun activated by **SCIOSTART** on this target and the version field is set to **SCSI\_VERSION\_0**, then an error log of type **INFO** is generated and dynamic tracking is temporarily disabled for this target until a corresponding **SCSI\_VERSION\_0 SCIOSTOP** is issued.

The **version** field for all ioctl structures should be set consistently. For example, if an **SCIOSTART** operation is performed with the version field set to **SCSI\_VERSION\_1**, but the **SCIOLINQU** or **SCIOSTOP** ioctl operations have the **version** field set to **SCSI\_VERSION\_0**, then the ioctl call fails if dynamic tracking is enabled because the version fields do not match.

If the FCP **SCIOSTART** ioctl operation completes successfully, then the **adap\_set\_flags** field might have the **SCIOL\_RET\_ID\_ALIAS** flag set. This field is set only if the **world\_wide\_name** field was provided in the ioctl call and the FC adapter driver detects that the **scsi\_id** field of this device has changed. The **scsi\_id** field contains the new **scsi\_id** value.

If the caller of the SAM **SCIOSTART** is a kernel extension, the **SCIOL\_RET\_HANDLE** flag can be set in the **adap\_set\_flags** field along with the **kernext\_handle** field. In this case the **kernext\_handle** field can be used for **scsi\_buf** structures issued to the adapter driver for this device.

A nonzero return value indicates an error has occurred and all operations to this SCSI/LUN pair should cease because the device is either already started or failed the start operation. Possible **errno** values are:

EIO	The command could not complete due to a system error.
EINVAL	Either the Logical Unit Number (LUN) ID or SCSI ID is invalid, or the adapter is already open.
ENOMEM	Indicates that system resources are not available to start this device.
ETIMEDOUT	Indicates that the command did not complete.
ENODEV	Indicates that no device responded to the explicit process login at this SCSI ID.
ECONNREFUSED	Indicates that the device at this SCSI ID rejected explicit process login. This could be due to the device rejecting the security password or the device does not support FCP.
EACCES	The adapter is not in normal mode.

## SCIOLSTOP

This operation closes a logical path to the device and causes the adapter device driver to deallocate all data areas that were allocated by the **SCIOLSTART** operation. This operation should only be issued after a successful **SCIOLSTART** operation to a device. The following is a typical call:

```
rc = fp_ioctl(fp, SCIOLSTOP, &sciolst);
```

where *fp* is a pointer to a file structure and *sciolst* is a **scsi\_sciolst** structure (defined in **/usr/include/sys/scsi\_buf.h**) that contains the SCSI or iSCSI device's SCSI ID alias, and Logical Unit Number (LUN) ID values of the device to be stopped.

A nonzero return value indicates an error has occurred. Possible **errno** values are:

EIO	An unrecoverable system error has occurred.
EINVAL	The adapter was not in open mode.

The **version** field of the **scsi\_sciolst** structure must be set to a minimum value of **SCSI\_VERSION\_1**, which is defined in the **/usr/include/sys/scsi\_buf.h** file. In addition, the following fields can be set:

- **world\_wide\_name** - The caller can set the **world\_wide\_name** field to the World Wide Name of the attached target device. If **Dynamic Tracking of FC devices** is enabled, the **world\_wide\_name** field must be set to ensure communication with the device because the **scsi\_id** field of a device can change after dynamic tracking events.
- **node\_name** - The caller can set the **node\_name** field to the Node Name of the attached target device. For AIX 5.2 through AIX 5.2.0.9, if the **world\_wide\_name** field and the **version** field are set to **SCSI\_VERSION\_1** but the **node\_name** field is not set, the **scsi\_id** field is used for device lookup instead of the **world\_wide\_name**. If **Dynamic Tracking of FC devices** is enabled, the **node\_name** field must be set to ensure communication with the device because the **scsi\_id** field of a device can change after dynamic tracking events.

For Virtual SCSI adapters, the **version** field of the **scsi\_sciolst** structure must be set to the value of **SCSI\_VERSION\_1** (defined in the **/usr/include/sys/scsi\_buf.h** file). In addition, Virtual SCSI adapters require the caller to set the **lun\_id** field to the Logical Unit Id (LUN) of the device being stopped.

This operation requires **SCIOLSTART** to be run first.

## SCIOLEVENT

This operation lets a device driver register a particular device instance for receiving asynchronous event status by calling the **SCIOLEVENT** ioctl operation for the adapter device driver. When an event covered by the **SCIOLEVENT** ioctl operation is detected by the adapter device driver, it builds an **scsi\_event\_info** structure and passes a pointer to the structure and to the asynchronous event-handler routine entry point, which was previously registered.

The information reported in the **scsi\_event\_info.events** field does not queue to the device driver, but is instead reported as one or more flags as they occur. Because the data does not queue, the adapter device driver writer can use a single **scsi\_event\_info** structure and pass it one at a time, by pointer, to each asynchronous event handler routine for the appropriate device instance. After determining for which device the events are being reported, the device driver must copy the **scsi\_event\_info.events** field into local space and must not modify the contents of the rest of the **scsi\_event\_info** structure.

Because the event status is optional, the device driver writer determines what action is necessary to take upon receiving event status. The writer might decide to save the status and report it back to the calling application, or the device driver or application level program can take error recovery actions.

This operation should only be issued after a successful **SCIOLEVENT** operation to a device. The following is a typical call:

```
rc = fp_ioctl(fp, SCIOLEVENT, &scevent);
```

where *fp* is a pointer to a file structure and *scevent* is a **scsi\_event\_struct** structure (defined in **/usr/include/sys/scsi\_buf.h**) that contains the SCSI and Logical Unit Number (LUN) ID values of the device to be started.

A non-zero return value indicates an error has occurred. Possible **errno** values are:

EIO	An unrecoverable system error has occurred.
EINVAL	The adapter was not in open mode.

For FCP adapters, the **version** field of the **scsi\_event\_struct** structure must be set to a minimum value of **SCSI\_VERSION\_1**, which is defined in the **/usr/include/sys/scsi\_buf.h** file. In addition, the following fields can be set:

- **world\_wide\_name** - The caller can set the **world\_wide\_name** field to the World Wide Name of the attached target device. If **Dynamic Tracking of FC devices** is enabled, the **world\_wide\_name** field must be set to ensure communication with the device because the **scsi\_id** field of a device can change after dynamic tracking events.
- **node\_name** - The caller can set the **node\_name** field to the Node Name of the attached target device. If the **world\_wide\_name** field and the **version** field are set to **SCSI\_VERSION\_1** but the **node\_name** field is not set, the **scsi\_id** field is used for device lookup instead of the **world\_wide\_name**. If **Dynamic Tracking of FC devices** is enabled, the **node\_name** field must be set to ensure communication with the device because the **scsi\_id** field of a device can change after dynamic tracking events.

This operation requires **SCIOLEVENT** to be run first.

If the FCP **SCIOLEVENT** ioctl operation completes successfully, then the **adap\_set\_flags** field might have the **SC\_RET\_ID** flag set. This field is set only if the **world\_wide\_name** and **node\_name** fields were provided in the ioctl call and the FC adapter driver detects that the **scsi\_id** field of this device has changed. The **scsi\_id** field contains the new **scsi\_id** value.

## SCIOLINQU

This operation issues an inquiry command to an device and is used to aid in device configuration. The following is a typical call:

```
rc = ioctl(adapter, SCIOLINQU, &inquiry_block);
```

where *adapter* is a file descriptor and *inquiry\_block* is a **scsi\_inquiry** structure as defined in the **/usr/include/sys/scsi\_buf.h** header file. The FCP SCSI ID or device's SCSI ID alias, and LUN must be placed in the **scsi\_inquiry** parameter block. Possible **errno** values are:

EIO	A system error has occurred. Consider retrying the operation several times, because another attempt might be successful.
EFAULT	A user process copy has failed.
EINVAL	The device is not opened.
EACCES	The adapter is in diagnostics mode.
ENOMEM	A memory request has failed.
ETIMEDOUT	The command has timed out. Consider retrying the operation several times, because another attempt might be successful.
ENODEV	The device is not responding. Possibly no LUNs exist on the present FCP ID.
ENOCONNECT	A transport fault occurred.

For all physical transport types, the **version** field of the **scsi\_inquiry** structure must be set to a minimum value of `SCSI_VERSION_1`, which is defined in the `/usr/include/sys/scsi_buf.h` file. In addition, the following fields can be set:

- **world\_wide\_name** - The caller can set the **world\_wide\_name** field to the World Wide Name of the attached target device. If **Dynamic Tracking of FC devices** is enabled, the **world\_wide\_name** field must be set to ensure communication with the device because the **scsi\_id** field of a device can change after dynamic tracking events.
- **node\_name** - The caller can set the **node\_name** field to the Node Name of the attached target device. If the **world\_wide\_name** field and the **version** field are set to `SCSI_VERSION_1` but the **node\_name** field is not set, the **scsi\_id** field is used for device lookup instead of the **world\_wide\_name**. If **Dynamic Tracking of FC devices** is enabled, the **node\_name** field must be set to ensure communication with the device because the **scsi\_id** field of a device can change after dynamic tracking events.

When the **SCIOLINQU** ioctl request with the **version** field set to `SCSI_VERSION_2` completes and the device did not fully satisfy the request, the **residual** field indicates left over data. If the request completes successfully, the **residual** field indicates the device does not have all the requested data. If the request did not complete successfully, check the **status\_validity** to see whether a valid SCSI problem exists. If a valid SCSI problem exists, the **residual** field indicates the number of bytes by which the device failed to complete the request.

This operation requires **SCIOLSTART** to be run first.

If the FCP **SCIOLINQU** ioctl operation completes successfully, then the **adap\_set\_flags** field might have the **SC\_RET\_ID** flag set. This field is set only if the **world\_wide\_name** and **node\_name** fields were provided in the ioctl call and the FC adapter driver detects that the **scsi\_id** field of this device has changed. The **scsi\_id** field contains the new **scsi\_id** value.

## SCIOLSTUNIT

This operation issues a start unit command to an device and is used to aid in device configuration. The following is a typical call:

```
rc = ioctl(adapter, SCIOLSTUNIT, &start_block);
```

where *adapter* is a file descriptor and *start\_block* is a **scsi\_startunit** structure as defined in the `/usr/include/sys/scsi_buf.h` header file. The FCP ID or iSCSI device's SCSI ID alias, and LUN should be placed in the **scsi\_startunit** parameter block. The **start\_flag** field designates the start option, which when set to true, makes the device available for use. When this field is set to false, the device is stopped.

The **immed\_flag** field supports overlapping start operations to several devices on the adapter. When this field is set to false, status is returned only when the operation has completed. When this field is set to true, status is returned as soon as the device receives the command. The **SCIOLTUR** operation can then be issued to check on completion of the operation on a particular device.

**Attention:** When the SAM adapter issues simultaneous start operations, it is important that a sufficient delay is buffered between successive **SCIOLSTUNIT** operations to devices sharing a common power supply because damage to the system or devices can occur.

Possible **errno** values are:

EIO	A system error has occurred. Consider retrying the operation several times, because another attempt might be successful.
EFAULT	A user process copy has failed.
EINVAL	The device is not opened.
EACCES	The adapter is in diagnostics mode.
ENOMEM	A memory request has failed.

ETIMEDOUT	The command has timed out. Consider retrying the operation several times, because another attempt might be successful.
ENODEV	The device is not responding. Possibly no LUNs exist on the present FCP ID.
ENOCONNECT	A transport fault occurred.

For all physical transport types, the **version** field of the **scsi\_startunit** structure must be set to a minimum value of **SCSI\_VERSION\_1**, which is defined in the **/usr/include/sys/scsi\_buf.h** file. In addition, the following fields can be set:

- **world\_wide\_name** - The caller can set the **world\_wide\_name** field to the World Wide Name of the attached target device. If **Dynamic Tracking of FC devices** is enabled, the **world\_wide\_name** field must be set to ensure communication with the device because the **scsi\_id** field of a device can change after dynamic tracking events.
- **node\_name** - The caller can set the **node\_name** field to the Node Name of the attached target device. If the **world\_wide\_name** field and the **version** field are set to **SCSI\_VERSION\_1** but the **node\_name** field is not set, the **scsi\_id** field is used for device lookup instead of the **world\_wide\_name**. If **Dynamic Tracking of FC devices** is enabled, the **node\_name** field must be set to ensure communication with the device because the **scsi\_id** field of a device can change after dynamic tracking events.

This operation requires **SCIOSTART** to be run first.

If the FCP **SCIOSTUNIT** ioctl operation completes successfully, then the **adap\_set\_flags** field might have the **SC\_RET\_ID** flag set. This field is set only if the **world\_wide\_name** and **node\_name** fields were provided in the ioctl call and the FC adapter driver detects that the **scsi\_id** field of this device has changed. The **scsi\_id** field contains the new **scsi\_id** value.

## SCIOLTUR

This operation issues a Test Unit Ready command to an adapter and aids in device configuration. The following is a typical call:

```
rc = ioctl(adapter, SCIOLTUR, &ready_struct);
```

where *adapter* is a file descriptor and *ready\_struct* is a **scsi\_ready** structure as defined in the **/usr/include/sys/scsi\_buf.h** header file. The FCP SCSI ID or iSCSI device's SCSI ID alias, and LUN should be placed in the **scsi\_ready** parameter block. The status of the device can be determined by evaluating the two output fields: **status\_validity** and **scsi\_status**. Possible **errno** values are:

EIO	A system error has occurred. Consider retrying the operation several (around three) times, because another attempt might be successful. If an EIO error occurs and the <b>status_validity</b> field is set to <b>SC_FCP_ERROR</b> , then the <b>scsi_status</b> field has a valid value and should be inspected.  If the <b>status_validity</b> field is zero and remains so on successive retries, then an unrecoverable error has occurred with the device.  If the <b>status_validity</b> field is <b>SC_FCP_ERROR</b> and the <b>scsi_status</b> field contains a Check Condition status, then the <b>SCIOLTUR</b> operation should be retried after several seconds.  If after successive retries, the Check Condition status remains, the device should be considered inoperable.
EFAULT	A user process copy has failed.
EINVAL	The device is not opened.
EACCES	The adapter is in diagnostics mode.
ENOMEM	A memory request has failed.
ETIMEDOUT	The command has timed out. Consider retrying the operation several times, because another attempt might be successful.
ENODEV	The device is not responding and possibly no LUNs exist on the present target.
ENOCONNECT	A transport fault occurred.



For FCP adapters, the **version** field of the **scsi\_ready** structure must be set to a minimum value of `SCSI_VERSION_1`, which is defined in the `/usr/include/sys/scsi_buf.h` file. In addition, the following fields can be set:

- **world\_wide\_name** - The caller can set the **world\_wide\_name** field to the World Wide Name of the attached target device. If **Dynamic Tracking of FC devices** is enabled, the **world\_wide\_name** field must be set to ensure communication with the device because the **scsi\_id** field of a device can change after dynamic tracking events.
- **node\_name** - The caller can set the **node\_name** field to the Node Name of the attached target device. If the **world\_wide\_name** field and the **version** field are set to `SCSI_VERSION_1` but the **node\_name** field is not set, the **scsi\_id** field is used for device lookup instead of the **world\_wide\_name**. If **Dynamic Tracking of FC devices** is enabled, the **node\_name** field must be set to ensure communication with the device because the **scsi\_id** field of a device can change after dynamic tracking events.

For Virtual SCSI adapters, the **version** field of the **scsi\_sciolst** structure must be set to the value of `SCSI_VERSION_1` (defined in the `/usr/include/sys/scsi_buf.h` file). In addition, Virtual SCSI adapters require the caller to set the **lun\_id** field to the Logical Unit Id (LUN) of the device being reset. Target Reset is not supported on `vscsi` devices, so the driver sends a Lun Reset regardless of the value of the `SCIOLRESET_LUN_RESET` flag.

This operation requires **SCIOLSTART** to be run first.

If the FCP **SCIOLTUR** ioctl operation completes successfully, then the **adap\_set\_flags** field might have the **SC\_RET\_ID** flag set. This field is set only if the **world\_wide\_name** and **node\_name** fields were provided in the ioctl call and the FC adapter driver detects that the **scsi\_id** field of this device has changed. The **scsi\_id** field contains the new **scsi\_id** value.

## SCIOLREAD

This operation issues a read command to an device and is used to aid in device configuration. The following is a typical call:

```
rc = ioctl(adapter, SCIOLREAD, &readblk);
```

where *adapter* is a file descriptor and *readblk* is a **scsi\_readblk** structure as defined in the `/usr/include/sys/scsi_buf.h` header file. The FCP SCSI ID or iSCSI device's SCSI ID alias, and the LUN should be placed in the **scsi\_readblk** parameter block. Possible **errno** values are:

EIO	A system error has occurred. Consider retrying the operation several times, because another attempt might be successful.
EFAULT	A user process copy has failed.
EINVAL	The device is not opened.
EACCES	The adapter is in diagnostics mode.
ENOMEM	A memory request has failed.
ETIMEDOUT	The command has timed out. Consider retrying the operation several times, because another attempt might be successful.
ENODEV	The device is not responding. Possibly no LUNs exist on the present target.
ENOCONNECT	A transport fault occurred.

For all transport types, the **version** field of the **scsi\_readblk** structure must be set to a minimum value of `SCSI_VERSION_1`, which is defined in the `/usr/include/sys/scsi_buf.h` file. In addition, the following fields can be set:

- **world\_wide\_name** - The caller can set the **world\_wide\_name** field to the World Wide Name of the attached target device. If **Dynamic Tracking of FC devices** is enabled, the **world\_wide\_name** field must be set to ensure communication with the device because the **scsi\_id** field of a device can change after dynamic tracking events.

- **node\_name** - The caller can set the **node\_name** field to the Node Name of the attached target device. If the **world\_wide\_name** field and the **version** field are set to `SCSI_VERSION_1` but the **node\_name** field is not set, the **scsi\_id** field is used for device lookup instead of the **world\_wide\_name**. If **Dynamic Tracking of FC devices** is enabled, the **node\_name** field must be set to ensure communication with the device because the **scsi\_id** field of a device can change after dynamic tracking events.

When the **SCIOLREAD** ioctl request with the **version** field set to `SCSI_VERSION_2` completes and the device did not fully satisfy the request, the **residual** field indicates left over data. If the request completes successfully, the **residual** field indicates the device does not have all the requested data. If the request did not complete successfully, check the **status\_validity** to see whether a valid SCSI bus problem exists. If a valid SCSI bus problem exists, the **residual** field indicates the number of bytes by which the device failed to complete the request.

This operation requires **SCIOLSTART** to be run first.

If the FCP **SCIOLREAD** ioctl operation completes successfully, then the **adap\_set\_flags** field might have the **SC\_RET\_ID** flag set. This field is set only if the **world\_wide\_name** and **node\_name** fields were provided in the ioctl call and the FC adapter driver detects that the **scsi\_id** field of this device has changed. The **scsi\_id** field contains the new **scsi\_id** value.

## SCIOLRESET

If the **SCIOLRESET\_LUN\_RESET** flag is not set in the flags field of the **scsi\_sciolst**, then this operation causes an FCP or iSCSI device, clear all current commands, and return to an initial state by issuing a Target Reset, which resets all LUNs associated with the specified FCP ID or iSCSI device's SCSI ID alias. If, for an FCP or iSCSI device the **SCIOLRESET\_LUN\_RESET** flag is set in the flags field of the **scsi\_sciolst**, then this operation causes an FCP or iSCSI device to clear all current commands, and return to an initial state by issuing a Lun Reset, which resets just the specified LUN associated with the specified FCP ID or iSCSI device's SCSI ID alias. For the SAS transport type, the Target Reset command is not defined; and thus for SAS, this ioctl always results in Lun Reset being issued without regard to the **SCIOLRESET\_LUN\_RESET** flag setting.

**Note:** Both the Target and LUN reset commands that are generated by the **SCIOLRESET** ioctl causes a device that honors a standard SCSI reservation to release that reservation. Whether a device driver for the device reestablishes the standard SCSI reservation is based on the following conditions:

- The configuration attributes of the device
- The condition when the reservation is initially established
- The mode in which the device is opened

The following is a typical call:

```
rc = fp_ioctl(fp, SCIOLRESET, &sciolst);
```

where *fp* is a pointer to a file structure and *sciolst* is a **scsi\_sciolst** structure (defined in **/usr/include/sys/scsi\_buf.h**) that contains the SCSI ID or iSCSI device's SCSI ID alias, and Logical Unit Number (LUN) ID values of the device to be started.

A nonzero return value indicates an error has occurred. Possible **errno** values are:

EIO	An unrecoverable system error has occurred.
EINVAL	The device is not opened.
EACCES	The adapter is in diagnostics mode.
ETIMEDOUT	The operation did not complete before the time-out value was exceeded.

For all transport types, the **version** field of the **scsi\_sciolst** structure must be set to a minimum value of `SCSI_VERSION_1`, which is defined in the **/usr/include/sys/scsi\_buf.h** file. In addition, the following fields can be set:

- **world\_wide\_name** - The caller can set the **world\_wide\_name** field to the World Wide Name of the attached target device. If **Dynamic Tracking of FC devices** is enabled, the **world\_wide\_name** field must be set to ensure communication with the device because the **scsi\_id** field of a device can change after dynamic tracking events.
- **node\_name** - The caller can set the **node\_name** field to the Node Name of the attached target device. If the **world\_wide\_name** field and the **version** field are set to `SCSI_VERSION_1` but the **node\_name** field is not set, the **scsi\_id** field is used for device lookup instead of the **world\_wide\_name**. If **Dynamic Tracking of FC devices** is enabled, the **node\_name** field must be set to ensure communication with the device because the **scsi\_id** field of a device can change after dynamic tracking events.

This operation requires **SCIOLSTART** to be run first.

If the FCP **SCIOLRESET** ioctl operation completes successfully, then the **adap\_set\_flags** field might have the **SCIOL\_RET\_ID\_ALIAS** flag set. This field is set only if the **world\_wide\_name** and **node\_name** fields were provided in the ioctl call and the FC adapter driver detects that the **scsi\_id** field of this device has changed. The **scsi\_id** field contains the new **scsi\_id** value.

## SCIOLHALT

This operation stops the current command of the selected device, clears the command queue of any pending commands, and brings the device to a halted state. The adapter sends a SCSI abort message to the device and is usually used by the device driver to abort the current operation instead of waiting for it to complete or time out.

The adapter also performs normal error recovery procedures during this command. The following is a typical call:

```
rc = fp_ioctl(fp, SCIOLHALT, &sciolst);
```

where *fp* is a pointer to a file structure and *sciolst* is a **scsi\_sciolst** structure (defined in `/usr/include/sys/scsi_buf.h`) that contains the SCSI ID or iSCSI device's SCSI ID alias, and Logical Unit Number (LUN) ID values of the device to be started.

A nonzero return value indicates an error has occurred. Possible **errno** values are:

EIO	An unrecoverable system error has occurred.
EINVAL	The device is not opened.
EACCES	The adapter is in diagnostics mode.
ETIMEDOUT	The operation did not complete before the time-out value was exceeded.

For all transport types, the **version** field of the **scsi\_sciolst** structure must be set to a minimum value of `SCSI_VERSION_1`, which is defined in the `/usr/include/sys/scsi_buf.h` file. In addition, the following fields can be set:

- **world\_wide\_name** - The caller can set the **world\_wide\_name** field to the World Wide Name of the attached target device. If **Dynamic Tracking of FC devices** is enabled, the **world\_wide\_name** field must be set to ensure communication with the device because the **scsi\_id** field of a device can change after dynamic tracking events.
- **node\_name** - The caller can set the **node\_name** field to the Node Name of the attached target device. If the **world\_wide\_name** field and the **version** field are set to `SCSI_VERSION_1` but the **node\_name** field is not set, the **scsi\_id** field is used for device lookup instead of the **world\_wide\_name**. If **Dynamic Tracking of FC devices** is enabled, the **node\_name** field must be set to ensure communication with the device because the **scsi\_id** field of a device can change after dynamic tracking events.

This operation requires **SCIOLSTART** to be run first.

If the FCP **SCIOLHALT** ioctl operation completes successfully, then the **adap\_set\_flags** field might have the **SCIOL\_RET\_ID\_ALIAS** flag set. This field is set only if the **world\_wide\_name** and **node\_name** fields

were provided in the `ioctl` call and the FC adapter driver detects that the `scsi_id` field of this device has changed. The `scsi_id` field contains the new `scsi_id` value.

## SCIOLCMD

After the SCSI device has been successfully started using `SCIOLSTART`, the `SCIOLCMD` `ioctl` operation provides the means for issuing any SCSI command to the specified device. The SAM adapter driver performs no error recovery or logging on failures of this `ioctl` operation. The following is a typical call:

```
rc = ioctl(adapter, SCIOLCMD, &iocmd);
```

where `adapter` is a file descriptor and `iocmd` is a `scsi_iocmd` structure as defined in the `/usr/include/sys/scsi_buf.h` header file. The SCSI ID or iSCSI device's SCSI ID alias, and LUN ID should be placed in the `scsi_iocmd` parameter block.

The SCSI status byte and the adapter status bytes are returned using the `scsi_iocmd` structure. If the `SCIOLCMD` operation returns a value of -1 and the `errno` global variable is set to a nonzero value, the requested operation has failed. In this case, the caller should evaluate the returned status bytes to determine why the operation failed and what recovery actions should be taken.

The `devinfo` structure defines the maximum transfer size for the command. If an attempt is made to transfer more than the maximum, a value of -1 is returned and the `errno` global variable set to a value of `EINVAL`. Refer to the *Small Computer System Interface (SCSI) Specification* for the applicable device to get request sense information.

Possible `errno` values are:

EIO	A system error has occurred. Consider retrying the operation several (around three) times, because another attempt might be successful. If an EIO error occurs and the <code>status_validity</code> field is set to <code>SC_SCSI_ERROR</code> , then the <code>scsi_status</code> field has a valid value and should be inspected.  If the <code>status_validity</code> field is zero and remains so on successive retries then an unrecoverable error has occurred with the device.  If the <code>status_validity</code> field is <code>SC_SCSI_ERROR</code> and the <code>scsi_status</code> field contains a Check Condition status, then a SCSI request sense should be issued using the <code>SCIOLCMD</code> <code>ioctl</code> to recover the the sense data.
EFAULT	A user process copy has failed.
EINVAL	The device is not opened.
EACCES	The adapter is in diagnostics mode.
ENOMEM	A memory request has failed.
ETIMEDOUT	The command has timed out. Consider retrying the operation several times, because another attempt might be successful.
ENODEV	The device is not responding.
ETIMEDOUT	The operation did not complete before the time-out value was exceeded.

For all transport types, the `version` field of the `scsi_iocmd` structure must be set to a minimum value of `SCSI_VERSION_1`, which is defined in the `/usr/include/sys/scsi_buf.h` file. In addition, the following fields can be set:

- **world\_wide\_name** - The caller can set the `world_wide_name` field to the World Wide Name of the attached target device. If **Dynamic Tracking of FC devices** is enabled, the `world_wide_name` field must be set to ensure communication with the device because the `scsi_id` field of a device can change after dynamic tracking events.
- **node\_name** - The caller can set the `node_name` field to the Node Name of the attached target device. If the `world_wide_name` field and the `version` field are set to `SCSI_VERSION_1` but the `node_name` field is not set, the `scsi_id` field is used for device lookup instead of the `world_wide_name`. If **Dynamic Tracking of FC devices** is enabled, the `node_name` field must be set to ensure communication with the device because the `scsi_id` field of a device can change after dynamic tracking events.

The **version** field of the **scsi\_iocmd** structure can be set to the value of `SCSI_VERSION_2`, and the user can provide the following fields:

- **variable\_cdb\_ptr** - pointer to a buffer that contains the SCSI variable *cdb*.
- **variable\_cdb\_length** - the length of the *cdb* variable to which the *variable\_cdb\_ptr* points.

**Note:** The SAS transport type does not support the *cdb*s variable length.

When the **SCIOLCMD** ioctl request with the **version** field set to `SCSI_VERSION_2` completes and the device did not fully satisfy the request, the **residual** field indicates left over data. If the request completes successfully, the **residual** field indicates the device does not have all the requested data. If the request did not complete successfully, check the **status\_validity** to see whether a valid SCSI problem exists. If a valid SCSI problem exists, the **residual** field indicates the number of bytes by which the device failed to complete the request.

This operation requires **SCIOLSTART** to be run first.

If the FCP **SCIOLCMD** ioctl operation completes successfully, then the **adap\_set\_flags** field might have the **SC\_RET\_ID** flag set. This field is set only if the **world\_wide\_name** and **node\_name** fields were provided in the ioctl call and the FC adapter driver detects that the **scsi\_id** field of this device has changed. The **scsi\_id** field contains the new **scsi\_id** value.

## SCIOLNMSRV

**Note:** **SCIOLNMSRV** is specific to the iSCSI and FCP transport types.

This operation returns a list of target devices and is used to aid in SCSI device configuration. The following is a typical call:

```
rc = ioctl(adapter, SCIOLNMSRV, &nmserv);
```

where *adapter* is a file descriptor and *nmserv* is a **scsi\_nmserv** structure as defined in the `/usr/include/sys/scsi_buf.h` header file. The caller of this ioctl, must allocate a buffer be referenced by the **scsi\_id\_list** field. In addition the caller must set the **list\_len** field to indicate the size of the buffer in bytes.

On successful completion, the **num\_ids** field indicates the number of SCSI IDs returned in the current list. If the more ids were found then could be placed in the list, then the adapter driver updates the **list\_len** field to indicate the length of buffer needed to receive all SCSI IDs.

Possible **errno** values are:

EIO	A system error has occurred. Consider retrying the operation several times, because another attempt might be successful.
EFAULT	A user process copy has failed.
EINVAL	The physical configuration does not support this request.
ENOMEM	A memory request has failed.
ETIMEDOUT	The command has timed out. Consider retrying the operation several times, because another attempt might be successful.
ENODEV	The device is not responding. Possibly no LUNs exist on the present target.

## SCIOLQWWN

**Note:** **SCIOLQWWN** is specific to the FCP transport type.

This operation issues a request to find the SCSI ID of a device for the specified world wide name. The following is a typical call:

```
rc = ioctl(adapter, SCIOIQWVN, &qrywwn);
```

where *adapter* is a file descriptor and *qrywwn* is a **scsi\_qry\_wwn** structure as defined in the **/usr/include/sys/scsi\_buf.h** header file. The caller of this `ioctl`, must specify the device's world wide name in the **world\_wide\_name** field. On successful completion, the **scsi\_id** field is returned with the SCSI ID of the device with this world wide name.

Possible **errno** values are:

EIO	A system error has occurred. Consider retrying the operation several times, because another attempt might be successful.
EFAULT	A user process copy has failed.
EINVAL	The physical configuration does not support this request.
ENOMEM	A memory request has failed.
ETIMEDOUT	The command has timed out. Consider retrying the operation several times, because another attempt might be successful.
ENODEV	The device is not responding. Possibly no LUNs exist on the present FCP ID.

## SCIOLPAYLD

**Note:** The **SCIOLPAYLD** operation is specific to the iSCSI and FCP transport types.

The **SCIOLPAYLD** `ioctl` is not supported by the Virtual SCSI adapter driver.

This operation provides the means for issuing a transport payload to the specified device. The SAM adapter driver performs no error recovery or logging on failures of this `ioctl` operation. The following is a typical call:

```
rc = ioctl(adapter, SCIOLPAYLD, &payload);
```

where *adapter* is a file descriptor and *payload* is a **scsi\_trans\_payld** structure as defined in the **/usr/include/sys/scsi\_buf.h** header file. The SCSI ID or SCSI ID alias should be placed in the **scsi\_trans\_payld**. In addition the user must allocate a payload buffer referenced by the **payload\_buffer** field and a response buffer referenced by the **response\_buffer** field. The fields **payload\_size** and **response\_size** specify the size in bytes of the payload buffer and response buffer, respectively. In addition the caller can also set **payload\_type** (for FC this is the FC-4 type), and **payload\_ctl** (for FC this is the router control field),.

If the **SCIOLPAYLD** operation returns a value of -1 and the **errno** global variable is set to a nonzero value, the requested operation has failed. In this case, the caller should evaluate the returned status bytes to determine why the operation failed and what recovery actions should be taken.

Possible **errno** values are:

EIO	A system error has occurred.
EFAULT	A user process copy has failed.
EINVAL	Payload and or response buffer are too large. For FCP, iSCSI, and Virtual SCSI Client the maximum size is 4096 bytes.
ENOMEM	A memory request has failed.
ETIMEDOUT	The command has timed out. Consider retrying the operation several times, because another attempt might be successful.
ENODEV	The device is not responding.
ETIMEDOUT	The operation did not complete before the time-out value was exceeded.

## SCIO LCHBA

**Note:** The **SCIO LCHBA** operation is specific to the iSCSI and FCP transport types.

The **SCIO LCHBA** ioctl is not supported by the Virtual SCSI adapter driver.

When the device has been successfully opened, the **SCIO LCHBA** operation provides the means for issuing one or more common host bus adapter (HBA) API commands to the adapter. The SAM adapter driver performs full error recovery on failures of this operation.

The **arg** parameter contains the address of a **scsi\_chba** structure, which is defined in the **/usr/include/sys/scsi\_buf.h** file.

The **cmd** field in the **scsi\_chba** structure determines the common HBA API operation that is performed.

If the **SCIO LCHBA** operation fails, the subroutine returns a value of -1 and sets the **errno** global variable to a nonzero value. In this case, the caller should evaluate the returned status bytes to determine why the operation was unsuccessful and what recovery actions should be taken.

If a **SCIO LCHBA** operation fails because a field in the **scsi\_chba** structure has an invalid value, the subroutine returns a value of -1 and set the **errno** global variable to **EINVAL**.

## SCIO LPASSTHRUHBA

**Note:** The **SCIO LPASSTHRUHBA** operation is specific to the iSCSI and FCP transport types.

The **SCIO LPASSTHRUHBA** ioctl is not supported by the Virtual SCSI adapter driver.

When the device has been successfully opened, the **SCIO LPASSTHRUHBA** operation provides the means for issuing **passthru** commands to the adapter. The SAM adapter driver performs full error recovery on failures of this operation.

The **arg** parameter contains the address of a **scsi\_passthru\_hba** structure, which is defined in the **/usr/include/sys/scsi\_buf.h** file.

The **cmd** field in the **scsi\_passthru\_hba** structure determines the type of **passthru** operation to be performed.

If the **SCIO LPASSTHRUHBA** operation fails, the subroutine returns a value of -1 and sets the **errno** global variable to a nonzero value. In this case, the caller should evaluate the returned status bytes to determine why the operation was unsuccessful and what recovery actions should be taken.

If a **SCIO LPASSTHRUHBA** operation fails because a field in the **scsi\_passthru\_hba** structure has an invalid value, the subroutine returns a value of -1 and set the **errno** global variable to **EINVAL**.

---

## SAM Subsystem Overview

This section frequently refers to both *device driver* and *adapter device driver*. These two distinct device drivers work together in a layered approach to support attachment of a range of devices. The adapter device driver is the *lower* device driver of the pair, and the device driver is the *upper* device driver.

### Responsibilities of the Adapter Device Driver

The adapter device driver is the software interface to the system hardware. This hardware includes the transport layer hardware, plus any other system I/O hardware required to run an I/O request. The adapter device driver hides the details of the I/O hardware from the device driver. The design of the software interface lets a user with limited knowledge of the system hardware write the upper device driver.

The adapter device driver manages the transport layer but not the devices. It can send and receive commands, but it cannot interpret the contents of the command. The lower driver also provides recovery and logging for errors related to the transport layer and system I/O hardware. Management of the device specifics is left to the device driver. The interface of the two drivers supports communication between the upper driver and the different transport layer adapters without requiring special code paths for each adapter.

## Responsibilities of the Device Driver

The device driver provides the rest of the operating system with the software interface to a given device or device class. The upper layer recognizes which commands are required to control a particular device or device class. The device driver builds I/O requests containing device commands, and sends them to the adapter device driver in the sequence needed to operate the device successfully. The device driver cannot manage adapter resources or give the command to the adapter. Specifics about the adapter and system hardware are left to the lower layer.

The device driver also provides recovery and logging for errors related to the device that it controls.

The operating system provides several kernel services that let the device driver communicate with adapter device driver entry points without having the actual name or address of those entry points. See “Logical File System Kernel Services” on page 63 for more information.

## Communication between Devices

When two devices communicate, one assumes the initiator-mode role, and the other assumes the target-mode role. The initiator-mode device generates the command, which requests an operation, and the target-mode device receives the command and acts. It is possible for a device to perform both roles simultaneously.

When writing a new adapter device driver, the writer must know which mode or modes must be supported to meet the requirements of the adapter and any interfaced device drivers.

## Initiator-Mode Support

The interface between the device driver and the adapter device driver for initiator-mode support (that is, the attached device acts as a target) is accessed through calls to the adapter device driver **open**, **close**, **ioctl**, and **strategy** subroutines. I/O requests are queued to the adapter device driver through calls to its strategy entry point.

Communication between the device driver and the adapter device driver for a particular initiator I/O request is made through the **scsi\_buf** structure, which is passed to and from the **strategy** subroutine in the same way a standard driver uses a **struct buf** structure.

## Fast I/O Failure for Fibre Channel Devices

AIX supports Fast I/O Failure for Fibre Channel (FC) devices after link events in a switched environment. If the FC adapter driver detects a link event, such as a lost link between a storage device and a switch, the FC adapter driver waits a short period of time, approximately 15 seconds, so that the fabric can stabilize. At that point, if the FC adapter driver detects that the device is not on the fabric, it begins failing all I/Os at the adapter driver. Any new I/O or future retries of the failed I/Os are failed immediately by the adapter until the adapter driver detects that the device has rejoined the fabric.

Fast Failure of I/O is controlled by a new **fscsi** device attribute, **fc\_err\_recov**. The default setting for this attribute is `delayed_fail`, which is the I/O failure behavior seen in previous versions of AIX. To enable Fast I/O Failure, set this attribute to `fast_fail`, as shown in the example:

```
chdev -l fscsi0 -a fc_err_recov=fast_fail
```



In this example, the **fscsi** device instance is `fscsi0`. Fast fail logic is called when the adapter driver receives an indication from the switch that there has been a link event involving a remote storage device port by way of a Registered State Change Notification (RSCN) from the switch.

Fast I/O Failure is useful in situations where multipathing software is used. Setting the **fc\_err\_recov** attribute to `fast_fail` can decrease the I/O fail times because of link loss between the storage device and switch. This would support faster failover to alternate paths.

In single-path configurations, especially configurations with a single path to a paging device, the `delayed_fail` default setting is recommended.

Fast I/O Failure requires the following:

- A switched environment. It is not supported in arbitrated loop environments, including public loop.
- FC 6227 adapter firmware, level 3.22A1 or higher.
- FC 6228 adapter firmware, level 3.82A1 or higher.
- FC 6239 adapter firmware, all levels.
- All subsequent FC adapter releases support Fast I/O Failure.

If any of these requirements is not met, the **fscsi** device logs an error log of type INFO indicating that one of these requirements is not met and that Fast I/O Failure is not enabled.

## Dynamic Tracking of Fibre Channel Devices

AIX supports dynamic tracking of Fibre Channel (FC) devices. Previous releases of AIX required a user to unconfigure FC storage device and adapter device instances before making changes on the system area network (SAN) that might result in an N\_Port ID (SCSI ID) change of any remote storage ports.

If dynamic tracking of FC devices is enabled, the FC adapter driver detects when the Fibre Channel N\_Port ID of a device changes. The FC adapter driver then reroutes traffic destined for that device to the new address while the devices are still online. Events that can cause an N\_Port ID to change include moving a cable between a switch and storage device from one switch port to another, connecting two separate switches using an inter-switch link (ISL), and possibly rebooting a switch.

Dynamic tracking of FC devices is controlled by a new **fscsi** device attribute, **dyntrk**. The default setting for this attribute is `no`. To enable dynamic tracking of FC devices, set this attribute to `dyntrk=yes`, as shown in the example:

```
chdev -l fscsi0 -a dyntrk=yes
```

In this example, the **fscsi** device instance is `fscsi0`. Dynamic tracking logic is called when the adapter driver receives an indication from the switch that there has been a link event involving a remote storage device port.

Dynamic tracking support requires the following:

- A switched environment. It is not supported in arbitrated loop environments, including public loop.
- FC 6227 adapter firmware, level 3.22A1 or higher.
- FC 6228 adapter firmware, level 3.82A1 or higher.
- FC 6239 adapter firmware, all levels.
- All subsequent FC adapter releases support Fast I/O Failure.
- The World Wide Name (Port Name) and Node Names devices must remain constant, and the World Wide Name device must be unique. Changing the World Wide Name or Node Name of an available or online device can result in I/O failures. In addition, each FC storage device instance must have **world\_wide\_name** and **node\_name** attributes. Updated filesets that contain the **sn\_location** attribute (see the following bullet) must also be updated to contain both of these attributes.

- The storage device must provide a reliable method to extract a unique serial number for each LUN. The AIX FC device drivers do not autodetect serial number location, so the method for serial number extraction must be explicitly provided by any storage vendor in order to support dynamic tracking for their devices. This information is conveyed to the drivers using the **sn\_location** ODM attribute for each storage device. If the disk or tape driver detects that the **sn\_location** ODM attribute is missing, an error log of type INFO is generated and dynamic tracking is not enabled.

**Note:** The **sn\_location** attribute might not be displayed, so running the **lsattr** command on an **hdisk**, for example, might not show the attribute even though it could be present in ODM.

- The FC device drivers can track devices on a *SAN fabric*, which is a fabric as seen from a single host bus adapter, if the **N\_Port** IDs on the fabric stabilize within about 15 seconds. If cables are not reseated or **N\_Port** IDs continue to change after the initial 15 seconds, I/O failures could result.
- Devices are not tracked across host bus adapters. Devices only track if they remain visible from the same HBA that they were originally connected to.  
For example, if device A is moved from one location to another on fabric A attached to host bus adapter A (in other words, its **N\_Port** on fabric A changes), the device is seamlessly tracked without any user intervention, and I/O to this device can continue.  
However, if a device A is visible from HBA A but not from HBA B, and device A is moved from the fabric attached to HBA A to the fabric attached to HBA B, device A is not accessible on fabric A nor on fabric B. User intervention would be required to make it available on fabric B by running the **cfgmgr** command. The AIX device instance on fabric A would no longer be usable, and a new device instance on fabric B would be created. This device would have to be added manually to volume groups, multipath device instances, and so on. This is essentially the same as removing a device from fabric A and adding a new device to fabric B.
- No dynamic tracking can be performed for FC dump devices while an AIX system dump is in progress. In addition, dynamic tracking is not supported when booting or running the **cfgmgr** command. SAN changes should not be made while any of these operations are in progress.
- After devices are tracked, ODM might contain stale information because Small Computer System Interface (SCSI) IDs in ODM no longer reflect actual SCSI IDs on the SAN. ODM remains in this state until **cfgmgr** is run manually or the system is rebooted (provided all drivers, including any third party FC SCSI target drivers, are dynamic-tracking capable). If **cfgmgr** is run manually, **cfgmgr** must be run on all affected **fscsi** devices. This can be accomplished by running **cfgmgr** without any options, or by running **cfgmgr** on each **fscsi** device individually.

**Note:** Running **cfgmgr** at run time to recalibrate the SCSI IDs might not update the SCSI ID in ODM for a storage device if the storage device is currently opened, such as when volume groups are varied on. The **cfgmgr** command must be run on devices that are not opened or the system must be rebooted to recalibrate the SCSI IDs. Stale SCSI IDs in ODM have no adverse affect on the FC drivers, and recalibration of SCSI IDs in ODM is not necessary for the FC drivers to function properly. Any applications that communicate with the adapter driver directly using **ioctl** calls and that use the SCSI ID values from ODM, however, must be updated (see the next bullet) to avoid using potentially stale SCSI IDs.

- All applications and kernel extensions that communicate with the FC adapter driver, either through **ioctl** calls or directly to the FC driver's entry points, must support the version 1 **ioctl** and **scsi\_buf** APIs of the FC adapter driver to work properly with FC dynamic tracking. Noncompliant applications or kernel extensions might not function properly or might even fail after a dynamic tracking event. If the FC adapter driver detects an application or kernel extension that is not adhering to the new version 1 **ioctl** and **scsi\_buf** API, an error log of type INFO is generated and dynamic tracking might not be enabled for the device that this application or kernel extension is trying to communicate with.  
For ISVs developing kernel extensions or applications that communicate with the AIX Fibre Channel Driver stack, refer to the "SAM Adapter Device Driver **ioctl** Commands" on page 301 and "Understanding the **scsi\_buf** Structure" on page 290 for changes necessary to support dynamic tracking.
- Even with dynamic tracking enabled, users should make SAN changes, such as moving or swapping cables and establishing ISL links, during maintenance windows. Making SAN changes during full

production runs is discouraged because the interval of time to perform any SAN changes is too short. Cables that are not reseated correctly, for example, could result in I/O failures. Performing these operations during periods of little or no traffic minimizes the impact of I/O failures.

The base AIX FC SCSI Disk and FC SCSI Tape and FastT device drivers support dynamic tracking. The IBM ESS, EMC Symmetrix, and HDS storage devices support dynamic tracking provided that the vendor provides the ODM filesets with the necessary **sn\_location** and **node\_name** attributes. Contact the storage vendor if you are not sure if your current level of ODM fileset supports dynamic tracking.

If vendor-specific ODM entries are not being used for the storage device, but the ESS, Symmetrix, or HDS storage subsystem is configured with the MPI0 Other FC SCSI Disk message, dynamic tracking is supported for the devices in this configuration. This supersedes the need for the **sn\_location** attribute. All current AIX Path Control Modules (PCM) shipped with the AIX base support dynamic tracking.

The STK tape device using the standard AIX device driver also supports dynamic tracking provided the STK fileset contains the necessary **sn\_location** and **node\_name** attributes.

**Note:** SAN changes involving tape devices should be made with no active I/O. Because of the serial nature of tape devices, a single I/O failure can cause an application to fail, including tape backups.

Devices configured with the Other FC SCSI Disk or Other FC SCSI Tape messages do not support dynamic tracking.

## Fast I/O Failure and Dynamic Tracking Interaction

Although Fast I/O Failure and dynamic tracking of FC Devices are technically separate features, the enabling of one can change the interpretation of the other in certain situations. The following table shows the behavior exhibited by the FC drivers with the various permutations of these settings:

<b>dyntrk</b>	<b>fc_err_recov</b>	<b>FC Driver Behavior</b>
no	delayed_fail	The default setting. This is legacy behavior existing in previous versions of AIX. The FC drivers do not recover if the SCSI ID of a device changes, and I/Os take longer to fail when a link loss occurs between a remote storage port and switch. This might be preferable in single-path situations if dynamic tracking support is not a requirement.
no	fast_fail	If the driver receives a RSCN from the switch, this could indicate a link loss between a remote storage port and switch. After an initial 15-second delay, the FC drivers query to see if the device is on the fabric. If not, I/Os are flushed back by the adapter. Future retries or new I/Os fail immediately if the device is still not on the fabric. If the FC drivers detect that the device is on the fabric but the SCSI ID has changed, the FC device drivers do not recover, and the I/Os fail with PERM errors.

dyntrk	fc_err_recov	FC Driver Behavior
yes	delayed_fail	If the driver receives a RSCN from the switch, this could indicate a link loss between a remote storage port and switch. After an initial 15-second delay, the FC drivers query to see if the device is on the fabric. If not, I/Os are flushed back by the adapter. Future retries or new I/Os fail immediately if the device is still not on the fabric, although the storage driver (disk, tape, FastT) drivers might inject a small delay (2-5 seconds) between I/O retries. If the FC drivers detect that the device is on the fabric but the SCSI ID has changed, the FC device drivers reroute traffic to the new SCSI ID.
yes	fast_fail	If the driver receives a Registered State Change Notification (RSCN) from the switch, this could indicate a link loss between a remote storage port and switch. After an initial 15-second delay, the FC drivers query to see if the device is on the fabric. If not, I/Os are flushed back by the adapter. Future retries or new I/Os fail immediately if the device is still not on the fabric. The storage driver (disk, tape, FastT) will likely not delay between retries. If the FC drivers detect the device is on the fabric but the SCSI ID has changed, the FC device drivers reroute traffic to the new SCSI ID.

When dynamic tracking is disabled, there is a marked difference between the `delayed_fail` and `fast_fail` settings of the `fc_err_recov` attribute. However, with dynamic tracking enabled, the setting of the `fc_err_recov` attribute is less significant. This is because there is some overlap in the dynamic tracking and fast fail error-recovery policies. Therefore, enabling dynamic tracking inherently enables some of the fast fail logic.

The general error recovery procedure when a device is no longer reachable on the fabric is the same for both `fc_err_recov` settings with dynamic tracking enabled. The minor difference is that the storage drivers can choose to inject delays between I/O retries if `fc_err_recov` is set to `delayed_fail`. This increases the I/O failure time by an additional amount, depending on the delay value and number of retries, before permanently failing the I/O. With high I/O traffic, however, the difference between `delayed_fail` and `fast_fail` might be more noticeable.

SAN administrators might want to experiment with these settings to find the correct combination of settings for their environment.

---

## SAM Asynchronous Event Handling

A device driver can register a particular device instance for receiving asynchronous event status by calling the **SCIOLEVENT** ioctl operation for the adapter device driver. When an event covered by the **SCIOLEVENT** ioctl operation is detected by the adapter device driver, it builds an **scsi\_event\_info** structure and passes a pointer to the structure and to the asynchronous event-handler routine entry point, which was previously registered. The fields in the structure are filled in by the adapter device driver as follows:

### **scsi\_id**

For initiator mode, this is set to the SCSI ID or SCSI ID alias of the attached target device. For target mode, this is set to the SCSI ID or SCSI ID alias of the attached initiator device.

### **lun\_id**

For initiator mode, this is set to the SCSI LUN of the attached target device. For target mode, this is set to 0.

**mode** Identifies whether the initiator or target mode device is being reported. The following values are possible:

#### **SCSI\_IM\_MODE**

An initiator mode device is being reported.

#### **SCSI\_TM\_MODE**

A target mode device is being reported.

### **events**

This field is set to indicate what event or events are being reported. The following values are possible, as defined in the **/usr/include/sys/scsi.h** file:

#### **SCSI\_FATAL\_HDW\_ERR**

A fatal adapter hardware error occurred.

#### **SCSI\_ADAP\_CMD\_FAILED**

An unrecoverable adapter command failure occurred.

#### **SCSI\_RESET\_EVENT**

A transport layer reset was detected.

#### **SCSI\_BUFS\_EXHAUSTED**

In target-mode, a maximum buffer usage event has occurred.

### **adap\_devno**

This field is set to indicate the device major and minor numbers of the adapter on which the device is located.

### **async\_correlator**

This field is set to the value passed to the adapter device driver in the **scsi\_event\_struct** structure. The device driver might optionally use this field to provide an efficient means of associating event status with the device instance it goes with. Alternatively, the device driver would use the combination of the **id**, **lun**, **mode**, and **adap\_devno** fields to identify the device instance.

The information reported in the **scsi\_event\_info.events** field does not queue to the device driver, but is instead reported as one or more flags as they occur. Because the data does not queue, the adapter device driver writer can use a single **scsi\_event\_info** structure and pass it one at a time, by pointer, to each asynchronous event handler routine for the appropriate device instance. After determining for which device the events are being reported, the device driver must copy the **scsi\_event\_info.events** field into local space and must not modify the contents of the rest of the **scsi\_event\_info** structure.

Because the event status is optional, the device driver writer determines which action is necessary to take upon receiving event status. The writer might decide to save the status and report it back to the calling application, or the device driver or application level program can take error recovery actions.

**Note:** SAM asynchronous event handling is not supported by all adapter device drivers.

## Defined Events and Recovery Actions

The adapter fatal hardware failure event is intended to indicate that no further commands to or from this device are likely to succeed, because the adapter to which it is attached, has failed. It is recommended that the application end the session with the device.

The unrecoverable adapter command failure event is not necessarily a fatal condition, but it can indicate that the adapter is not functioning properly. Possible actions by the application program include:

- Ending of the session with the device in the near future.
- Ending of the session after multiple (two or more) such events.
- Attempt to continue the session indefinitely.

The SCSI Reset detection event is mainly intended as information only, but can be used by the application to perform further actions, if necessary.

The maximum buffer usage detected event only applies to a given target-mode device; it will not be reported for an initiator-mode device. This event indicates to the application that this particular target-mode device instance has filled its maximum allotted buffer space. The application should perform **read** system calls fast enough to prevent this condition. If this event occurs, data is not lost, but it is delayed to prevent further buffer usage. Data reception will be restored when the application empties enough buffers to continue reasonable operations. The **num\_bufs** attribute might need to be increased to help minimize this problem. Also, it is possible that regardless of the number of buffers, the application simply is not processing received data fast enough. This might require some fine tuning of the application's data processing routines.

## Asynchronous Event-Handling Routine

The device driver asynchronous event-handling routine is typically called directly from the hardware interrupt-handling routine for the adapter device driver. The device driver writer must be aware of how this affects the design of the device driver.

Because the event handling routine is running on the hardware interrupt level, the device driver must be careful to limit operations in that routine. Processing should be kept to a minimum. In particular, if any error recovery actions are performed, it is recommended that the event-handling routine set state or status flags only and allow a process level routine to perform the actual operations.

The device driver must be careful to disable interrupts at the correct level in places where the device driver's lower execution priority routines manipulate variables that are also modified by the event-handling routine. To allow the device driver to disable at the correct level, the adapter device driver writer must provide a configuration database attribute that defines the interrupt class, or priority, it runs on. This attribute must be named **intr\_priority** so that the device driver configuration method knows which attribute of the parent adapter to query. The device driver configuration method should then pass this interrupt priority value to the device driver along with other configuration data for the device instance.

The SAM device driver writer must follow any other general system rules for writing a routine that must execute in an interrupt environment. For example, the routine must not attempt to sleep or wait on I/O operations. It can perform wakeups to allow the process level to handle those operations.

Because the device driver copies the information from the **scsi\_event\_info.events** field on each call to its asynchronous event-handling routine, there is no resource to free and no information that must be passed back later to the adapter device driver.

---

## SAM Error Recovery

If the device is in initiator mode, the error-recovery process varies depending on whether or not the device is supporting command queuing. Also some devices might support NACA=1 error recovery. Thus, error recovery needs to deal with the two following concepts.

### Autosense Data

When a device returns a check condition or command terminated (the **scsi\_buf.scsi\_status** will have the value of `SC_CHECK_CONDITION` or `SC_COMMAND_TERMINATED`, respectively), it will also return the request sense data.

**Note:** Subsequent commands to the device will clear the request sense data.

If the device driver has specified a valid autosense buffer (**scsi\_buf.autosense\_length** > 0 and the **scsi\_buf.autosense\_buffer\_ptr** field is not NULL), then the adapter device driver will copy the returned autosense data into the buffer referenced by **scsi\_buf.autosense\_buffer\_ptr**. When this occurs, the adapter device driver will set the **SC\_AUTOSENSE\_DATA\_VALID** flag in the **scsi\_buf.adap\_set\_flags**.

When the device driver receives the SCSI status of check condition or command terminated (the **scsi\_buf.scsi\_status** will have the value of `SC_CHECK_CONDITION` or `SC_COMMAND_TERMINATED`, respectively), it should then determine if the **SC\_AUTOSENSE\_DATA\_VALID** flag is set in the **scsi\_buf.adap\_set\_flags**. If so then it should process the autosense data and not send a SCSI request sense command.

### NACA=1 error recovery

Some devices support setting the NACA (Normal Auto Contingent Allegiance) bit to a value of one (NACA=1) in the control byte of the SCSI command . If a device returns a check condition or command terminated (the **scsi\_buf.scsi\_status** will have the value of `SC_CHECK_CONDITION` or `SC_COMMAND_TERMINATED`, respectively) for a command with NACA=1 set, then the device will require a Clear ACA task management request to clear the error condition on the drive. The device driver can issue a Clear ACA task management request by sending a transaction with the **SC\_CLEAR\_ACA** flag in the **sc\_buf.flags** field. The **SC\_CLEAR\_ACA** flag can be used in conjunction with the **SC\_Q\_CLR** and **SC\_Q\_RESUME** flag in the **sc\_buf.flags** to clear or resume the queue of transactions for this device, respectively. For more information, see “Initiator-Mode Recovery During Command Tag Queuing” on page 286.

---

## SAM Initiator-Mode Recovery When Not Command Tag Queuing

If an error such as a check condition or hardware failure occurs, the transaction active during the error is returned with the **scsi\_buf.bufstruct.b\_error** field set to `EIO`. Other transactions in the queue might be returned with the **scsi\_buf.bufstruct.b\_error** field set to `ENXIO`. If the adapter driver decides not to return other outstanding commands it has queued to it, then the failed transaction will be returned to the device driver with an indication that the queue for this device is not cleared by setting the **SC\_DID\_NOT\_CLEAR\_Q** flag in the **scsi\_buf.adap\_q\_status** field. The device driver should process or recover the condition, rerunning any mode selects or device reservations to recover from this condition properly. After this recovery, it should reschedule the transaction that had the error. In many cases, the device driver only needs to retry the unsuccessful operation.

The adapter device driver should never retry a SCSI command on error after the command has successfully been given to the adapter. The consequences for retrying a command at this point range from minimal to catastrophic, depending upon the type of device. Commands for certain devices cannot be retried immediately after a failure (for example, tapes and other sequential access devices). If such an error occurs, the failed command returns an appropriate error status with an **iodone** call to the device driver for error recovery. Only the device driver that originally issued the command knows if the command can be retried on the device. The adapter device driver must only retry commands that were never

successfully transferred to the adapter. In this case, if retries are successful, the **scsi\_buf** status should not reflect an error. However, the adapter device driver should perform error logging on the retried condition.

The first transaction passed to the adapter device driver during error recovery must include a special flag. This **SC\_RESUME** flag in the **scsi\_buf.flags** field must be set to inform the adapter device driver that the device driver has recognized the fatal error and is beginning recovery operations. Any transactions passed to the adapter device driver, after the fatal error occurs and before the **SC\_RESUME** transaction is issued, should be flushed; that is, returned with an error type of ENXIO through an **iodone** call.

**Note:** If a device driver continues to pass transactions to the adapter device driver after the adapter device driver has flushed the queue, these transactions are also flushed with an error return of ENXIO through the **iodone** service. This gives the device driver a positive indication of all transactions flushed.

---

## Initiator-Mode Recovery During Command Tag Queuing

If the device driver is queuing multiple transactions to the device and either a check condition error or a command terminated error occurs, the adapter driver does not clear all transactions in its queues for the device. It returns the failed transaction to the device driver with an indication that the queue for this device is not cleared by setting the **SC\_DID\_NOT\_CLEAR\_Q** flag in the **scsi\_buf.adap\_q\_status** field. The adapter driver halts the queue for this device awaiting error recovery notification from the device driver. The device driver then has three options to recover from this error:

- Send one error recovery command (request sense) to the device.
- Clear the adapter driver's queue for this device.
- Resume the adapter driver's queue for this device.

When the adapter driver's queue is halted, the device driver can get sense data from a device by setting the **SC\_RESUME** flag in the **scsi\_buf.flags** field and the **SC\_NO\_Q** flag in **scsi\_buf.q\_tag\_msg** field of the request-sense **scsi\_buf**. This action notifies the adapter driver that this is an error-recovery transaction and should be sent to the device while the remainder of the queue for the device remains halted. When the request sense completes, the device driver needs to either clear or resume the adapter driver's queue for this device.

The device driver can notify the adapter driver to clear its halted queue by sending a transaction with the **SC\_Q\_CLR** flag in the **scsi\_buf.flags** field. This transaction must not contain a command because it is cleared from the adapter driver's queue without being sent to the adapter. However, this transaction must have the SCSI ID field ( **scsi\_buf.scsi\_id** ) and the LUN field ( **scsi\_buf.lun\_id** ) filled in with the device's SCSI ID and logical unit number (LUN), respectively. Upon receiving an **SC\_Q\_CLR** transaction, the adapter driver flushes all transactions for this device and sets their **scsi\_buf.bufstruct.b\_error** fields to ENXIO. The device driver must wait until the **scsi\_buf** with the **SC\_Q\_CLR** flag set is returned before it resumes issuing transactions. The first transaction sent by the device driver after it receives the returned **SC\_Q\_CLR** transaction must have the **SC\_RESUME** flag set in the **scsi\_buf.flags** fields.

If the device driver wants the adapter driver to resume its halted queue, it must send a transaction with the **SC\_Q\_RESUME** flag set in the **scsi\_buf.flags** field. This transaction can contain an actual command, but it is not required. However, this transaction must have the SCSI ID field (**scsi\_buf.scsi\_id**) and the LUN field (**scsi\_buf.lun\_id**) filled in with the device's SCSI ID and logical unit number (LUN). If this is the first transaction issued by the device driver after receiving the error (indicating that the adapter driver's queue is halted), then the **SC\_RESUME** flag must be set as well as the **SC\_Q\_RESUME** flag.

## Analyzing Returned Status

The following order of precedence should be followed by device drivers when analyzing the returned status:



1. If the **scsi\_buf.bufstruct.b\_flags** field has the **B\_ERROR** flag set, then an error has occurred and the **scsi\_buf.bufstruct.b\_error** field contains a valid **errno** value.
 

If the **b\_error** field contains the **ENXIO** value, either the command needs to be restarted or it was canceled at the request of the device driver.

If the **b\_error** field contains the **EIO** value, then either one or no flag is set in the **scsi\_buf.status\_validity** field. If a flag is set, an error in either the **scsi\_status** or **adapter\_status** field is the cause.

If the **status\_validity** field is 0, then the **scsi\_buf.bufstruct.b\_resid** field should be examined to see if the command issued was in error. The **b\_resid** field can have a value without an error having occurred. To decide whether an error has occurred, the device driver must evaluate this field with regard to the command being sent and the device being driven.

If the **SC\_CHECK\_CONDITION** or **SC\_COMMAND\_TERMINATED** is set in **scsi\_status**, then a device driver must analyze the value of **scsi\_buf.adap\_set\_flags** to determine if autosense data was returned from the device.

If the **SC\_AUTOSENSE\_DATA\_VALID** flag is set in the **scsi\_buf.adap\_set\_flags** field for a device, then the device returned autosense data in the buffer referenced by **scsi\_buf.autosense\_buffer\_ptr**. In this situation the device driver does not need to issue a SCSI request sense to determine the appropriate error recovery for the devices.

If the device driver is queuing multiple transactions to the device and if either **SC\_CHECK\_CONDITION** or **SC\_COMMAND\_TERMINATED** is set in **scsi\_status**, then the value of **scsi\_buf.adap\_q\_status** must be analyzed to determine if the adapter driver has cleared its queue for this device. If the adapter driver has not cleared its queue after an error, then it holds that queue in a halted state.

If **scsi\_buf.adap\_q\_status** is set to 0, the adapter driver has cleared its queue for this device and any transactions outstanding are flushed back to the device driver with an error of **ENXIO**.

If the **SC\_DID\_NOT\_CLEAR\_Q** flag is set in the **scsi\_buf.adap\_q\_status** field, the adapter driver has not cleared its queue for this device. When this condition occurs, the adapter driver allows the device driver to send one error recovery transaction (request sense) that has the field **scsi\_buf.q\_tag\_msg** set to **SC\_NO\_Q** and the field **scsi\_buf.flags** set to **SC\_RESUME**. The device driver can then notify the adapter driver to clear or resume its queue for the device by sending a **SC\_Q\_CLR** or **SC\_Q\_RESUME** transaction.

If the device driver does not queue multiple transactions to the device (that is, the **SC\_NO\_Q** is set in **scsi\_buf.q\_tag\_msg**), then the adapter clears its queue on error and sets **scsi\_buf.adap\_q\_status** to 0.
2. If the **scsi\_buf.bufstruct.b\_flags** field does not have the **B\_ERROR** flag set, then no error is being reported. However, the device driver should examine the **b\_resid** field to check for cases where less data was transferred than expected. For some commands, this occurrence might not represent an error. The device driver must determine if an error has occurred.
 

If a nonzero **b\_resid** field does represent an error condition, then the device queue is not halted by the adapter device driver. It is possible for one or more succeeding queued commands to be sent to the adapter (and possibly the device). Recovering from this situation is the responsibility of the device driver.
3. In any of the above cases, if **scsi\_buf.bufstruct.b\_flags** field has the **B\_ERROR** flag set, then the queue of the device in question has been halted. The first **scsi\_buf** structure sent to recover the error (or continue operations) must have the **SC\_RESUME** bit set in the **scsi\_buf.flags** field.

---

## A Typical Initiator-Mode SAM Driver Transaction Sequence

A simplified sequence of events for a transaction between a device driver and an adapter device driver follows. In this sequence, routine names preceded by **dd\_** are part of the device driver, and those preceded by **scsi\_** are part of the adapter device driver.

1. The device driver receives a call to its **dd\_strategy** routine; any required internal queuing occurs in this routine. The **dd\_strategy** entry point then triggers the operation by calling the **dd\_start** entry

- point. The **dd\_start** routine invokes the **scsi\_strategy** entry point by calling the **devstrategy** kernel service with the relevant **scsi\_buf** structure as a parameter.
2. The **scsi\_strategy** entry point initially checks the **scsi\_buf** structure for validity. These checks include validating the **devno** field, matching the SCSI ID or the LUN to internal tables for configuration purposes, and validating the request size.
  3. Although the adapter device driver cannot reorder transactions, it does perform queue chaining. If no other transactions are pending for the requested device, the **scsi\_strategy** routine immediately calls the **scsi\_start** routine with the new transaction. If there are other transactions pending, the new transaction is added to the tail of the device chain.
  4. At each interrupt, the **scsi\_intr** interrupt handler verifies the current status. The adapter device driver fills in the **scsi\_buf status\_validity** field, updating the **scsi\_status** and **adapter\_status** fields as required. The adapter device driver also fills in the **bufstruct.b\_resid** field with the number of bytes not transferred from the request. If all the data was transferred, the **b\_resid** field is set to a value of 0. If the SCSI adapter driver is a adapter driver and autosense data is returned from the device, then the adapter driver will also fill in the **adap\_set\_flags** and **autosense\_buffer\_ptr** fields of the **scsi\_buf** structure. When a transaction completes, the **scsi\_intr** routine causes the **scsi\_buf** entry to be removed from the device queue and calls the **iodone** kernel service, passing the just dequeued **scsi\_buf** structure for the device as the parameter. The **scsi\_start** routine is then called again to process the next transaction on the device queue. The **iodone** kernel service calls the device driver **dd\_iodone** entry point, signaling the device driver that the particular transaction has completed.
  5. The device driver **dd\_iodone** routine investigates the I/O completion codes in the **scsi\_buf** status entries and performs error recovery, if required. If the operation completed correctly, the device driver dequeues the original buffer structures. It calls the **iodone** kernel service with the original buffer pointers to notify the originator of the request.

---

## Understanding the Execution of SAM Initiator I/O Requests

During normal processing, many transactions are queued in the device driver. As the device driver processes these transactions and passes them to the adapter device driver, the device driver moves them to the in-process queue. When the adapter device driver returns through the **iodone** service with one of these transactions, the device driver either recovers any errors on the transaction or returns using the **iodone** kernel service to the calling level.

The device driver can send only one **scsi\_buf** structure per call to the adapter device driver. Thus, the **scsi\_buf.bufstruct.av\_forw** pointer should be null when given to the adapter device driver, which indicates that this is the only request. The device driver can queue multiple **scsi\_buf** requests by making multiple calls to the adapter device driver strategy routine.

## Spanned (Consolidated) Commands

Some kernel operations may be composed of sequential operations to a device. For example, if consecutive blocks are written to disk, blocks might or might not be in physically consecutive buffer pool blocks.

To enhance the transport layer performance, the device driver should consolidate multiple queued requests when possible into a single command. To allow the adapter device driver the ability to handle the scatter and gather operations required, the **scsi\_buf.bp** should always point to the first **buf** structure entry for the spanned transaction. A null-terminated list of additional **struct buf** entries should be chained from the first field through the **buf.av\_forw** field to give the adapter device driver enough information to perform the DMA scatter and gather operations required. This information must include at least the buffer's starting address, length, and cross-memory descriptor.

The spanned requests should always be for requests in either the read or write direction but not both, since the adapter device driver must be given a single command to handle the requests. The spanned

request should always consist of complete I/O requests (including the additional **struct buf** entries). The device driver should not attempt to use partial requests to reach the maximum transfer size.

The maximum transfer size is actually adapter-dependent. The **IOCINFO** ioctl operation can be used to discover the adapter device driver's maximum allowable transfer size.

If a transfer size larger than the supported maximum is attempted, the adapter device driver returns a value of **EINVAL** in the **scsi\_buf.bufstruct.b\_error** field.

Due to system hardware requirements, the device driver must consolidate only commands that are memory page-aligned at both their starting and ending addresses. Specifically, this applies to the consolidation of *inner* memory buffers. The ending address of the first buffer and the starting address of all subsequent buffers should be memory page-aligned. However, the starting address of the first memory buffer and the ending address of the last do not need to be aligned so.

The purpose of consolidating transactions is to decrease the number of commands and transport layer phases required to perform the required operation. The time required to maintain the simple chain of **buf** structure entries is significantly less than the overhead of multiple (even two) transport layer transactions.

## Fragmented Commands

Single I/O requests larger than the maximum transfer size must be divided into smaller requests by the device driver. For calls to a device driver's character I/O (read/write) entry points, the **uphysio** kernel service can be used to break up these requests. For a *fragmented command* such as this, the **scsi\_buf.bp** field should be **null** so that the adapter device driver uses only the information in the **scsi\_buf** structure to prepare for the DMA operation.

---

## SAM Command Tag Queuing

Command tag queuing refers to queuing multiple commands to a device. Queuing to the device can improve performance because the device itself determines the most efficient way to order and process commands. Devices that support command tag queuing can be divided into two classes: those that clear their queues on error and those that do not. Devices that do not clear their queues on error resume processing of queued commands when the error condition is cleared (either by receiving the next command for **NACA=0** error recovery or by receiving a Clear ACA task management command for **NACA=1** error recovery). Devices that do clear their queues flush all commands currently outstanding.

Command tag queuing requires the adapter, the device, the device driver, and the adapter driver to support this capability. For a device driver to queue multiple commands to a device (that supports command tag queuing), it must be able to provide at least one of the following values in the **scsi\_buf.q\_tag\_msg**:

- **SC\_SIMPLE\_Q**
- **SC\_HEAD\_OF\_Q**
- **SC\_ORDERED\_Q**

The disk device driver and adapter driver do support this capability. This implementation provides some queuing-specific changeable attributes for disks that can queue commands. With this information, the disk device driver attempts to queue to the disk, first by queuing commands to the adapter driver. The adapter driver then queues these commands to the adapter, providing that the adapter supports command tag queuing. If the adapter does not support command tag queuing, then the adapter driver sends only one command at a time to the adapter and so multiple commands are not queued to the disk.

**Note:** This operation might not be supported by all adapter device drivers.

---

## Understanding the `scsi_buf` Structure

The `scsi_buf` structure is used for communication between the device driver and the adapter device driver during an initiator I/O request. This structure is passed to and from the strategy routine in the same way a standard driver uses a `struct buf` structure.

### Fields in the `scsi_buf` Structure

The `scsi_buf` structure contains certain fields used to pass a command and associated parameters to the adapter device driver. Other fields within this structure are used to pass returned status back to the device driver. The `scsi_buf` structure is defined in the `/usr/include/sys/scsi_buf.h` file.

Fields in the `scsi_buf` structure are used as follows:

- Reserved fields should be set to a value of 0, except where noted.
- The `bufstruct` field contains a copy of the standard `buf` buffer structure that documents the I/O request. Included in this structure, for example, are the buffer address, byte count, and transfer direction. The `b_work` field in the `buf` structure is reserved for use by the adapter device driver. The current definition of the `buf` structure is in the `/usr/include/sys/buf.h` include file.
- The `bp` field points to the original buffer structure received by the device driver from the caller, if any. This can be a chain of entries in the case of spanned transfers (commands that transfer data from or to more than one system-memory buffer). A `null` pointer indicates a nonspanned transfer. The `null` value specifically tells the adapter device driver that all the information needed to perform the DMA data transfer is contained in the `bufstruct` fields of the `scsi_buf` structure.
- The `scsi_command` field, defined as a `scsi_cmd` structure, contains, for example, the SCSI command length, SCSI command, and a flag variable:

- The `scsi_length` field is the number of bytes in the actual SCSI command. This is normally 6, 10, 12, or 16 (decimal).
- The `FCP_flags` field contains the following bit flags:

#### `SC_NODISC`

Do not allow the target to disconnect during this command.

#### `SC_ASYNC`

Do not allow the adapter to negotiate for synchronous transfer to the device.

During normal use, the `SC_NODISC` bit should not be set. Setting this bit allows a device running commands to monopolize the transport layer. Sometimes it is desirable for a particular device to maintain control of the transport layer once it has successfully arbitrated for it; for instance, when this is the only device on the transport layer or the only device that will be in use. For performance reasons, it might not be desirable to go through selections again to save transport layer overhead on each command.

Also during normal use, the `SC_ASYNC` bit must not be set. It should be set only in cases where a previous command to the device ended in an unexpected transport free condition. This condition is noted as `SCSI_TRANSPORT_FAULT` in the `adapter_status` field of the `scsi_cmd` structure. Because other errors might also result in the `SCSI_TRANSPORT_FAULT` flag being set, the `SC_ASYNC` bit should only be set on the last retry of the failed command.

- The `scsi_cdb` structure contains the physical SCSI command block. The 6 to 16 bytes of a single SCSI command are stored in consecutive bytes, with the op code identified individually. The `scsi_cdb` structure contains the following fields:

#### `scsi_op_code`

This field specifies the standard SCSI op code for this command.

#### `scsi_bytes`

This field contains the remaining command-unique bytes of the command block. The actual number of bytes depends on the value in the `scsi_op_code` field.

- The **timeout\_value** field specifies the time-out limit (in seconds) to be used for completion of this command. A time-out value of 0 means no time-out is applied to this I/O request.
- The **status\_validity** field contains an output parameter that can have one of the following bit flags as a value:

**SC\_SCSI\_ERROR**

The scsi\_status field is valid.

**SC\_ADAPTER\_ERROR**

The adapter\_status field is valid.

- The **scsi\_status** field in the **scsi\_buf** structure is an output parameter that provides valid command completion status when its **status\_validity** bit is nonzero. The **scsi\_buf.bufstruct.b\_error** field should be set to **EIO** any time the **scsi\_status** field is valid. Typical status values include:

**SC\_GOOD\_STATUS**

The target successfully completed the command.

**SC\_CHECK\_CONDITION**

The target is reporting an error, exception, or other conditions.

**SC\_BUSY\_STATUS**

The target is currently transporting and cannot accept a command now.

**SC\_RESERVATION\_CONFLICT**

The target is reserved by another initiator and cannot be accessed.

**SC\_COMMAND\_TERMINATED**

The target terminated this command after receiving a terminate I/O process message from the adapter.

**SC\_QUEUE\_FULL**

The target's command queue is full, so this command is returned.

**SC\_ACA\_ACTIVE**

The device has an ACA (auto contingent allegiance) condition that requires a Clear ACA request to clear it.

- The **adapter\_status** field is an output parameter that is valid when its **status\_validity** bit is nonzero. The **scsi\_buf.bufstruct.b\_error** field should be set to **EIO** any time the **adapter\_status** field is valid. This field contains generic adapter card status. It is intentionally general in coverage so that it can report error status from any typical adapter.

If an error is detected while an command is running, and the error prevented the command from actually being sent to the transport layer by the adapter, then the error should be processed or recovered, or both, by the adapter device driver.

If it is recovered successfully by the adapter device driver, the error is logged, as appropriate, but is not reflected in the **adapter\_status** byte. If the error cannot be recovered by the adapter device driver, the appropriate **adapter\_status** bit is set and the **scsi\_buf** structure is returned to the device driver for further processing.

If an error is detected after the command was actually sent to the device, then it should be processed or recovered, or both, by the device driver.

For error logging, the adapter device driver logs transport layer and adapter-related conditions, and the device driver logs device-related errors. In the following description, a capital letter (A) after the error name indicates that the adapter device driver handles error logging. A capital letter (H) indicates that the device driver handles error logging.

Some of the following error conditions indicate a device failure. Others are transport layer or adapter-related.

**SCSI\_HOST\_IO\_BUS\_ERR (A)**

The system I/O transport layer generated or detected an error during a DMA or Programmed I/O (PIO) transfer.

**SCSI\_TRANSPORT\_FAULT (H)**

The transport protocol or hardware was unsuccessful.

**SCSI\_CMD\_TIMEOUT (H)**

The command timed out before completion.

**SCSI\_NO\_DEVICE\_RESPONSE (H)**

The target device did not respond to selection phase.

**SCSI\_ADAPTER\_HDW\_FAILURE (A)**

The adapter indicated an onboard hardware failure.

**SCSI\_ADAPTER\_SFW\_FAILURE (A)**

The adapter indicated microcode failure.

**SCSI\_FUSE\_OR\_TERMINAL\_PWR (A)**

The adapter indicated a blown terminator fuse or bad termination.

**SCSI\_TRANSPORT\_RESET (A)**

The adapter indicated the transport layer has been reset.

**SCSI\_WW\_NAME\_CHANGE (A)**

The adapter indicated the device at this SCSI ID has a new world wide name. For AIX 5.2 with 5200-01 and later, if **Dynamic Tracing of FC Devices** is enabled, the adapter driver has detected a change to the **scsi\_id** field for this device and a **scsi\_buf** structure with the **SC\_DEV\_RESTART** flag can be sent to the device. For more information, see 294.

**Note:** When **Dynamic Tracing of FC Devices** is enabled, an adapter status of **SCSI\_WW\_NAME\_CHANGE** might mean that the SCSI ID of a given world wide name on the fabric has changed, and not that the world wide name changed.

An adapter status of **SCSI\_WW\_NAME\_CHANGE** should be interpreted more generally as a situation where the SCSI ID-to-WWN mapping has changed when dynamic tracking is enabled as opposed to interpreting this literally as a world wide name change for this SCSI ID.

If dynamic tracking is *disabled*, the FC adapter driver assumes that the SCSI ID-to-WWN mapping cannot change. If a cable is moved from remote target port *A* to target port *B*, and target port *B* assumes the SCSI ID that previously belonged to target port *A*, then from the perspective of the driver with dynamic tracking disabled, the world wide name at this SCSI ID has changed.

With dynamic tracking *enabled*, the general error recovery logic in this case is different. The SCSI ID is considered volatile, so devices are tracked by world wide name. As such, all queries after events such as those described in the above text, are based on world wide name. The situation described in the previous paragraph would most likely result in a **SCSI\_NO\_DEVICE\_RESPONSE** status, since it would be determined that the world wide name of port *A* is no longer reachable. If a cable connected to port *A* was instead moved from one switch port to another, the SCSI ID of port *A* on the remote target might change. The FC adapter driver will return **SCSI\_WW\_NAME\_CHANGE** in this case, even though the SCSI ID is what actually changed, and not the world wide name.

**SCSI\_TRANSPORT\_BUSY (A)**

The adapter indicated the transport layer is busy.

**SCSI\_TRANSPORT\_DEAD (A)**

The adapter indicated the transport layer currently inoperative and is likely to remain this way for an extended time.

- The **add\_status** field contains additional device status. For devices, this field contains the Response code returned.
- When the device driver queues multiple transactions to a device, the **adap\_q\_status** field indicates whether or not the adapter driver has cleared its queue for this device after an error has occurred. The

flag of `SC_DID_NOT_CLEAR_Q` indicates that the adapter driver has not cleared its queue for this device and that it is in a halted state (so none of the pending queued transactions are sent to the device).

- The `q_tag_msg` field indicates if the adapter can attempt to queue this transaction to the device. This information causes the adapter to fill in the Queue Tag Message Code of the queue tag message for a command. The following values are valid for this field:

#### `SC_NO_Q`

Specifies that the adapter does not send a queue tag message for this command, and so the device does not allow more than one command on its command queue. This value must be used for all commands sent to devices that do not support command tag queuing.

#### `SC_SIMPLE_Q`

Specifies placing this command in the device's command queue. The device determines the order that it executes commands in its queue. The SCSI-2 specification calls this value the "Simple Queue Tag Message".

#### `SC_HEAD_OF_Q`

Specifies placing this command first in the device's command queue. This command does not preempt an active command at the device, but it is run before all other commands in the command queue. The SCSI-2 specification calls this value the "Head of Queue Tag Message".

#### `SC_ORDERED_Q`

Specifies placing this command in the device's command queue. The device processes these commands in the order that they are received. The SCSI-2 specification calls this value the "Ordered Queue Tag Message".

#### `SC_ACA_Q`

Specifies placing this command in the device's command queue, when the device has an ACA (Auto Contingent Allegiance) condition. The SCSI-3 Architecture Model calls this value the "ACA task attribute".

**Note:** Commands with the value of `SC_NO_Q` for the `q_tag_msg` field (except for request sense commands) should not be queued to a device whose queue contains a command with another value for `q_tag_msg`. If commands with the `SC_NO_Q` value (except for request sense) are sent to the device, then the device driver must make sure that no active commands are using different values for `q_tag_msg`. Similarly, the device driver must also make sure that a command with a `q_tag_msg` value of `SC_ORDERED_Q`, `SC_HEAD_Q`, or `SC_SIMPLE_Q` is not sent to a device that has a command with the `q_tag_msg` field of `SC_NO_Q`.

- The `flags` field contains bit flags sent from the device driver to the adapter device driver. The following flags are defined:

#### `SC_CLEAR_ACA`

When set, means the SCSI adapter driver should issue a Clear ACA task management request for this ID/LUN. This flag should be used in conjunction with either the `SC_Q_CLR` or `SC_Q_RESUME` flags to clear or resume the SCSI adapter driver's queue for this device. If neither of these flags is used, then this transaction is treated as if the `SC_Q_RESUME` flag is also set. The transaction containing the `SC_CLEAR_ACA` flag setting does not require an actual SCSI command in the `sc_buf`. If this transaction contains a SCSI command then it will be processed depending on whether `SC_Q_CLR` or `SC_Q_RESUME` is set.

This transaction must have the SCSI ID field (`scsi_buf.scsi_id`) and the LUN field (`scsi_buf.lun_id`) filled in with the device's SCSI ID and logical unit number (LUN). This flag is valid only during error recovery of a check condition or command terminated at a command tag queuing.

#### `SC_DELAY_CMD`

When set, means the adapter device driver should delay sending this command (following a

reset or BDR to this device) by at least the number of seconds specified to the adapter device driver in its configuration information. For devices that do not require this function, this flag should not be set.

#### **SC\_DEV\_RESTART**

If a **scsi\_buf** request fails with a status of **SCSI\_WW\_NAME\_CHANGE**, a **scsi\_buf** request with the **SC\_DEV\_RESTART** flag can be sent if the device driver is dynamic tracking capable.

For AIX 5.2 with 5200-01 and later, if **Dynamic Tracking of FC Devices** is enabled, a **scsi\_buf** request with **SC\_DEV\_RESTART** performs a handshake, indicating that the device driver acknowledges the device address change and that the FC adapter driver can proceed with tracking operations. If the **SC\_DEV\_RESTART** flag is set, then the **SC\_Q\_CLR** flag must also be set. In addition, no scsi command can be included in this **scsi\_buf** structure. Failure to meet these two criteria will result in a failure with adapter status of **SCSI\_ADAPTER\_SFW\_FAILURE**.

After the **SC\_DEV\_RESTART** call completes successfully, the device driver performs device validation procedures, such as those performed during an open (Test Unit Ready, Inquiry, Serial Number validation, etc.), in order to confirm the identity of the device after the fabric event.

If an **SC\_DEV\_RESTART** call fails with any adapter status, the **SC\_DEV\_RESTART** call can be retried as deemed appropriate by the device driver, because a future retry might succeed.

#### **SC\_LUN\_RESET**

When set, means the SCSI adapter driver should issue a Lun Reset task management request for this ID/LUN. This flag should be used in conjunction with the **SC\_Q\_CLR** flag. The transaction containing this flag setting does allow an actual command to be sent to the adapter driver. However, this transaction must have the SCSI ID field (**scsi\_buf.scsi\_id**) and the LUN field (**scsi\_buf.lun\_id**) filled in with the device's SCSI ID and logical unit number (LUN). If the transaction containing this flag setting is the first issued by the device driver after it receives an error (indicating that the adapter driver's queue is halted), then the **SC\_RESUME** flag must be set also.

#### **SC\_Q\_CLR**

When set, means the adapter driver should clear its transaction queue for this ID/LUN. The transaction containing this flag setting does not require an actual command in the **scsi\_buf** because it is flushed back to the device driver with the rest of the transactions for this ID/LUN. However, this transaction must have the SCSI ID field (**scsi\_buf.scsi\_id**) and the LUN field (**scsi\_buf.lun\_id**) filled in with the device's SCSI ID and logical unit number (LUN). This flag is valid only during error recovery of a check condition or command ended at a command tag queuing device when the **SC\_DID\_NOT\_CLR\_Q** flag is set in the **scsi\_buf.adap\_q\_status** field.

#### **SC\_Q\_RESUME**

When set, means that the adapter driver should resume its halted transaction queue for this ID/LUN. The transaction containing this flag setting does not require an actual command to be sent to the adapter driver. However, this transaction must have the SCSI ID field (**scsi\_buf.scsi\_id**) and the LUN field (**scsi\_buf.lun\_id**) filled in with the device's SCSI ID and logical unit number (LUN). If the transaction containing this flag setting is the first issued by the device driver after it receives an error (indicating that the adapter driver's queue is halted), then the **SC\_RESUME** flag must be set also.

#### **SC\_RESUME**

When set, means the adapter device driver should resume transaction queuing for this ID/LUN. Error recovery is complete after a **SCIOHALT** operation, check condition, or severe transport error. This flag is used to restart the adapter device driver following a reported error.

#### **SC\_TARGET\_RESET**

When set, means the SCSI adapter driver should issue a Target Reset task management request for this ID/LUN. This flag should be used in conjunction with the **SC\_Q\_CLR** flag. The transaction containing this flag setting does allow an actual command to be sent to the adapter driver. However, this transaction must have the SCSI ID field (**scsi\_buf.scsi\_id**) filled in



with the device's SCSI ID. If the transaction containing this flag setting is the first issued by the device driver after it receives an error (indicating that the adapter driver's queue is halted), then the **SC\_RESUME** flag must be set also.

- The **dev\_flags** field contains additional values sent from the device driver to the adapter device driver. The following values are defined:

#### **FC\_CLASS1**

When set, this tells the SCSI adapter driver that it should issue this request as a Fibre Channel Class 1 request. If the SCSI adapter driver does not support this class, then it will fail the **scsi\_buf** with an error of EINVAL. If no Fibre Channel Class is specified in the **scsi\_buf** then the SCSI adapter will default to a Fibre Channel Class.

#### **FC\_CLASS2**

When set, this tells the SCSI adapter driver that it should issue this request as a Fibre Channel Class 2 request. If the SCSI adapter driver does not support this class, then it will fail the **scsi\_buf** with an error of EINVAL. If no Fibre Channel Class is specified in the **scsi\_buf** then the SCSI adapter will default to a Fibre Channel Class.

#### **FC\_CLASS3**

When set, this tells the SCSI adapter driver that it should issue this request as a Fibre Channel Class 3 request. If the SCSI adapter driver does not support this class, then it will fail the **scsi\_buf** with an error of EINVAL. If no Fibre Channel Class is specified in the **scsi\_buf** then the SCSI adapter will default to a Fibre Channel Class.

#### **FC\_CLASS4**

When set, this tells the SCSI adapter driver that it should issue this request as a Fibre Channel Class 4 request. If the SCSI adapter driver does not support this class, then it will fail the **scsi\_buf** with an error of EINVAL. If no Fibre Channel Class is specified in the **scsi\_buf** then the SCSI adapter will default to a Fibre Channel Class.

- The **add\_work** field is reserved for use by the adapter device driver.
- The **adap\_set\_flags** field contains an output parameter that can have one of the following bit flags as a value:

#### **SC\_AUTOSENSE\_DATA\_VALID**

Autosense data was placed in the autosense buffer referenced by the **autosense\_buffer\_ptr** field.

- The **autosense\_length** field contains the length in bytes of the SCSI device driver's sense buffer, which is referenced via the **autosense\_buffer\_ptr** field. For devices this field must be non-zero, otherwise the autosense data will be lost.
- The **autosense\_buffer\_ptr** field contains the address of the SCSI devices driver's autosense buffer for this command. For devices this field must be non-NULL, otherwise the autosense data will be lost.
- The **dev\_burst\_len** field contains the burst size if this write operation in bytes. This should only be set by the device driver if it has negotiated with the device and it allows burst of write data without transfer ready. For most operations, this should be set to 0.
- The **scsi\_id** field contains the 64-bit SCSI ID for this device. This field must be set for devices.
- The **lun\_id** field contains the 64-bit lun ID for this device. This field must be set for devices.
- The **kernext\_handle** field contains the pointer returned from the **kernext\_handle** field of the **scsi\_sciolst** argument for the **SCIOLSTART** ioctl operation. For AIX 5.2 with 5200-01 and later, if **Dynamic Tracking of FC Devices** is enabled, the **kernext\_handle** field must be set for all **scsi\_buf** calls that are sent to the the adapter driver. Failure to do so results in a failure with an adapter status of **SCSI\_ADAPTER\_SFW\_FAILURE**.
- The **version** field contains the version of this **scsi\_buf** structure. Beginning with AIX 5.2, this field should be set to a value of **SCSI\_VERSION\_1**. The **version** field of the **scsi\_buf** structure should be consistent with the version of the **scsi\_sciolst** argument used for the **SCIOLSTART** ioctl operation.

---

## Other SAM Design Considerations

The following topics cover design considerations of device and adapter device drivers:

- Responsibilities of the Device Driver
- Options to the `openx` Subroutine
- Using the `SC_FORCED_OPEN` Option
- Using the `SC_RETAIN_RESERVATION` Option
- Using the `SC_DIAGNOSTIC` Option
- Using the `SC_NO_RESERVE` Option
- Using the `SC_SINGLE` Option
- Closing the Device
- Error Processing
- Device Driver and Adapter Device Driver Interfaces
- Performing Dumps

### Responsibilities of the Device Driver

SAM device drivers are responsible for the following actions:

- Interfacing with block I/O and logical-volume device-driver code in the operating system.
- Translating I/O requests from the operating system into commands suitable for the particular device. These commands are then given to the adapter device driver for execution.
- Issuing any and all commands to the attached device. The adapter device driver sends no commands except those it is directed to send by the calling device driver.
- Managing device reservations and releases. In the operating system, it is assumed that other initiators might be active on the transport layer. Usually, the device driver establishes a reservation at open time and releases it at close time (except when told to do otherwise through parameters in the device driver interface). Depending upon the type of reservations that are used, it might be necessary for a device driver to monitor conditions such as Unit Attention, which indicates that a change is the state of the previously established reservation.

### Options to the `openx` Subroutine

Device drivers must support eight defined extended options in their open routine (that is, an **openx** subroutine). Additional extended options to the open are also allowed, but they must not conflict with predefined open options. The defined extended options are bit flags in the `ext` open parameter. These options can be specified singly or in combination with each other. The required `ext` options are defined in the `/usr/include/sys/scsi.h` header file and can have one of the following values:

#### **SC\_FORCED\_OPEN**

Does not honor device reservation-conflict status and can include a target-level device reset.

#### **SC\_RETAIN\_RESERVATION**

Does not release device on close and can include a logical unit number (LUN)-level device reset.

#### **SC\_DIAGNOSTIC**

Enters diagnostic mode for this device.

#### **SC\_NO\_RESERVE**

Prevents the reservation of the device during an **openx** subroutine call to that device. Allows multiple hosts to share a device.

#### **SC\_SINGLE**

Places the selected device in Exclusive Access mode.

## Using the **SC\_FORCED\_OPEN** Option

The **SC\_FORCED\_OPEN** option forces access to a device by taking action to remove any type of reservation on the device that can inhibit access. The type of action to remove the reservation depends upon the specific type of the established reservation. You can use a **SCIORESET** ioctl to perform a target-level reset of the device. After the action to remove the reservation is completed, other commands are sent as in a normal open. If any of the commands fail due to a reservation conflict, the open registers the failure as an **EBUSY** status. This is also the result if a reservation conflict occurs during a normal open. The device driver must require the caller to have appropriate authority to request the **SC\_FORCED\_OPEN** option because this request can force a device to drop a reservation. If the caller attempts to initiate this system call without the proper authority, the device driver should return a value of -1, with the **errno** global variable set to a value of **EPERM**.

## Using the **SC\_FORCED\_OPEN\_LUN** Option

The **SC\_FORCED\_OPEN\_LUN** option forces access to a device by taking action to remove any type of reservation on the device that can inhibit access. The type of action needed to remove the reservation depends on the specific type of reservation established. You can use a **SCIORESET** ioctl to perform a LUN level reset of the device. After the action to remove the reservation is completed, other commands are sent as in a normal open. If any of the commands fail due to a reservation conflict, the open registers the failure as an **EBUSY** status. This is also the result if a reservation conflict occurs during a normal open. The device driver must require the caller to have appropriate authority to request the **SC\_FORCED\_OPEN\_LUN** option because this request can force a device to drop a reservation. If the caller tries to initiate this system call without the proper authority, the device driver returns a value of -1, with the **errno** global variable set to a value of **EPERM**.

## Using the **SC\_RETAIN\_RESERVATION** Option

The **SC\_RETAIN\_RESERVATION** option causes the device driver to hold any established reservation as part of an open during the close of the device. For shared devices (for example, disk or CD-ROM), the device driver must OR together this option for all opens to a given device. If any caller requests this option, the **close** routine does not issue the release even if other opens to the device do not set **SC\_RETAIN\_RESERVATION**. The device driver should require the caller to have appropriate authority to request the **SC\_RETAIN\_RESERVATION** option because this request can allow a program to monopolize a device (for example, if this is a nonshared device). If the caller attempts to initiate this system call without the proper authority, the device driver should return a value of -1, with the **errno** global variable set to a value of **EPERM**.

## Using the **SC\_DIAGNOSTIC** Option

The **SC\_DIAGNOSTIC** option causes the device driver to enter Diagnostic mode for the given device. This option directs the device driver to perform only minimal operations to open a logical path to the device. No commands should be sent to the device in the **open** or **close** routine when the device is in Diagnostic mode. One or more ioctl operations should be provided by the device driver to allow the caller to issue commands to the attached device for diagnostic purposes.

The **SC\_DIAGNOSTIC** option gives the caller an exclusive open to the selected device. This option requires appropriate authority to run. If the caller attempts to execute this system call without the proper authority, the device driver should return a value of -1, with the **errno** global variable set to a value of **EPERM**. The **SC\_DIAGNOSTIC** option may be executed only if the device is not already opened for normal operation. If this ioctl operation is attempted when the device is already opened, or if an **openx** call with the **SC\_DIAGNOSTIC** option is already in progress, a return value of -1 should be passed, with the **errno** global variable set to a value of **EACCES**. Once successfully opened with the **SC\_DIAGNOSTIC** flag, the device driver is placed in Diagnostic mode for the selected device.

## Using the SC\_NO\_RESERVE Option

The **SC\_NO\_RESERVE** option causes the device driver to prohibit any reservation commands during the open of the device. This facilitates the sharing of the device by multiple hosts. The device driver should require the caller to have appropriate authority to request the **SC\_NO\_RESERVE** option, because this request allows other hosts to modify data on the device. If a caller does this kind of request then the caller must ensure data integrity between multiple hosts. If the caller attempts to execute this system call without the proper authority, the device driver should return a value of -1, with the **errno** global variable set to a value of EPERM.

## Using the SC\_SINGLE Option

The **SC\_SINGLE** option causes the device driver to issue a normal open, but does not allow another caller to issue another open until the first caller has closed the device. This request gives the caller an exclusive open to the selected device. If this **openx** is attempted when the device is already open, a return value of -1 is passed, with the **errno** global variable set to a value of EBUSY.

Once successfully opened, the device is placed in Exclusive Access mode. If another caller tries to do any type of **open**, a return value of -1 is passed, with the **errno** global variable set to a value of EACCES.

The following table shows how the various combinations of *ext* options should be handled in the device driver.

EXT OPTIONS <b>openx</b> <i>ext</i> option	Device Driver Action	
	Open	Close
None	Normal	Normal
diag	No commands	No commands
diag + force	Removes a reservation; otherwise, no commands issued	No commands
diag + force + no_reserve	Removes a reservation; otherwise, no commands issued	No commands
diag + force + no_reserve + single	Removes a reservation; otherwise, no commands issued	No commands
diag + force + retain	Removes a reservation; otherwise, no commands issued	No commands
diag + force + retain + no_reserve	Removes a reservation; otherwise, no commands issued	No commands
diag + force + retain + no_reserve + single	Removes a reservation; otherwise, no commands issued	No commands
diag + force + retain + single	Removes a reservation; otherwise, no commands issued	No commands
diag + force + single	Removes a reservation; otherwise, no commands issued	No commands
diag + no_reserve	No commands	No commands
diag + retain	No commands	No commands
diag + retain + no_reserve	No commands	No commands
diag + retain + no_reserve + single	No commands	No commands
diag + retain + single	No commands	No commands
diag + single	No commands	No commands
diag + single + no_reserve	No commands	No commands

EXT OPTIONS openx ext option	Device Driver Action	
	Open	Close
force	Normal, except that a reservation is removed before any commands are issued	Normal
force + no_reserve	Normal, except that a reservation is removed before any commands are issued	Normal except no RELEASE
force + retain	Normal, except that a reservation is removed before any commands are issued	No RELEASE
force + retain + no_reserve	Normal, except that a reservation is removed before any commands are issued	No RELEASE
force + retain + no_reserve + single	Normal, except that a reservation is removed before any commands are issued	No RELEASE
force + retain + single	Normal, except that a reservation is removed before any commands are issued	No RELEASE
force + single	Normal, except that a reservation is removed before any commands are issued	Normal
force + single + no_reserve	Normal, except that a reservation is removed before any commands are issued	No RELEASE
no_reserve	No RESERVE	No RELEASE
retain	Normal	No RELEASE
retain + no_reserve	No RESERVE	No RELEASE
retain + single	Normal	No RELEASE
retain + single + no_reserve	Normal, except no RESERVE command issued.	No RELEASE
single	Normal	Normal
single + no_reserve	No RESERVE	No RELEASE

## Closing the Device

When a device driver is preparing to close a device through the adapter device driver, it must ensure that all transactions are complete. When the adapter device driver receives a **SCIOISTOP** ioctl operation and there are pending I/O requests, the ioctl operation does not return until all have completed. New requests received during this time are rejected from the adapter device driver's **ddstrategy** routine.

## Error Processing

It is the responsibility of the device driver to process check conditions and other returned errors properly. The adapter device driver only passes commands without otherwise processing them and is not responsible for device error recovery.

## Device Driver and Adapter Device Driver Interfaces

The device drivers can have both character (raw) and block special files in the `/dev` directory. The adapter device driver has only character (raw) special files in the `/dev` directory and has only the **ddconfig**, **ddopen**, **ddclose**, **dddump**, and **ddioctl** entry points available to operating system programs. The **ddread** and **ddwrite** entry points are not implemented.

Internally, the **devsw** table has entry points for the **ddconfig**, **ddopen**, **ddclose**, **dddump**, **ddioctl**, and **ddstrat** routines. The device drivers pass their commands to the adapter device driver by calling the adapter device driver **ddstrat** routine. (This routine is unavailable to other operating system programs due to the lack of a block-device special file.)

Access to the adapter device driver's **ddconfig**, **ddopen**, **ddclose**, **dddump**, **ddioctl**, and **ddstrat** entry points by the device drivers is performed through the kernel services provided. These include such services as **fp\_opendev**, **fp\_close**, **fp\_ioctl**, **devdump**, and **devstrat**.

## Performing Dumps

A adapter device driver must have a **dddump** entry point if it is used to access a system dump device. A device driver must have a **dddump** entry point if it drives a dump device. Examples of dump devices are disks and tapes.

**Note:** Adapter-device-driver writers should be aware that system services providing interrupt and timer services are unavailable for use in the dump routine. Kernel DMA services are assumed to be available for use by the dump routine. The adapter device driver should be designed to ignore extra **DUMPINIT** and **DUMPSTART** commands to the **dddump** entry point.

The **DUMPQUERY** option should return a minimum transfer size of 0 bytes, and a maximum transfer size equal to the maximum transfer size supported by the adapter device driver.

Calls to the adapter device driver **DUMPWRITE** option should use the *arg* parameter as a pointer to the **scsi\_buf** structure to be processed. Using this interface, a **write** command can be executed on a previously started (opened) target device. The *uiop* parameter is ignored by the adapter device driver during the **DUMPWRITE** command. Spanned, or consolidated, commands are not supported using the **DUMPWRITE** option. Gathered **write** commands are also not supported using the **DUMPWRITE** option. No queuing of **scsi\_buf** structures is supported during dump processing because the dump routine runs essentially as a subroutine call from the caller's dump routine. Control is returned when the entire **scsi\_buf** structure has been processed.

**Note:** Also, both adapter-device-driver and device-driver writers should be aware that any error occurring during the **DUMPWRITE** option is considered unsuccessful. Therefore, no error recovery is employed during the **DUMPWRITE**. Return values from the call to the **dddump** routine indicate the specific nature of the failure.

Successful completion of the selected operation is indicated by a 0 return value to the subroutine. Unsuccessful completion is indicated by a return code set to one of the following values for the **errno** global variable. The various **scsi\_buf** status fields, including the **b\_error** field, are not set by the adapter device driver at completion of the **DUMPWRITE** command. Error logging is, of necessity, not supported during the dump.

- An **errno** value of **EINVAL** indicates that a request that was not valid passed to the adapter device driver, such as to attempt a **DUMPSTART** command before successfully executing a **DUMPINIT** command.
- An **errno** value of **EIO** indicates that the adapter device driver was unable to complete the command due to a lack of required resources or an I/O error.
- An **errno** value of **ETIMEDOUT** indicates that the adapter did not respond with completion status before the passed command time-out value expired.

---

## SAM Adapter Device Driver ioctl Commands

Various ioctl operations must be performed for proper operation of the adapter device driver.

The ioctl operations described in the following topics are the primary set of ioctl commands that the adapter device driver must implement to support device drivers. Many of these ioctl operations are relevant to all SAM physical transport types, however, some operations are relevant only to a specific transport type. Other ioctl operations might be required in the adapter device driver to support, for example, system management facilities and diagnostics.

**Attention:** The adapter device driver ioctl operations can only be called from the process level. They cannot be run from a call on any more favored priority levels. Attempting to call them from a more favored priority level can result in the system ending abnormally.

### Related Information

“SAM Adapter ioctl Operations” on page 263

---

### Related Information

Logical File System Kernel Services.

scdisk SCSI Device Driver in *AIX Version 6.1 Technical Reference: Kernel and Subsystems Volume 2*.





---

## Chapter 14. Integrated Device Electronics (IDE) Subsystem

This overview describes the interface between an Integrated Device Electronics (IDE) device driver and an IDE adapter device driver. It is directed toward those designing and writing an IDE device driver that interfaces with an existing IDE adapter device driver. It is also meant for those designing and writing an IDE adapter device driver that interfaces with existing IDE device drivers.

The main topics covered in this overview are:

- Responsibilities of the IDE Adapter Device Driver
- Responsibilities of the IDE Device Driver
- Communication Between IDE Device Drivers and IDE Adapter Device Drivers

This section frequently refers to both an IDE device driver and an IDE adapter device driver. These two distinct device drivers work together in a layered approach to support attachment of a range of IDE devices. The IDE adapter device driver is the lower device driver of the pair, and the IDE device driver is the upper device driver.

---

### Responsibilities of the IDE Adapter Device Driver

The IDE adapter device driver is the software interface to the system hardware. This hardware includes the IDE bus hardware plus any other system I/O hardware required to run an I/O request. The IDE adapter device driver hides the details of the I/O hardware from the IDE device driver. The design of the software interface allows a user with limited knowledge of the system hardware to write the upper device driver.

The IDE adapter device driver manages the IDE bus, but not the IDE devices. It can send and receive IDE commands, but it cannot interpret the contents of the command. The lower driver also provides recovery and logging for errors related to the IDE bus and system I/O hardware. Management of the device specifics is left to the IDE device driver. The interface of the two drivers allows the upper driver to communicate with different IDE bus adapters without requiring special code paths for each adapter.

---

### Responsibilities of the IDE Device Driver

The IDE device driver provides the rest of the operating system with the software interface to a given IDE device or device class. The upper layer recognizes which IDE commands are required to control a particular IDE device or device class. The IDE device driver builds I/O requests containing device IDE commands and sends them to the IDE adapter device driver in the sequence needed to operate the device successfully. The IDE device driver cannot manage adapter resources. Specifics about the adapter and system hardware are left to the lower layer.

The IDE device driver also provides command retries and logging for errors related to the IDE device it controls.

The operating system provides several kernel services allowing the IDE device driver to communicate with IDE adapter device driver entry points without having the actual name or address of those entry points. See “Logical File System Kernel Services” on page 63 for more information.

---

### Communication Between IDE Device Drivers and IDE Adapter Device Drivers

The interface between the IDE device driver and the IDE adapter device driver is accessed through calls to the IDE adapter device driver **open**, **close**, **ioctl**, and **strategy** subroutines. I/O requests are queued to the IDE adapter device driver through calls to its **strategy** subroutine entry point.

Communication between the IDE device driver and the IDE adapter device driver for a particular I/O request uses the `ataide_buf` structure, which is passed to and from the `strategy` subroutine in the same way a standard driver uses a `struct buf` structure. The `ataide_buf.ata` structure represents the **ATA** or **ATAPI** command that the adapter driver must send to the specified IDE device. The `ataide_buf.status_validity` field in the `ataide_buf.ata` structure contains completion status returned to the IDE device driver.

---

## IDE Error Recovery

If an error, such as a check condition or hardware failure occurs, the transaction active during the error is returned with the `ataide_buf.bufstruct.b_error` field set to **EIO**. The IDE device driver will process the error by gathering hardware and software status. In many cases, the IDE device driver only needs to retry the unsuccessful operation.

The IDE adapter driver should never retry an IDE command on error after the command has successfully been given to the adapter. The consequences for the adapter driver retrying an IDE command at this point range from minimal to catastrophic, depending upon the type of device. Commands for certain devices cannot be retried immediately after a failure (for example, tapes and other sequential access devices). If such an error occurs, the failed command returns an appropriate error status with an `iodone` call to the IDE device driver for error recovery. Only the IDE device driver that originally issued the command knows if the command can be retried on the device. The IDE adapter driver must only retry commands that were never successfully transferred to the adapter. In this case, if retries are successful, the `ataide_buf` status should not reflect an error. However, the IDE adapter driver should perform error logging on the retried condition.

## Analyzing Returned Status

The following order of precedence should be followed by IDE device drivers when analyzing the returned status:

1. If the `ataide_buf.bufstruct.b_flags` field has the **B\_ERROR** flag set, then an error has occurred and the `ataide_buf.bufstruct.b_error` field contains a valid **errno** value.  
If the `b_error` field contains the **ENXIO** value, either the command needs to be restarted or it was canceled at the request of the IDE device driver.  
If the `b_error` field contains the **EIO** value, then either one or no flag is set in the `ataide_buf.status_validity` field. If a flag is set, an error in either the `ata.status` or `ata.errval` field is the cause.
2. If the `ataide_buf.bufstruct.b_flags` field does not have the **B\_ERROR** flag set, then no error is being reported. However, the IDE device driver should examine the `b_resid` field to check for cases where less data was transferred than expected. For some IDE commands, this occurrence might not represent an error. The IDE device driver must determine if an error has occurred.  
There is a special case when `b_resid` will be nonzero. The DMA service routine might not be able to map all virtual to real memory pages for a single DMA transfer. This might occur when sending close to the maximum amount of data that the adapter driver supports. In this case, the adapter driver transfers as much of the data that can be mapped by the DMA service. The unmapped size is returned in the `b_resid` field, and the `status_validity` will have the **ATA\_IDE\_DMA\_NORES** bit set. The IDE device driver is expected to send the data represented by the `b_resid` field in a separate request.  
If a nonzero `b_resid` field does represent an error condition, recovering is the responsibility of the IDE device driver.

---

## A Typical IDE Driver Transaction Sequence

A simplified sequence of events for a transaction between an IDE device driver and an IDE adapter driver follows. In this sequence, routine names preceded by a `dd_` are part of the IDE device driver, while those preceded by an `eide_` are part of the IDE adapter driver.

1. The IDE device driver receives a call to its **dd\_strategy** routine; any required internal queuing occurs in this routine. The **dd\_strategy** entry point then triggers the operation by calling the **dd\_start** entry point. The **dd\_start** routine invokes the **eide\_strategy** entry point by calling the **devstrat** kernel service with the relevant **ataide\_buf** structure as a parameter.
2. The **eide\_strategy** entry point initially checks the **ataide\_buf** structure for validity. These checks include validating the **devno** field, matching the IDE device ID to internal tables for configuration purposes, and validating the request size.
3. The IDE adapter driver does not queue transactions. Only a single transaction is accepted per device (one master, one slave). If no transaction is currently active, the **eide\_strategy** routine immediately calls the **eide\_start** routine with the new transaction. If there is a current transaction for the same device, the new transaction is returned with an error indicated in the **ataide\_buf** structure. If there is a current transaction for the other device, the new transaction is queued to the inactive device.
4. At each interrupt, the **eide\_intr** interrupt handler verifies the current status. The IDE adapter driver fills in the **ataide\_buf** **status\_validity** field, updating the **ata.status** and **ata.errval** fields as required. The IDE adapter driver also fills in the **bufstruct.b\_resid** field with the number of bytes not transferred from the transaction. If all the data was transferred, the **b\_resid** field is set to a value of 0. When a transaction completes, the **eide\_intr** routine causes the **ataide\_buf** entry to be removed from the device queue and calls the **iodone** kernel service, passing the just dequeued **ataide\_buf** structure for the device as the parameter. The **eide\_start** routine is then called again to process the next transaction on the device queue. The **iodone** kernel service calls the IDE device driver **dd\_iodone** entry point, signaling the IDE device driver that the particular transaction has completed.
5. The IDE device driver **dd\_iodone** routine investigates the I/O completion codes in the **ataide\_buf** status entries and performs error recovery, if required. If the operation completed correctly, the IDE device driver dequeues the original buffer structures. It calls the **iodone** kernel service with the original buffer pointers to notify the originator of the request.

---

## IDE Device Driver Internal Commands

During initialization, error recovery, and open or close operations, IDE device drivers initiate some transactions not directly related to an operating system request. These transactions are called internal commands and are relatively simple to handle.

Internal commands differ from operating system-initiated transactions in several ways. The primary difference is that the IDE device driver is required to generate a **struct buf** that is not related to a specific request. Also, the actual IDE commands are typically more control oriented than data transfer related.

The only special requirement for commands is that the IDE device driver must have pinned the transfer data buffers. However, due to system hardware considerations, additional precautions must be taken for data transfers into system memory pages. The problem is that any system memory area with a DMA data operation in progress causes the entire memory page that contains it to become inaccessible.

As a result, an IDE device driver that initiates an internal command must have preallocated and pinned an area of some multiple of system page size. The driver must not place in this area any other data that it might need to access while I/O is being performed into or out of that page. Memory pages allocated must be avoided by the device driver from the moment the transaction is passed to the adapter driver until the device driver **iodone** routine is called for the transaction.

---

## Execution of I/O Requests

During normal processing, many transactions are queued in the IDE device driver. As the IDE device driver processes these transactions and passes them to the IDE adapter driver, the IDE device driver moves them to the in-process queue. When the IDE adapter device driver returns through the **iodone** service with one of these transactions, the IDE device driver either recovers any errors on the transaction or returns using the **iodone** kernel service to the calling level.

The IDE device driver can send only one **ataide\_buf** structure per call to the IDE adapter driver. Thus, the **ataide\_buf.bufstruct.av\_forw** pointer must be null when given to the IDE adapter driver, which indicates that this is the only request. The IDE adapter driver does not support queuing multiple requests to the same device.

## Spanned (Consolidated) Commands

Some kernel operations might be composed of sequential operations to a device. For example, if consecutive blocks are written to disk, blocks might or might not be in physically consecutive buffer pool blocks.

To enhance IDE bus performance, the IDE device driver should consolidate multiple queued requests when possible into a single IDE command. To allow the IDE adapter driver the ability to handle the scatter and gather operations required, the **ataide\_buf.bp** should always point to the first **buf** structure entry for the spanned transaction. A null-terminated list of additional **struct buf** entries should be chained from the first field through the **buf.av\_forw** field to give the IDE adapter driver enough information to perform the DMA scatter and gather operations required. This information must include at least the buffer's starting address, length, and cross-memory descriptor.

The spanned requests should always be for requests in either the read or write direction but not both, because the IDE adapter driver must be given a single IDE command to handle the requests. The spanned request should always consist of complete I/O requests (including the additional **struct buf** entries). The IDE device driver should not attempt to use partial requests to reach the maximum transfer size.

The maximum transfer size is actually adapter-dependent. The **IOCINFO** ioctl operation can be used to discover the IDE adapter driver's maximum allowable transfer size. If a transfer size larger than the supported maximum is attempted, the IDE adapter driver returns a value of **EINVAL** in the **ataide\_buf.bufstruct.b\_error** field.

Due to system hardware requirements, the IDE device driver must consolidate only commands that are memory page-aligned at both their starting and ending addresses. Specifically, this applies to the consolidation of memory buffers. The ending address of the first buffer and the starting address of all subsequent buffers should be memory page-aligned. However, the starting address of the first memory buffer and the ending address of the last do not need to be aligned.

The purpose of consolidating transactions is to decrease the number of IDE commands and bus phases required to perform the required operation. The time required to maintain the simple chain of **buf** structure entries is significantly less than the overhead of multiple (even two) IDE bus transactions.

## Fragmented Commands

Single I/O requests larger than the maximum transfer size must be divided into smaller requests by the IDE device driver. For calls to an IDE device driver's character I/O (read/write) entry points, the **uphysio** kernel service can be used to break up these requests. For a fragmented command such as this, the **ataide\_buf.bp** field should be NULL so that the IDE adapter driver uses only the information in the **ataide\_buf** structure to prepare for the DMA operation.

---

## ataide\_buf Structure

The **ataide\_buf** structure is used for communication between the IDE device driver and the IDE adapter driver during an initiator I/O request. This structure is passed to and from the **strategy** routine in the same way a standard driver uses a **struct buf** structure.

## Fields in the `ataide_buf` Structure

The `ataide_buf` structure contains certain fields used to pass an IDE command and associated parameters to the IDE adapter driver. Other fields within this structure are used to pass returned status back to the IDE device driver. The `ataide_buf` structure is defined in the `/usr/include/sys/ide.h` file.

Fields in the `ataide_buf` structure are used as follows:

1. Reserved fields should be set to a value of 0, except where noted.
2. The `bufstruct` field contains a copy of the standard `buf` buffer structure that documents the I/O request. Included in this structure, for example, are the buffer address, byte count, and transfer direction. The `b_work` field in the `buf` structure is reserved for use by the IDE adapter driver. The current definition of the `buf` structure is in the `/usr/include/sys/buf.h` include file.
3. The `bp` field points to the original buffer structure received by the IDE device driver from the caller, if any. This can be a chain of entries in the case of spanned transfers (IDE commands that transfer data from or to more than one system-memory buffer). A null pointer indicates a nonspanned transfer. The null value specifically tells the IDE adapter driver all the information needed to perform the DMA data transfer is contained in the `bufstruct` fields of the `ataide_buf` structure. If the `bp` field is set to a non-null value, the `ataide_buf.sg_ptr` field must have a value of null, or else the operation is not allowed.
4. The `ata` field, defined as an `ata_cmd` structure, contains the IDE command (ATA or ATAPI), status, error indicator, and a flag variable:

- a. The `flags` field contains the following bit flags:

### **ATA\_CHS\_MODE**

Execute the command in cylinder head sector mode.

### **ATA\_LBA\_MODE**

Execute the command in logical block addressing mode.

### **ATA\_BUS\_RESET**

Reset the ATA bus, ignore the current command.

- b. The `command` field is the IDE ATA command opcode. For ATAPI packet commands, this field must be set to **ATA\_ATAPI\_PACKET\_COMMAND** (0xA0).
- c. The `device` field is the IDE indicator for either the master (0) or slave (1) IDE device.
- d. The `sector_cnt_cmd` field is the number of sectors affected by the command. A value of zero usually indicates 256 sectors.
- e. The `startblk` field is the starting LBA or CHS sector.
- f. The `feature` field is the ATA feature register.
- g. The `status` field is a return parameter indicating the ending status for the command. This field is updated by the IDE adapter driver upon completion of a command.
- h. The `errval` field is the error type indicator when the `ATA_ERROR` bit is set in the `status` field. This field has slightly different interpretations for ATA and ATAPI commands.
- i. The `sector_cnt_ret` field is the number of sectors not processed by the device.
- j. The `endblk` field is the completion LBA or CHS sector.
- k. The `atapi` field is defined as an `atapi_command` structure, which contains the IDE **ATAPI** command. The 12 or 16 bytes of a single ATAPI command are stored in consecutive bytes, with the opcode identified individually. The `atapi_command` structure contains the following fields:
  - l. The `length` field is the number of bytes in the actual ATAPI command. This is normally 12 or 16 (decimal).
  - m. The `packet.op_code` field specifies the standard ATAPI opcode for this command.
  - n. The `packet.bytes` field contains the remaining command-unique bytes of the ATAPI command block. The actual number of bytes depends on the value in the `length` field.

- o. The `ataide_buf.bufstruct.b_un.b_addr` field normally contains the starting system-buffer address and is ignored and can be altered by the IDE adapter driver when the `ataide_buf` is returned. The `ataide_buf.bufstruct.b_bcount` field should be set by the caller to the total transfer length for the data.
- p. The `timeout_value` field specifies the time-out limit (in seconds) to be used for completion of this command. A time-out value of 0 means no time-out is applied to this I/O request.
- q. The `status_validity` field contains an output parameter that can have the following bit flags as a value:

**ATA\_IDE\_STATUS**

The `ata.status` field is valid.

**ATA\_ERROR\_VALID**

The `ata.errval` field contains a valid error indicator.

**ATA\_CMD\_TIMEOUT**

The IDE adapter driver caused the command to time out.

**ATA\_NO\_DEVICE\_RESPONSE**

The IDE device is not ready.

**ATA\_IDE\_DMA\_ERROR**

The IDE adapter driver encountered a DMA error.

**ATA\_IDE\_DMA\_NORES**

The IDE adapter driver was not able to transfer entire request. The `bufstruct.b_resid` contains the count not transferred.

If an error is detected while an IDE command is being processed, and the error prevented the IDE command from actually being sent to the IDE bus by the adapter, then the error should be processed or recovered, or both, by the IDE adapter driver.

If it is recovered successfully by the IDE adapter driver, the error is logged, as appropriate, but is not reflected in the `ata.errval` byte. If the error cannot be recovered by the IDE adapter driver, the appropriate `ata.errval` bit is set and the `ataide_buf` structure is returned to the IDE device driver for further processing.

If an error is detected after the command was actually sent to the IDE device, then the adapter driver will return the command to the device driver for error processing and possible retries.

For error logging, the IDE adapter driver logs IDE bus- and adapter-related conditions, where as the IDE device driver logs IDE device-related errors. In the following description, a capital letter "A" after the error name indicates that the IDE adapter driver handles error logging. A capital letter "H" indicates that the IDE device driver handles error logging.

Some of the following error conditions indicate an IDE device failure. Others are IDE bus- or adapter-related.

**ATA\_IDE\_DMA\_ERROR (A)**

The system I/O bus generated or detected an error during a DMA transfer.

**ATA\_ERROR\_VALID (H)**

The request sent to the device failed.

**ATA\_CMD\_TIMEOUT (A) (H)**

The command timed out before completion.

**ATA\_NO\_DEVICE\_RESPONSE (A)**

The target device did not respond.

## ATA\_IDE\_BUS\_RESET (A)

The adapter indicated the IDE bus reset failed.

---

## Other IDE Design Considerations

The following topics cover design considerations of IDE device and adapter drivers:

- IDE Device Driver Tasks
- Closing the IDE Device
- IDE Error Processing
- Device Driver and adapter driver Interfaces
- Performing IDE Dumps

## IDE Device Driver Tasks

IDE device drivers are responsible for the following actions:

- Interfacing with block I/O and logical volume device driver code in the operating system.
- Translating I/O requests from the operating system into IDE commands suitable for the particular IDE device. These commands are then given to the IDE adapter driver for execution.
- Issuing any and all IDE commands to the attached device. The IDE adapter driver sends no IDE commands except those it is directed to send by the calling IDE device driver.

## Closing the IDE Device

When an IDE device driver is preparing to close a device through the IDE adapter driver, it must ensure that all transactions are complete. When the IDE adapter driver receives an **IDEIOSTOP** ioctl operation and there are pending I/O requests, the ioctl operation does not return until all have completed. New requests received during this time are rejected from the adapter driver's **ddstrategy** routine.

## IDE Error Processing

It is the responsibility of the IDE device driver to properly process IDE check conditions and other returned device errors. The IDE adapter driver only passes IDE commands to the device without otherwise processing them and is not responsible for device error recovery.

## Device Driver and Adapter Driver Interfaces

The IDE device drivers can have both character (raw) and block special files in the **/dev** directory. The IDE adapter driver has only character (raw) special files in the **/dev** directory and has only the **ddconfig**, **ddopen**, **ddclose**, **dddump**, and **ddioctl** entry points available to operating system programs. The **ddread** and **ddwrite** entry points are not implemented.

Internally, the devsw table has entry points for the **ddconfig**, **ddopen**, **ddclose**, **dddump**, **ddioctl**, and **ddstrategy** routines. The IDE device drivers pass their IDE commands to the IDE adapter driver by calling the IDE adapter driver **ddstrategy** routine. (This routine is unavailable to other operating system programs due to the lack of a block-device special file.)

Access to the IDE adapter driver's **ddconfig**, **ddopen**, **ddclose**, **dddump**, **ddioctl**, and **ddstrategy** entry points by the IDE device drivers is performed through the kernel services provided. These include such kernel services as **fp\_opendev**, **fp\_close**, **fp\_ioctl**, **devdump**, and **devstrat**.

## Performing IDE Dumps

An IDE adapter driver must have a **dddump** entry point if it is used to access a system dump device. An IDE device driver must have a **dddump** entry point if it drives a dump device. Examples of dump devices are disks and tapes.

**Note:** IDE adapter driver writers should be aware that system services providing interrupt and timer services are unavailable for use while executing the **dump** routine. Kernel DMA services are assumed to be available for use by the **dump** routine. The IDE adapter driver should be designed to ignore extra **DUMPINIT** and **DUMPSTART** commands to the **dddump** entry point while processing the **dump** routine.

The **DUMPQUERY** option should return a minimum transfer size of 0 bytes, and a maximum transfer size equal to the maximum transfer size supported by the IDE adapter driver.

Calls to the IDE adapter driver **DUMPWRITE** option should use the **arg** parameter as a pointer to the **ataide\_buf** structure to be processed. Using this interface, an IDE write command can be executed on a previously started (opened) target device. The **uiop** parameter is ignored by the IDE adapter driver during the **DUMPWRITE** command. Spanned or consolidated commands are not supported using the **DUMPWRITE** option. Gathered write commands are also not supported using the **DUMPWRITE** option. No queuing of **ataide\_buf** structures is supported during dump processing because the **dump** routine runs essentially as a subroutine call from the caller's dump routine. Control is returned when the entire **ataide\_buf** structure has been processed.

**Note:** No error recovery techniques are used during the **DUMPWRITE** option because *any* error occurring during **DUMPWRITE** is a real problem as the system is already unstable. Return values from the call to the **dddump** routine indicate the specific nature of the failure.

Successful completion of the selected operation is indicated by a 0 return value to the subroutine. Unsuccessful completion is indicated by a return code set to one of the following values for the **errno** global variable. The various **ataide\_buf** status fields, including the **b\_error** field, are not set by the IDE adapter driver at completion of the **DUMPWRITE** command. Error logging is, of necessity, not supported during the dump.

- An **errno** value of **EINVAL** indicates that an invalid request (unknown command or bad parameter) was passed to the IDE adapter driver, such as to attempt a **DUMPSTART** command before successfully executing a **DUMPINIT** command.
- An **errno** value of **EIO** indicates that the IDE adapter driver was unable to complete the command due to a lack of required resources or an I/O error.
- An **errno** value of **ETIMEDOUT** indicates that the adapter did not respond to a command that was put in its register before the passed command time-out value expired.

---

## Required IDE Adapter Driver ioctl Commands

Various ioctl operations must be performed for proper operation of the IDE adapter driver. The ioctl operations described here are the minimum set of commands the IDE adapter driver must implement to support IDE device drivers. Other operations might be required in the IDE adapter driver to support, for example, system management facilities. IDE device driver writers also need to understand these ioctl operations.

Every IDE adapter driver must support the **IOCINFO** ioctl operation. The structure to be returned to the caller is the **devinfo** structure, including the **ide** union definition for the IDE adapter found in the **/usr/include/sys/devinfo.h** file. The IDE device driver should request the **IOCINFO** ioctl operation (probably during its open routine) to get the maximum transfer size of the adapter.

**Note:** The IDE adapter driver ioctl operations can only be called from the process level. They cannot be executed from a call on any more favored priority levels. Attempting to call them from a more favored priority level can result in a system crash.



## ioctl Commands

The following **IDEIOSTART** and **IDEIOSTOP** operations must be sent by the IDE device driver (for the open and close routines, respectively) for each device. They cause the IDE adapter driver to allocate and initialize internal resources. The **IDEIORESET** operation is provided for clearing device hard errors.

Except where noted otherwise, the *arg* parameter for each of the ioctl operations described here must contain a long integer. In this field, the least significant byte is the IDE device ID value. (The upper three bytes are reserved and should be set to 0.) This provides the information required to allocate or deallocate resources and perform IDE bus operations for the ioctl operation requested.

The following information is provided on the various ioctl operations:

### IDEIOSTART

This operation allocates and initializes IDE device-dependent information local to the IDE adapter driver. Run this operation only on the first open of a device. Subsequent **IDEIOSTART** commands to the same device fail unless an intervening **IDEIOSTOP** command is issued.

For more information, see **IDEIOSTART (Start IDE) IDE Adapter Device Driver ioctl Operation** in *AIX Version 6.1 Technical Reference: Kernel and Subsystems Volume 1*.

### IDEIOSTOP

This operation deallocates resources local to the IDE adapter driver for this IDE device. This should be run on the last close of an IDE device. If an **IDEIOSTART** operation has not been previously issued, this command is unsuccessful.

For more information, see **IDEIOSTOP (Stop) IDE Adapter Device Driver ioctl Operation** in *AIX Version 6.1 Technical Reference: Kernel and Subsystems Volume 1*.

### IDEIORESET

This operation causes the IDE adapter driver to send an ATAPI device reset to the specified IDE device ID.

The IDE device driver should use this command only when directed to do a forced open. This occurs in for the situation when the device needs to be reset to clear an error condition.

**Note:** In normal system operation, this command should not be issued, as it would reset all devices connected to the controller. If an **IDEIOSTART** operation has not been previously issued, this command is unsuccessful.

### IDEIOINQU

This operation allows the caller to issue an **IDE device inquiry** command to a selected device.

For more information, see **IDEIOINQU (Inquiry) IDE Adapter Device Driver ioctl Operation** in *AIX Version 6.1 Technical Reference: Kernel and Subsystems Volume 1*.

### IDEIOSTUNIT

This operation allows the caller to issue an **IDE Start Unit** command to a selected IDE device. For the **IDEIOSTUNIT** operation, the *arg* parameter operation is the address of an **ide\_startunit** structure. This structure is defined in the `/usr/include/sys/ide.h` file.

For more information, see **IDEIOSTUNIT (Start Unit) IDE Adapter Device Driver ioctl Operation** in *AIX Version 6.1 Technical Reference: Kernel and Subsystems Volume 1*.

### IDEIOTUR

This operation allows the caller to issue an **IDE Test Unit Ready** command to a selected IDE device.

For more information, see **IDEIOTUR (Test Unit Ready) IDE Adapter Device Driver ioctl Operation** in *AIX Version 6.1 Technical Reference: Kernel and Subsystems Volume 1*.

## IDEIOREAD

This operation allows the caller to issue an **IDE device read** command to a selected device.

For more information, see IDEIOREAD (Read) IDE Adapter Device Driver ioctl Operation in *AIX Version 6.1 Technical Reference: Kernel and Subsystems Volume 1*.

## IDEIOIDENT

This operation allows the caller to issue an **IDE identify device** command to a selected device.

For more information, see IDEIOIDENT (Identify Device) IDE Adapter Device Driver ioctl Operation in *AIX Version 6.1 Technical Reference: Kernel and Subsystems Volume 1*.

---

## Related Information

Logical File System Kernel Services

## Technical References

The **ddconfig**, **ddopen**, **ddclose**, **dddump**, **ddioctl**, **ddread**, **ddstrategy**, **ddwrite** entry points in *AIX Version 6.1 Technical Reference: Kernel and Subsystems Volume 2*.

The **fp\_opendev**, **fp\_close**, **fp\_ioctl**, **devdump**, **devstrat** kernel services in *AIX Version 6.1 Technical Reference: Kernel and Subsystems Volume 2*.

IDE Adapter Device Driver, idecdrom IDE Device Driver, idedisk IDE Device Driver, IDEIOIDENT (Identify Device) IDE Adapter Device Driver ioctl Operation, IDEIOINQU (Inquiry) IDE Adapter Device Driver ioctl Operation, IDEIOREAD (Read) IDE Adapter Device Driver ioctl Operation, IDEIOSTART (Start IDE) Adapter Device Driver ioctl Operation, IDEIOSTOP (Stop) Device IDE Adapter Device Driver ioctl Operation, IDEIOSTUNIT (Start Unit) IDE Adapter Device Driver ioctl Operation, and IDEIOTUR (Test Unit Ready) IDE Adapter Device Driver ioctl Operation in *AIX Version 6.1 Technical Reference: Kernel and Subsystems Volume 2*.

---

## Chapter 15. Serial Direct Access Storage Device Subsystem

With *sequential* access to a storage device, such as with tape, a system enters and retrieves data based on the location of the data, and on a reference to information previously accessed. The closer the physical location of information on the storage device, the quicker the information can be processed.

In contrast, with *direct* access, entering and retrieving information depends only on the location of the data and not on a reference to data previously accessed. Because of this, access time for information on direct access storage devices (DASDs) is effectively independent of the location of the data.

Direct access storage devices (DASDs) include both fixed and removable storage devices. Typically, these devices are hard disks. A *fixed* storage device is any storage device defined during system configuration to be an integral part of the system DASD. If a fixed storage device is not available at some time during normal operation, the operating system detects an error.

A *removable* storage device is any storage device you define during system configuration to be an optional part of the system DASD. Removable storage devices can be removed from the system at any time during normal operation. As long as the device is logically unmounted before you remove it, the operating system does not detect an error.

The following types of devices are not considered DASD and are not supported by the logical volume manager (LVM):

- Diskettes
- CD-ROM (compact disk read-only memory)
- DVD-ROM (DVD read-only memory)
- WORM (write-once read-mostly)

---

### DASD Device Block Level Description

The DASD *device block* (or *sector*) level is the level at which a processing unit can request low-level operations on a device block address basis. Typical low-level operations for DASD are read-sector, write-sector, read-track, write-track, and format-track.

By using direct access storage, you can quickly retrieve information from random addresses as a stream of one or more blocks. Many DASDs perform best when the blocks to be retrieved are close in physical address to each other.

A DASD consists of a set of flat, circular rotating platters. Each platter has one or two sides on which data is stored. Platters are read by a set of nonrotating, but positionable, read or read/write heads that move together as a unit.

The following terms are used when discussing DASD device block operations:

<b>sector</b>	An addressable subdivision of a track used to record one block of a program or data. On a DASD, this is a contiguous, fixed-size block. Every sector of every DASD is exactly 512 bytes.
<b>track</b>	A circular path on the surface of a disk on which information is recorded and from which recorded information is read; a contiguous set of sectors. A track corresponds to the surface area of a single platter swept out by a single head while the head remains stationary.  A DASD contains at least 17 sectors per track. Otherwise, the number of sectors per track is not defined architecturally and is device-dependent. A typical DASD track can contain 17, 35, or 75 sectors.  A DASD can contain 1024 tracks. The number of tracks per DASD is not defined architecturally and is device-dependent.

**head** A head is a positionable entity that can read and write data from a given track located on one side of a platter. Usually a DASD has a small set of heads that move from track to track as a unit.

There must be at least 43 heads on a DASD. Otherwise, the number is not defined architecturally and is device-dependent. A typical DASD has 8 heads.

**cylinder** The tracks of a DASD that can be accessed without repositioning the heads. If a DASD has  $n$  number of vertically aligned heads, a cylinder has  $n$  number of vertically aligned tracks.

## Related Information

Programming in the Kernel Environment Overview

Understanding Physical Volumes and the Logical Volume Device Driver

Special Files Overview in *AIX Version 6.1 Files Reference*.

---

## Chapter 16. Debug Facilities

You can use the available procedures for debugging a device driver that is under development. The procedures you can use include the following types:

- The System Dump Facility, which results in a system restart, creates a snapshot of kernel data structures and heaps when a severe error occurs.
- The Live Dump Facility, which does not entail a system restart, creates small snapshots of memory relevant to a component.
- The Component Trace Facility captures information about a specific kernel component, kernel extension, or device driver.
- Error logging records device-specific hardware or software abnormalities.
- The Debug and Performance Tracing monitors entry and exit of device drivers and selectable system events.
- The Memory Overlay Detection System (MODS) helps detect memory overlay problems in the kernel, kernel extensions, and device drivers.

---

### System Dump Facility

Your system generates a system dump when a severe error occurs. System dumps can also be user-initiated by users with root user authority. A system dump creates a picture of your system's memory contents. System administrators and programmers can generate a dump and analyze its contents when debugging new applications.

System dumps can be assisted by firmware. Different from traditional system dumps that are generated before the partition is reinitialized, firmware-assisted system dumps take place when the partition is restarting. Firmware-assisted system dumps can be one of these types:

#### Selective memory dump

Selective memory dumps are triggered by or use AIX instances that must be dumped.

#### Full memory dump

The whole partition memory is dumped without any interaction with an AIX instance that is failing.

By default, the system generates a traditional system dump. After you explicitly enable a firmware-assisted system dump, it becomes the preferred system dump. However, if the configuration of the firmware-assisted system dump fails, the system generates a traditional system dump. A firmware-assisted system dump takes place under the following conditions:

- The platform supports firmware-assisted system dumps. AIX retrieves the property of firmware-assisted system dumps in the device tree to get the information.
- The memory size at system startup is equal to or greater than 1.5GB.
- You have not configured a traditional system dump.

A firmware-assisted dump cannot copy dump tables. Because the data is written on the next restart of the system, the dump tables, which are used to refer to the data, cannot be preserved.

RAS infrastructure components can be system-dump aware, allowing granular control of the amount of data that is dumped in a system dump by infrastructure components. Components that are system-dump aware can be excluded from a system dump to reduce the dump size. You can use the **dumpctrl** command to obtain information about which infrastructure components are registered for a system dump.

Use the **ras\_register** kernel service to make an infrastructure component dump-aware. Use the **RASCD\_SET\_SDMP\_ON** command that is passed to the **ras\_control** kernel service to make an infrastructure component system-dump aware. See “Callback Commands for System Dumps” on page 323.

If your system stops with an 888 number flashing in the operator panel display, the system has generated a dump and saved it to a dump device (the condition only occurs with traditional system dumps).

To generate a system dump see the following topics:

- Configure a Dump Device
- Start a System Dump
- Check the Status of a System Dump
- Copy a System Dump
- Increase the Size of a Dump Device
- A Sample Kernel Extension

Some of the error log and dump commands are delivered in an optionally installable package called **bos.sysmgt.serv\_aid**. System dump commands included in the **bos.sysmgt.serv\_aid** include the **sysdumpstart** command. See the Software Service Aids Package for more information.

## Configuring a Dump Device

When an unexpected system halt occurs, the system dump facility automatically copies selected areas of kernel data to the primary dump device. These areas include kernel segment 0, as well as other areas registered in the Master Dump Table by kernel modules or kernel extensions. An attempt is made to dump to a secondary dump device if it has been defined.

**Restriction:** You cannot use firmware-assisted system dumps for the secondary dump device.

If your system has less than 4 GB of memory, the default dump device is the **/dev/hd6** logical volume, which is a paging logical volume. In addition, if the configuration is for a firmware-assisted system dump, you cannot configure the **/dev/hd6** logical volume, or any paging space, as the dump device.

Beginning with AIX 6.1 with the 6100-01 Technology Level, you can configure an iSCSI device as the dump device for a firmware-assisted system dump.

**Restriction:** You can only use an iSCSI-dump logical volume that resides on the iSCSI boot disk.

Beginning with AIX 6.1 Technology Level 6100-04, you can perform remote dumps on thin servers. To perform a remote dump on a thin server, you must define the relative dump resource on the NIM master and allocate the related dump resource to the thin server. The dump resource appears as an iSCSI disk to the NIM client and can only be used to configure the primary dump device. You can only configure firmware-assisted system dumps on primary dump devices. See Defining a dump resource for more information.

If a dump occurs to paging space (traditional system dumps only), the system will automatically copy the dump when the system is rebooted. By default, the dump is copied to the **/var/adm/ras** directory in the root volume group. See the **sysdumpdev** command for details on how to control dump copying.

Beginning with AIX 6.1, you can only use compressed dumps. See the **sysdumpdev**, **uncompress** command, and **dumpuncompress** commands for more details.

The **dumpcheck** facility notifies you if your dump device needs to be larger, or if the file system containing the copy directory is too small. This notification appears in the system error log. If you need to increase the size of your dump device, see “Increasing the Size of a Dump Device” on page 323.

For maximum effectiveness, dumpcheck should be run when the system is most heavily loaded. At such times, the system dump is most likely to be at its maximum size. Also, even with dumpcheck watching the dump size, it may still happen that the dump won't fit on the dump device or in the copy directory at the time it happens. This could occur if there is a peak in system load right at dump time.

## Including Device Driver Data

To have your device driver data areas included in a system dump, you must register the data areas in the master dump table. Use the **dmp\_ctl** kernel service to add an entry to the master dump table or to delete an entry.

The syntax is as follows:

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/dump.h>

int dmp_ctl(op, data)
int op;
struct dmpctl_data *data;
```

Use the **ras\_register** kernel service and the **ras\_control** kernel service to register a component for a system dump, as shown in the following example:

```
/* Register a component as dump aware */
rv = ras_register(&Rascb, (char*)Compname, (ras_block_t)0, RAS_TYPE_OTHER,
    "sample component", RASF_DUMP_AWARE, sample_callback, NULL);
if (rv) return(KERROR2ERRNO(rv));

/* Make the component system dump aware. */
rv = ras_control(Rascb, RASCD_SET_SDMP_ON, 0, 0);
if (rv) return(KERROR2ERRNO(rv));
```

## Starting a System Dump

**Attention:** Do not start a system dump if the flashing 888 number shows in your operator panel display. This number indicates your system has already created a system dump and written the information to your primary dump device. If you start your own dump before copying the information in your dump device, your new dump will overwrite the existing information. For more information, see “Checking the Status of a System Dump” on page 319.

A user-initiated dump is different from a dump initiated by an unexpected system halt because the user can designate which dump device to use. When the system halts unexpectedly, a system dump is initiated automatically to the primary dump device.

You can start a system dump by using one of the methods listed below.

You have access to the **sysdumpstart** command and can start a dump using one of these methods:

- Using the Command Line
- Using SMIT
- Using the Reset Button
- Using Special Key Sequences

### Using the Command Line

Use the following steps to choose a dump device, specify the dump type, initiate the system dump, and determine the status of the system dump.

**Requirement:** You must have root user authority to start a dump using the **sysdumpstart** command.

1. Check which dump device is appropriate for your system (the primary or secondary device) by using the following **sysdumpdev** command:

```
sysdumpdev -l
```

This command lists the current dump devices. You can use the **sysdumpdev** command to change device assignments.

2. Start the system dump using the **sysdumpstart** command. The **sysdumpstart** command starts a system dump on the default primary dump device. Use the **-p** flag to start a system dump to the primary dump device. Use the **-s** flag to start a system dump to the secondary dump device. Use the **-t** flag to specify a traditional system dump (**-t traditional**). You cannot force a firmware-assisted system dump if the traditional system dump is configured.

```
sysdumpstart -p
```

**Restriction:** You cannot start a firmware-assisted system dumps on the secondary dump device.

3. If a code shows in the operator panel display, refer to “Checking the Status of a System Dump” on page 319. If the operator panel display is blank, the dump was not started. Try again using the Reset button.

## Using SMIT

Use the following SMIT commands to choose a dump device and start the system dump:

**Note:** You must have root user authority to start a dump using SMIT. SMIT uses the **sysdumpstart** command to start a system dump.

1. Check which dump device is appropriate for your system (the primary or secondary device) by using the following SMIT fast path command:

```
smit dump
```

2. Choose the **Show Current Dump Devices** option and write the available devices on notepaper.
3. Enter the following SMIT fast path command again:

```
smit dump
```

4. Choose either the primary (the first example option) or secondary (the second example option) dump device to hold your dump information:

**Start a Dump to the Primary Dump Device**

OR

**Start a Dump to the Secondary Dump Device**

Base your decision on the list of devices you made in step 2.

5. Choose the type of dump you want, if you choose the primary dump device to hold your dump information.

To start a system dump according to the current dump configuration, choose the following type:

**Start a System Dump to the Primary Dump Device**

To force a full memory dump, a selective memory dump, or a traditional system dump, choose the corresponding type:

**Start a Full Memory Dump to the Primary Dump Device**

**Start a Selective Memory Dump to the Primary Dump Device**

**Start a Traditional System Dump to the Primary Dump Device**

6. Click **Reset** to start a dump again if the dump was not started (the operator panel display is blank). See “Checking the Status of a System Dump” on page 319 if a value shows in the operator panel display.

To start a dump with the **Reset** button or a key sequence you must have the key switch, or mode switch, in the Service position, or have set the Always Allow System Dump value to true. To do this:

- a. Use the following SMIT fast path command:

```
smit dump
```



- b. Set the Always Allow System Dump value to true. This is essential on systems that do not have a mode switch.

## Using the Reset Button

In AIX 5.3 and subsequent releases, pressing the reset button always dumps to the primary dump device. This is also true for LPAR systems running AIX 5.2.

Start a system dump with the Reset button by doing the following (this procedure works for all system configurations and will work in circumstances where other methods for starting a dump will not work):

1. If your machine has a key mode switch, do one of the following:
  - Turn the key mode switch to the Service position.
  - Set Always Allow System Dump to true.
  - Press the Reset button.
2. If your machine does not have a key mode switch, set Always Allow System Dump to true and press the Reset button.

Your system writes the dump information to the primary dump device.

**Note:** The procedure for using the reset button can vary, depending upon your hardware configuration.

## Using Special Key Sequences

Start a system dump with special key sequences by doing the following:

1. Turn your machine's mode switch to the Service position, or set Always Allow System Dump to true.
2. Press the Ctrl-Alt 1 key sequence to write the dump information to the primary dump device, or press the Ctrl-Alt 2 key sequence to write the dump information to the secondary dump device..

You can start a system dump by this method *only* on the native keyboard.

### Tips:

- You can use the **Dump** option of the **Restart Partition** function to initiate an AIX stand-alone system dump on POWER5 logical partitions that are managed by a Hardware Management Console (HMC). For more information about initiating a system dump during the restart of logical partitions using HMC, go to the IBM Systems Hardware Information Center available at <http://publib.boulder.ibm.com/eserver/> and search for *restart AIX logical partitions*.
- You can initiate a system dump remotely using a modem or terminal server after enabling the AIX remote-restart facility using the **smitty rrbtty** fast path. However, the AIX remote restart facility does not work for a system (integrated serial) port on a POWER5 system. Instead, enable serial port snoop. For more information about how to enable serial port snoop, go to the IBM Systems Hardware Information Center and search for *enable serial port snoop*.
- While a logical partition is dumping, dump progress indicators (0c0, 0c2, 0c9, and so on ) will appear on the HMC and in the LCD display. For more information about the dump status codes, go to the IBM Systems Hardware Information Center and search for *dump progress indicators*.

## Checking the Status of a System Dump

When a system dump is taking place, status and completion codes are displayed in the operator panel display on the operator panel. When the dump is complete, a 0cx status code displays if the dump was user initiated, a flashing 888 displays if the dump was system initiated.

You can check whether the dump was successful, and if not, what caused the dump to fail. If a 0cx is displayed, see "Status Codes" on page 320 below.

**Note:** If the dump fails and upon reboot you see an error log entry with the label DSI\_PROC or ISI\_PROC, and the Detailed Data area shows an **EXVAL** of 000 0005, this is probably a paging space I/O error.

If the paging space (probably/**dev/hd6**) is the dump device or on the same hard drive as the dump device, your dump may have failed because of a problem with that hard drive. You should run diagnostics against that disk.

## Status Codes

Find your status code in the following list, and follow the instructions:

- 000** The kernel debugger is started. If there is an ASCII terminal attached to one of the native serial ports, enter `q` at the debugger prompt (`>`) on that terminal and then wait for flashing 888s to appear in the operator panel display. After the flashing 888 appears, go to “Checking the Status of a System Dump” on page 319.
- 0c0** The dump completed successfully. Go to “Copying a System Dump.”
- 0c1** An I/O error occurred during the dump. Go to “System Dump Facility” on page 315.
- 0c2** A user-requested dump is not finished. Wait at least 1 minute for the dump to complete and for the operator panel display value to change. If the operator panel display value changes, find the new value on this list. If the value does not change, then the dump did not complete because of an unexpected error.
- 0c4** The dump ran out of space . A partial dump was written to the dump device, but there is not enough space on the dump device to contain the entire dump. To prevent this problem from occurring again, you must increase the size of your dump media. Go to “Increase the Size of a Dump Device” on page 322.
- 0c5** The dump failed because of an internal error.
- 0c8** The dump device has been disabled. The current system configuration does not designate a device for the requested dump. Enter the **sysdumpdev** command to configure the dump device.
- 0c9** A dump started by the system did not complete. Wait at least 1 minute for the dump to complete and for the operator panel display value to change. If the operator panel display value changes, find the new value on the list. If the value does not change, then the dump did not complete because of an unexpected error.
- 0ca** A firmware-assisted system dump is not finished yet. System startup resumes after the dump completes.
- 0cb** A firmware-assisted selective memory dump is started. Wait at least 10 minutes for the dump to complete and for the operator-panel display value to change. If the operator-panel display value changes, find the new value on the list. If the value does not change, the dump did not complete because of an unexpected error.
- 0cc** An error occurred dumping to the primary device; the dump has switched over to the secondary device. Wait at least 1 minute for the dump to complete and for the three-digit display value to change. If the three-digit display value changes, find the new value on this list. If the value does not change, then the dump did not complete because of an unexpected error.
- c20** The kernel debugger exited without a request for a system dump. Enter the **quit dump** subcommand. Read the new three-digit value from the LED display.

## Copying a System Dump

Your dump device holds the information that a system dump generates, whether the information is generated by the system or a user. You can copy this information to tape and deliver the material to your service department for analysis.

**Note:** If you intend to use a tape to send a snap image to IBM for software support, the tape must be one of the following formats: **8 mm, 2.3 Gb** capacity, **8 mm, 5.0 Gb** capacity, or **4 mm, 4.0 Gb** capacity. Using other formats will prevent or delay software support from being able to examine the contents.

There are two procedures for copying a system dump, depending on whether you are using a dataless workstation or a non-dataless machine:

- Copying a System Dump on a Dataless or Diskless Workstation
- Copying a System Dump on a Non-Dataless or Non-Diskless Machine

### Copying a System Dump on a Dataless or Diskless Workstation

On a dataless or diskless workstation, the dump is automatically copied to the server when the workstation is rebooted after the dump. The dump remains available to the workstation.

Copy a system dump on a dataless or diskless workstation by performing the following tasks:

1. Reboot in Normal mode.

2. Locate the System Dump.
3. Copy the System Dump from the Server.

**Reboot in Normal mode:** To reboot in normal mode:

1. Switch off the power on your machine.
2. Turn the mode switch to the Normal position.
3. Switch on the power on your machine.

**Locate the System Dump:** To locate the dump:

1. Log on to the server.
2. Use the **lsnim** command to find the dump object for the workstation. For this example, the workstation's object name on the server is worker.

```
lsnim -l worker
```

The dump object is displayed on the line:

```
dump = dumpobject
```

3. Use the **lsnim** command again to determine the path of the object:

```
lsnim -l dumpobject
```

The path name displayed is the directory containing the dump. The name of the dump is the same as the object for the dataless or diskless workstation.

**Copy the System Dump from the Server:** You can copy the dump like any other file, or you can aggregate the dump to workstation data by using the **snap** command.

To copy the dump to tape, use the **tar** command:

```
tar -c
```

To copy to a tape other than the **/dev/rmt0** tape:

```
tar -cftapedevice
```

To copy the dump back from the external media (such as a tape drive), use the **tar** command. Enter the following to copy the dump from the **/dev/rmt0** tape:

```
tar -x
```

To copy the dump from any other media, enter:

```
tar -xftapedevice
```

To aggregate the dump to workstation data, run the **nim -o snap** command or select the snap operation in the **smit nim\_mac\_op** panel. Both methods remotely execute the **snap** command. The collected data is located in the same directory as the dump directory. You can save the collected data to a tape, external media, or any other media as previously described.

## Copying a System Dump on a Non-Dataless or Non-Diskless Machine

Copy a system dump on a non-dataless or a non-diskless machine by performing the following tasks:

1. Reboot Your Machine.
2. Copy the System Dump using one of the following methods:
  - Copy a System Dump after Rebooting in Normal Mode
  - Copy a System Dump after Booting from Maintenance Mode

**Reboot Your Machine:** Reboot in Normal mode using the following steps:

1. Turn off the power on your machine.
2. Turn the mode switch to the Normal position.

3. Turn on the power on your machine.

If your system shows the login prompt, go to “Copy a System Dump after Rebooting in Normal Mode.”

If your system stops with a number in the operator panel display instead of showing the login prompt, reboot your machine from Maintenance mode, then go to “Copy a System Dump after Booting from Maintenance Mode.”

**Copy a System Dump after Rebooting in Normal Mode:** After rebooting in Normal mode, copy a system dump by doing the following:

1. Log in to your system as root user.
2. Copy the system dump to tape using the following **snap** command, where the number sign (#) is the number of your available tape device (the most common is **/dev/rmt0**):

```
/usr/sbin/snap -gfkD -o /dev/rmt#
```

To find the correct number of the available tape device, enter the following **lsdev** command, and look for the tape device listed as Available:

```
lsdev -C -c tape -H
```

**Note:** If your dump went to a paging space logical volume, it has been copied to a directory in your root volume group, **/var/adm/ras**. For more information, see Configure a Dump Device and the **sysdumpdev** command. These dumps are still copied by the **snap** command. The **sysdumpdev -L** command lists the exact location of the dump.

3. To copy the dump back from the external media (such as a tape drive), use the **pax** command. Enter the following to copy the dump from **/dev/rmt0**:

```
pax -rf/dev/rmt0
```

To copy the dump from any other media, enter:

```
tar -xftapedevice
```

**Copy a System Dump after Booting from Maintenance Mode:**

**Note:** Use this procedure only if you cannot boot your machine in Normal mode.

1. After booting from Maintenance mode, copy a system dump or tape using the following **snap** command:

```
/usr/sbin/snap -gfkD -o /dev/rmt#
```

2. To copy the dump back from the external media (such as a tape drive), use the **tar** command. Enter the following to copy the dump from the **/dev/rmt0** tape:

```
tar -x
```

To copy the dump from any other media, enter:

```
tar -xftapedevice
```

## Increase the Size of a Dump Device

Refer to the following to determine the appropriate size for your dump logical volume and to increase the size of either a logical volume or a paging space logical volume.

- Determining the Size of a Dump Device
- Determining the Type of Logical Volume
- Increasing the Size of a Dump Device

## Determining the Size of a Dump Device

The size required for a dump is not a constant value because the system does not dump paging space; only data that resides in real memory can be dumped. Paging space logical volumes will generally hold the system dump. However, because an incomplete dump may not be usable, follow the procedure below to make sure that you have enough dump space.

When a system dump occurs, all of the kernel segment that resides in real memory is dumped (the kernel segment is segment 0). Memory resident user data (such as u-blocks) are also dumped.

The minimum size for the dump space can best be determined using the **sysdumpdev -e** command. This gives an estimated dump size taking into account the memory currently in use by the system. If dumps are being compressed, then the estimate shown is for the compressed size of the dump, not the original size. In general, compressed dump size estimates will be much higher than the actual size. This occurs because of the unpredictability of the compression algorithm's efficiency. You should still ensure your dump device is large enough to hold the estimated size in order to avoid losing dump data.

For example, enter:

```
sysdumpdev -e
```

If **sysdumpdev -e** returns the message, Estimated dump size in bytes: 9830400, then the dump device should be at least 9830400 bytes or 12MB (if you are using three 4MB partitions for the disk).

## Determining the Type of Logical Volume

1. Enter the **sysdumpdev** command to list the dump devices. The logical volume of the primary dump device might be **/dev/lg\_dumplv** or **/dev/hd7**.

**Note:** You can also determine the dump devices using SMIT. Select the **Show Current Dump Devices** option from the System Dump SMIT menu.

2. Determine your logical volume type by using SMIT. Enter the SMIT fast path **smit lvm** or **smitty lvm**. You will go directly to Logical Volumes. Select the **List all Logical Volumes by Volume Group** option. Find your dump volume in the list and note its Type (in the second column). For example, the logical volume type might be **sysdump** if the dump device is **lg\_dumplv**, or **paging** if the dump device is **hd6**.

## Increasing the Size of a Dump Device

If you have confirmed that your dump device is a paging space, refer to Changing or Removing a Paging Space in *Operating system and device management* for more information.

If you have confirmed that your dump device type is **sysdump**, refer to the **extendlv** command for more information.

## Callback Commands for System Dumps

A RAS infrastructure component can participate in a system dump. The component uses the **RASCD\_SET\_SDMP\_ON** command that are passed to the **ras\_control** kernel service, and then handles the appropriate commands in the callback routine. Upon the receipt of the callback commands (listed in the first column of the table), the callback routine issues the commands that has **\_SET** in its name (listed in the second column of the table) to perform the actions.

Callback commands	Commands that are used by callbacks to perform actions	Description
RASCD_SDMP_ON	RASCD_SET_SDMP_ON	Enables a system dump for the component.
RASCD_SDMP_OFF	RASCD_SET_SDMP_OFF	Disables a system dump for the component.

Callback commands	Commands that are used by callbacks to perform actions	Description
RASCD_SDMP_LVL	RASCD_SET_SDMP_LVL	<p>Sets the system dump detail level for this component.</p> <p>The <i>arg</i> parameter of the <b>ras_control</b> service must be set with a value ranging from 0 through 9. If the level is not set, the component is dumped at the default detail level (CD_LVL_NORMAL).</p>
	RASCD_SET_SDMP_CALLONRESTART	<p>Causes the dump to call the callback again if the dump is restarted.</p> <p>If this attribute is not set, the callback is called only once during a system dump to collect data, the location of the returned dump table is saved, and the table is used as if the dump must be restarted. This attribute is assumed if the callback returns an unlimited dump table or requests staging buffer space.</p>
	RASCD_SET_SDMP_STAGING	<p>Reserves private staging buffer space.</p> <p>The <i>arg</i> parameter must contain the number of bytes to reserve. The maximum value is 32 MB. The <b>ENOMEM_RAS_SDMP_STAGING</b> error code is returned if the requested storage is unavailable.</p> <p>You can specify a size of 0 to remove an existing staging buffer. If you need to switch from using a staging buffer to using the shared staging buffer, you must first remove the existing buffer. To remove the existing buffer, issue a RASCD_SET_SDMP_STAGING request with a size of 0, and then issue a RASCD_SET_SDMP_SHARED_STAGING request.</p>

Callback commands	Commands that are used by callbacks to perform actions	Description
	RASCD_SET_SDMP_SHARED_STAGING	<p>Reserves shared staging buffer space. The system has one shared buffer for the system dump. This buffer is as large as the largest shared buffer request.</p> <p>Do not use this buffer to communicate between dump components. The buffer must not contain data that must appear in the dump because the firmware-assisted dump records the data to be copied to the dump device at the next restart. Therefore, the memory that is dumped from the specified address is the last data that existed at the address.</p> <p>Use the buffer only for dump metadata such as dump tables or for keeping the state between the RASCD_SDMP_START and RASCD_SDMP_AGAIN invocations.</p> <p>The <i>arg</i> parameter must contain the number of bytes to reserve. The maximum value is 32 MB. The <b>ENOMEM_RAS_SDMP_STAGING</b> error code is returned if the requested storage is unavailable.</p> <p>You can specify a size of 0 to remove an existing shared staging-buffer requirement. If you need to switch from using the shared staging buffer to using your own staging buffer, you must first remove the existing requirement. To remove the existing requirement, issue a RASCD_SET_SDMP_SHARED_STAGING request with a size of 0, and then issue a RASCD_SET_SDMP_STAGING request.</p>
RASCD_SDMP_START	N/A	<p>Provides data for the system dump.</p> <p>This is equivalent to the DMPRTN_START call that is passed to functions that are registered with the <b>dmp_add</b> or <b>dmp_ctl</b> kernel service. The <i>arg</i> parameter is of the <b>sdmp_start_t</b> type. The callback must set its dump table address to the <i>sdmpst_table</i> pointer. If the callback returns a negative value, the component is not included in the dump.</p>
RASCD_SDMP_AGAIN	N/A	<p>Provides more data for the unlimited dump table.</p> <p>This is equivalent to the DMPRTN_AGAIN call that is passed to functions that are registered with the <b>dmp_add</b> or <b>dmp_ctl</b> kernel service. The <i>arg</i> parameter is of the <b>sdmp_start_t</b> type. If the callback returns a negative value, no more data for that component is dumped. The component will not get an RASCD_SDMP_FINISHED call.</p>

Callback commands	Commands that are used by callbacks to perform actions	Description
RASCD_SDMP_FINISHED	N/A	Indicates that the system dump is completed.  This is equivalent to the DMPRTN_DONE call that is passed to functions that are registered with the <b>dmp_add</b> or <b>dmp_ctl</b> kernel service. The RASCD_SDMP_FINISHED call is not issued if a prior RASCD_SDMP_START or RASCD_SDMP_AGAIN call returned a negative return value.
RASCD_SDMP_ESTIMATE	N/A	Provides an estimate of how much data will be dumped.  The <i>arg</i> parameter is of the <b>sysdump_estimate_t</b> type. The value that is returned in the <i>se_value</i> parameter, is the same as that for the DMPRTN_ESTIMATE invocation in AIX 5.3. The <b>sysdump_estimate_t</b> structure contains the detail level at which to estimate.
	RASCD_SET_SDMP_SERIALIO	Enables serialized I/O during dump time. The need for the flag is device-specific. Only the developer of the device can determine whether this flag needs to be set. Use the flag only for the devices that can be on the dump I/O path. Serializing I/O during dump time might degrade dump performance. By default, without the flag, I/O can occur in parallel with function calls of component-dump tables. The action must be performed before the <b>ras_customize</b> routine is called, or an error is returned. <b>Note:</b> Beginning with AIX 6.1 with the 6100-02 Technology Level, the <b>ras_control</b> kernel service supports the RASCD_SET_SDMP_SERIALIO action flag.

---

## Live Dump Facility

Live dumps are small dumps that do not require a system restart. Live dumps replace system dumps when your system is running.

Only the components that are registered for live dumps are dumped. Use the **dumpctrl** command to obtain information about which components are registered for live dumps.

**Important:** The term component in this chapter refers to a component that is specified using the RAS infrastructure (created with the **ras\_register** kernel service). Only infrastructure components can be included in a live dump. See 1 on page 328. The **dmp\_ctl** and **dmp\_add** kernel services only apply to system dumps.

Live dumps can be initiated by software programs or by users with root user authority. Software programs use live dumps as part of recovery actions, or when the runtime error-checking value for the error disposition is ERROR\_LIVE\_DUMP. See “Initiating Live Dumps from Software Programs” on page 327 and “Sample Kernel Extension” on page 330. If you have root user authority, you can initiate live dumps when a subsystem does not respond or behaves erroneously. For more information about how to initiate and manage live dumps, see the **livedumpstart** and **dumpctrl** commands.



Unlike system dumps, which are written to a dedicated dump device, live dumps are written to the file system. When you install the operating system, a file system is created to contain live dumps. By default, live dumps are placed in the `/var/adm/ras/livedump` directory. You can change the directory using the `dumpctrl` command.

In AIX Version 6.1, only serialized live dumps are available. A serialized live dump causes a system to be frozen or suspended, when data is being dumped. When the system is frozen, the data is copied into the pinned kernel memory. The data is written to the file system only after the system is unfrozen. A component participating in a live dump must have a callback routine to handle the commands that are passed to the `ras_control` kernel service. See “Callback Commands for Live Dumps” on page 328 for details.

The default live-dump heap size is 64 MB or 1/64 of the size of physical memory, whichever is less. You can change the heap size using the `dumpctrl` command.

Duplicate live dumps that reoccur rapidly are eliminated to prevent system overload and to save file system space. Eliminating duplicate dumps requires periodic (once every 5 minutes) scans of the live dump repository through a cron job.

Each live dump has a data priority. A live dump of **info** priority is for informational purposes, and a live dump of **critical** priority is used to debug a problem. The size of a serialized live dump can be limited by the dump detail level. See “Live Dump Detail Levels” on page 330.

Data structures that are only related to live dumps are listed in the `/usr/include/sys/livedump.h` file.

You can disable all live dumps using the `dumpctrl ldmpoff` command.

## Initiating Live Dumps from Software Programs

A live dump can be initiated from software programs by the kernel or by a kernel extension. Components to be included in the dump must have been registered with the kernel, using the `ras_register` kernel service. The components must have indicated that they handle live dumps using the `RASCD_SET_LDMP_ON ras_control` service. See 1 on page 328.

To perform a live dump from software programs, follow these steps:

1. Initialize an `ldmp_parms_t` item using the `ldmp_setupparms` kernel service. This step sets up the data structure, filling in all default values including the eye catcher and version fields.
2. Specify infrastructure components using the `dmp_compspec` kernel service, and specify pseudo components using the pseudo components functions.
3. Take the live dump using the `livedump` kernel service.

### Pseudo components

A dump pseudo component is a service routine that is used to dump data that is not associated with a component. Such pseudo components are strictly used within a dump. The following pseudo components are provided.

<code>dmp_eaddr</code>	Dumps memory by effective addresses.
<code>dmp_context</code>	Dumps the kernel context.
<code>dmp_tid</code>	Dumps a thread.
<code>dmp_pid</code>	Dumps a process.
<code>dmp_errbuf</code>	Dumps the error logging buffer of the kernel.
<code>dmp_mtrc</code>	Dumps entries from the lightweight memory trace buffers.
<code>dmp_systrace</code>	Dumps entries from the system trace buffers.
<code>dmp_ct</code>	Dumps component trace entries.

## Examples

The following examples show how to register and include a component for live dumps.

1. In the following example, a component is registered for live dumps with the **ras\_register** and **ras\_control** kernel services:

```
/* Register a component as dump aware */
rv = ras_register(&Rascb, (char*)Compname, (ras_block_t)0, RAS_TYPE_OTHER,
    "sample component", RASF_DUMP_AWARE, sample_callback, NULL);
if (rv) return(KERROR2ERRNO(rv));

/* Make the component live dump aware. */
rv = ras_control(Rascb, RASCD_SET_LDMP_ON, 0, 0);
if (rv) return(KERROR2ERRNO(rv));
```

2. In the following example, a component is included in a live dump from software programs:

```
{
    ...
    ldmp_parms_t parms;
    extern ras_block_t Rascb; /* The ras_block_t from above */

    /* Setup the live dump parms structure. */
    if (ldmp_setupparms(&parm)) {
        /* serious error */
        ...
    }

    /* Each live dump must have a symptom. */
    parms.ldp_symptom = "sample dump";

    /* Include sample_comp in this dump as the failing component. */
    if (dmp_compspec(DCF_FAILING|DCF_BYCB, Rascb, &parm, NULL, NULL)) {
        /* error */
        ...
    }

    /* Add other components and/or pseudo components. */
    ...

    /* Take the dump. */
    if (livedump(&parm)) {
        /* error */
        ...
    }

    ...
}
```

## Callback Commands for Live Dumps

A component participating in a live dump must have a callback routine to handle the following commands that are passed to the **ras\_control** kernel service. Upon the receipt of the callback commands (listed in the first column in the table), the callback routine issues the commands that has **\_SET** in its name (listed in the second column in the table) to perform the actions.

Callback commands	Commands that are used by callbacks to perform actions	Description
RASCD_LDMP_ON	RASCD_SET_LDMP_ON	Enables a live dump for the component.
RASCD_LDMP_OFF	RASCD_SET_LDMP_OFF	Disables a live dump for the component.

Callback commands	Commands that are used by callbacks to perform actions	Description
RASCD_LDMP_LVL	RASCD_SET_LDMP_LVL	Sets the live dump level of the component.  The <i>arg</i> parameter of the <b>ras_control</b> service must be set with a value ranging from 0 through 9. If the level is not set, the component is dumped at the default detail level (CD_LVL_NORMAL).
RASCD_LDMP_PREPARE	N/A	Prepares to take a live dump.  The callback receives this call when it has been asked to participate in a live dump. The callback uses the <b>dmp_compspec</b> kernel service to specify other components to include in the dump, if necessary. The callback can also specify pseudo components.
RASCD_LDMP_START	N/A	Dumps data.  The callback stores its dump table address in the <b>ldmpst_table</b> field of the <b>ldmp_start_t</b> data item that is received as an argument.  When performing a serialized dump, the callback can use only the services listed in “Kernel Services for a Serialized Dump” on page 330.
RASCD_LDMP_AGAIN	N/A	Provides more data for the unlimited dump table (the <b>cdt_nn_u</b> type).  The return code is similar to that of the RASCD_LDMP_START command, except that if the return value is less than 0, no further data is dumped for the component. Data that was dumped by previous RASCD_LDMP_START and RASCD_LDMP_AGAIN calls appears in the dump.  When performing a serialized dump, the callback can use only the services listed in “Kernel Services for a Serialized Dump” on page 330.
RASCD_LDMP_FINISHED	N/A	Indicates that the live dump is completed.  When performing a serialized dump, the callback can use only the services listed in “Kernel Services for a Serialized Dump” on page 330.
RASCD_DMP_PASS_THROUGH	N/A	Passes arbitrary text data to the callback.  This command applies to the entire dump domain (there is only one pass-through for the domain containing live and system dumps). You can pass data to a component RASCD_DMP_PASS_THROUGH handler using the <b>dumpctrl</b> command. For example, use the <code>dumpctrl -l foo "pass through text"</code> command to the RASCD_DMP_PASS_THROUGH handler for the component with the alias of <i>foo</i> .
RASCD_LDMP_ESTIMATE	N/A	Provides an estimate of how much data will be dumped. The RASCD_LDMP_ESTIMATE call is similar to the RASCD_LDMP_PREPARE call.

## Kernel Services for a Serialized Dump

When data is provided for a serialized live dump, only the following services can be used. The services only apply to the RASCD\_LDMP\_START, RASCD\_LDMP\_AGAIN, and RASCD\_LDMP\_FINISHED calls.

- **ldmp\_bufest**, **ldmp\_timeleft**, **ldmp\_xmalloc**, **ldmp\_xmfree**, and **ldmp\_errstr**
- **vm\_att**, **vm\_det**, and **vm\_vmid**
- **lqra** and **lra**
- **raschk\_safe\_read**
- **disable\_lock**, **unlock\_enable**, **simple\_lock**, **simple\_lock\_try**, and **simple\_unlock**
- **i\_disable** and **i\_enable**
- Lightweight memory trace
- Component Trace
- **sprintf** and **sscanf**
- **printf** (debug only)
- The pinned string functions: **atoi**, **bcmp**, **memccpy**, **memchr**, **memcmp**, **memset**, **bzero**, **bcopy**, **memcpy**, **memmove**, **strcat**, **strchr**, **strcmp**, **strcpy**, **strcspn**, **strlen**, **strncat**, **strncmp**, **strncpy**, **strpbrk**, **strrchr**, **strspn**, **strstr**, **strtok\_r**, and **strtok**

A component can specify any data to be dumped, however, in a serialized dump, only memory-resident data is dumped. For a system dump, each data area in the dump has an associated bit map that indicates whether the data is in the dump, or the data cannot be included in the dump because it is not memory resident. While participating in a serialized live dump, a component must not directly refer to any storage that is not memory resident. To ensure safe access to data, use the **raschk\_safe\_read** kernel service with the **RAS\_SR\_NOPAGEIN** flag.

Do not use the system trace. If a system trace buffer fills, entries are lost until the system is unfrozen. Use the lightweight memory trace and component trace.

## Live Dump Detail Levels

In a serialized live dump, the amount of information that is dumped for a component is limited according to the dump detail level. Unless otherwise specified, the dump detail level of a component is **CD\_LVL\_NORMAL**. The detail level ranges from **CD\_LEVEL\_0** through **CD\_LEVEL\_9**. Three levels are used frequently: **CD\_LVL\_MINIMAL** (**CD\_LEVEL\_1**), **CD\_LVL\_NORMAL** (**CD\_LEVEL\_3**), and **CD\_LVL\_DETAIL** (**CD\_LEVEL\_7**). A component can query the value using the **rasrb\_ldmp\_level(rasb)** service, and can set the detail level with the **RASCD\_SET\_LDMP\_LVL** command in the **ras\_control** kernel service.

The following table shows the data limits for a component. If the component exceeds the limit, its data is truncated and only the data entries before the one that causes the limit to be exceeded are dumped.

Live dump detail level	Maximum size
< <b>CD_LVL_NORMAL</b>	2 MB
≥ <b>CD_LVL_NORMAL</b> & < <b>CD_LVL_DETAIL</b>	4 MB
≥ <b>CD_LVL_DETAIL</b> & < <b>CD_LEVEL_9</b>	8 MB
<b>CD_LEVEL_9</b>	unlimited

## Sample Kernel Extension

The following sample is a kernel extension that takes a live dump and a system dump. The **sample\_callback** function takes a live dump and a system dump using the commands that are sent by the system and passed to the **ras\_control** kernel service. The sample only shows the handling of the dump commands. Normally, the callback must handle component trace and error checking commands.

```

/*
 * This sample creates a component, makes it dump-aware, and handles both live
 * and system dump.
 */
#include <sys/types.h>
#include <sys/syspest.h>
#include <sys/uio.h>
#include <sys/processor.h>
#include <sys/systemcfg.h>
#include <sys/malloc.h>
#include <sys/ras.h>
#include <sys/livedump.h>
#include <sys/kernnodefs.h>
#include <sys/eyec.h>
#include <sys/raschk.h>
#include <sys/param.h>

/* Component name and handle */
const char Compname[] = "sample_comp";
ras_block_t Rascb=NULL;

/*
 * The sample data dumped consists of a header plus an unlimited number
 * of chained control blocks.
 * The header is dumped first.
 * Because we are using an unlimited dump table, the callback then gets an
 * "AGAIN" call to dump the rest of the data, the control blocks.
 * We'll dump these 3 at a time until finished.
 * The staging area is used for the dump table.
 */
#define NENTRIES 3
/* Staging area size for dumping the header. */
#define SZ1 (sizeof(struct cdt_nn_head) + sizeof(struct cdt_entry_u))
/* Staging area size for 3 control blocks. */
#define SZ2 (sizeof(struct cdt_nn_head) + DMP_DUL_SIZE(NENTRIES))
/* Staging buffer size is the MAX of SZ1 and SZ2 above. */
size_t Sbufsz;
/*
 * To estimate the dump table space, we need SZ1 plus the unlimited entry
 * for how many control blocks we have at dump time.
 */
#define TBLESTSZ(n) (SZ1 + DMP_DUL_SIZE(n))

/*
 * This is the sample data.
 * It would normally be protected with a lock, however, that is not shown here.
 */
typedef struct sample_cb {
    eye_catch8b_t scb_eyec; /* must be EYEC_SCB */
    struct sample_cb *scb_next; /* list ptr, terminated with
        INVALID_SCB_PTR */
    long scb_flags;
    long scb_fld1;
} sample_cb_t;
#define EYEC_SCB __EYEC8('s','a','m','p','l','A','B','C')
#define INVALID_SCB_PTR INVALID_PTR(EYEC_SCB)

typedef struct sample_hdr {
    eye_catch8b_t sh_eyec; /* must be EYEC_SH */
    long sh_flags;
    sample_cb_t *sh_cbp; /* ptr to first cb */
    long sh_fld1;
} sample_hdr_t;
#define EYEC_SH __EYEC8('s','a','h','d','r','A','B','C')
sample_hdr_t Sample_hdr = {EYEC_SH, 0, INVALID_SCB_PTR, 0};

kernno_t sample_callback(ras_block_t cb, ras_cmd_t cmd, void *arg, void *priv);

```

```

int alloc_sample(int n);
void free_sample();
static sample_cb_t *get_scbp(sample_cb_t *addr, dumpid_t id, int ldmpflag);

/*
 * Entry point called when this kernel extension is loaded.
 *
 * Input:
 *   cmd - unused (typically 1=config, 2=unconfig)
 *   uiop - points to the uio structure.
 */
int
sampleext(int cmd, struct uio *uiop)
{
    kerrno_t rv = 0;
    int rc;

    /* cmd should be 1 or 2 */
    if (cmd == 2) {
        /* Unloading */
        if (Rascb) ras_unregister(Rascb);
        free_sample();
        return(0);
    }
    if (cmd != 1) return(EINVAL);

    /* Set up local variables. */
    Sbufsz = MAX(SZ1, SZ2);

    /* The extension is being loaded, set up the sample data. */
    rc = alloc_sample(NENTRIES);
    if (rc) return(rc);

    /* Register the component as dump aware */
    rv = ras_register(&Rascb,
        (char*)Compname,
        (ras_block_t)0,
        RAS_TYPE_OTHER,
        "sample component",
        RASF_DUMP_AWARE,
        sample_callback,
        NULL);
    if (rv) return(KERROR2ERRNO(rv));

    /* Make the component system and live dump aware. */
    rv = ras_control(Rascb, RASCD_SET_SDMP_ON, 0, 0);
    if (rv) return(KERROR2ERRNO(rv));
    rv = ras_control(Rascb, RASCD_SET_LDMP_ON, 0, 0);
    if (rv) return(KERROR2ERRNO(rv));

    /*
     * System dump staging buffer space must be set up before a
     * system dump occurs.
     * Staging buffer space for live dumps is set up by the callback at
     * live dump time.
     */
    rv = ras_control(Rascb, RASCD_SET_SDMP_STAGING,
        (void*)Sbufsz, 0);
    if (rv) return(KERROR2ERRNO(rv));

    /*
     * The component must be customized.
     * It uses the default level, CD_LVL_NORMAL.
     */
    rv = ras_customize(Rascb);
    if (rv) return(KERROR2ERRNO(rv));
}

```

```

return(0);
}

/*
 * Sample callback that is called for live and system dumps.
 *
 * The data to dump consists of a header and
 * control blocks. The data is dumped using an unlimited dump table.
 * The header is dumped first, followed by the control blocks, dumped 3 at a
 * time until all have been dumped.
 *
 * Input:
 * cb - Contains the component's ras_block_t
 * cmd - ras_control command
 * arg - command argument
 * priv - private data, unused
 */
kernno_t
sample_callback(ras_block_t cb, ras_cmd_t cmd, void *arg, void *priv)
{
    kernno_t rv = 0;
    sample_cb_t *cbp, **wkptr;

    switch(cmd) {
        /* Live dump */
        case RASCD_LDMP_ON: {
            /* Turn live dump on. */
            rv = ras_control(cb, RASCD_SET_LDMP_ON, 0, 0);
            break;
        }
        case RASCD_LDMP_OFF: {
            /* Turn live dump off. */
            rv = ras_control(cb, RASCD_SET_LDMP_OFF, 0, 0);
            break;
        }
        case RASCD_LDMP_LVL: {
            /* Set the detail level at which this component will dump. */
            rv = ras_control(cb, RASCD_SET_LDMP_LVL, arg, 0);
            break;
        }
        case RASCD_LDMP_ESTIMATE: /* fall through */
        case RASCD_LDMP_PREPARE: {
            /*
             * An estimate, as a prepare, is done in the same way.
             * The estimate is received if the livedumpstart command is used
             * and the -e flag is specified.
             * The prepare is received when the component is participating
             * in a live dump. The prepare call is used to request
             * staging buffer space and provide an estimate of the amount
             * of data to be dumped. The sample also requests that the component
             * trace be dumped at this time.
             */
            ldmp_prepare_t *p = (ldmp_prepare_t*)arg;
            int n = 0;
            /* Staging buffer used for dump table */
            p->ldpr_sbufsz = Sbufsz;
            /* Data size - need # cbs */
            for (cbp=get_scbp(Sample_hdr.sh_cbp, p->ldpr_dumpid, 1), n=0;
                 cbp!=INVALID_SCB_PTR;
                 cbp=get_scbp(cbp->scb_next, p->ldpr_dumpid, 1), n++);
            p->ldpr_datasize = TBLESTSZ(n) + sizeof(Sample_hdr) +
                n*sizeof(sample_cb_t);
            /* Dump all of our component trace. */
            rv = dmp_ct(0, p, "", Rascb, 0);
            /*
             * If an error occurred, ldmp_errstr() puts the message in
             * the dump. The sample returns 0, so the dump

```

```

* proceeds normally.
* If the sample returned a value > 0,
* livedump() would put a generic error containing the
* return value in the dump.
*/
if (rv) {
    char str[40];
    sprintf(str, "dmp_ct returned 0x%lx.\n", rv);
    rv = 0;
    (void)ldmp_errstr(p->ldpr_dumpid, Rascb, str);
}
break;
}
case RASCD_LDMP_START:{
/*
* This is received to provide the dump table.
* Because the table is an unlimited table, subsequent
* RASCD_LDMP_AGAIN calls will be received.
*/
ldmp_start_t *p = (ldmp_start_t*)arg;
struct cdt_nn_head_u *hp;
struct cdt_entry_u *ep;
/* The dump table goes in the staging area */
hp = (struct cdt_nn_head_u*)p->ldmpst_buffer;
ep = (struct cdt_entry_u*)((struct cdt_nn_u*)hp)->cdtnu_entry;
/* Set up cdt_nn_head_u */
hp->cdtnu_magic = DMP_MAGIC_NU;
hp->cdtnu_nentries = 1;
/* Set up cdt_entry_u */
ep->du_magic = DMP_MAGIC_UD;
strcpy(ep->du_name, "header");
ep->du_len = sizeof(Sample_hdr);
ep->du_ptr = &Sample_hdr;
ep->du_segval = DUMP_GEN_SEGVAL;
ep->du_xmemp = NULL;
p->ldmpst_table = hp;
/*
* There is a work area in the ldmp_prepare_t and ldmp_start_t
* data areas for use in keeping the state across dump calls,
* generally between RASCD_LDMP_START and RASCD_LDMP_AGAIN
* commands.
* In our case, we'll keep a pointer to the next cb to dump,
* NULL initially.
*/
wkptr = (sample_cb_t**)p->ldmpst_wk;
*wkptr = NULL;
break;
}
case RASCD_LDMP_AGAIN:{
/*
* This is similar to the RASCD_LDMP_START command, but is received to dump
* subsequent data for an unlimited dump table.
*/
int i;
ldmp_start_t *p = (ldmp_start_t*)arg;
struct cdt_nn_head_u *hp;
struct cdt_entry_u *ep;
sample_cb_t *cbp;
/* The dump table goes in the staging area */
hp = (struct cdt_nn_head_u*)p->ldmpst_buffer;
ep = (struct cdt_entry_u*)((struct cdt_nn_u*)hp)->cdtnu_entry;
/* Point to the first/next cb */
wkptr = (sample_cb_t**)p->ldmpst_wk;
/* For the first AGAIN call, cbp will be Sample_hdr.sh_cbp. */
cbp = (*wkptr)? *wkptr: Sample_hdr.sh_cbp;
/* Validate the pointer. */
cbp = get_scbp(cbp, p->ldmpst_dumpid, 1);

```



```

/* Set up cdt_nn_head_u */
hp->cdtnu_magic = DMP_MAGIC_NU;
hp->cdtnu_nentries = 1;
/* Set up cdt_entry_u */
up->dul_magic = DMP_MAGIC_UL;
strcpy(up->dul_name, "cb");
up->dul_nentries = 0;
up->dul_len = sizeof(sample_cb_t);
/* Dump up to 3, NENTRIES, control blocks. */
for (i=0; iscb_next, p->ldmpst_dumpid, 1)) {
    up->dul_entry[i].dle_vmhandle = DUMP_GEN_SEGVAL;
    up->dul_entry[i].dle_ptr = cbp;
    up->dul_nentries++;
}
/* Save address of next cb. */
*wkptr = cbp;
/* Set the table address to NULL when finished. */
p->ldmpst_table = (!up->dul_nentries)? NULL: p->ldmpst_buffer;
break;
}
case RASCD_LDMP_FINISHED:
/* Nothing to do here. */
break;

/* System dump */
case RASCD_SDMP_ON: {
/* Turn system dump on. */
rv = ras_control(cb, RASCD_SET_SDMP_ON, 0, 0);
break;
}
case RASCD_SDMP_OFF: {
/* Turn system dump off. */
rv = ras_control(cb, RASCD_SET_SDMP_OFF, 0, 0);
break;
}
case RASCD_SDMP_LVL: {
/* Set the detail level at which this component will dump. */
rv = ras_control(cb, RASCD_SET_SDMP_LVL, arg, 0);
break;
}
case RASCD_SDMP_ESTIMATE:{
sysdump_estimate_t *p = (sysdump_estimate_t*)arg;
int n;
/* Data size - need # cbs */
for (cbp=get_scbp(Sample_hdr.sh_cbp, 0, 0), n=0;
     cbp!=INVALID_SCB_PTR;
     cbp=get_scbp(cbp->scb_next, 0, 0), n++);
p->se_value = TBLESTSZ(n) + sizeof(Sample_hdr) +
n*sizeof(sample_cb_t);
break;
}
case RASCD_SDMP_START:{
/*
 * This is received to provide the dump table.
 * Since the table is an unlimited table, subsequent
 * RASCD_SDMP_AGAIN calls will be received.
 */
sdmp_start_t *p = (sdmp_start_t*)arg;
struct cdt_nn_head_u *hp;
struct cdt_entry_u *ep;
/* The dump table goes in the staging area */
hp = (struct cdt_nn_head_u*)p->sdmpst_buffer;
ep = (struct cdt_entry_u*)((struct cdt_nn_u*)hp)->cdtnu_entry;
/* Set up cdt_nn_head_u */
hp->cdtnu_magic = DMP_MAGIC_NU;
hp->cdtnu_nentries = 1;
/* Set up cdt_entry_u */

```

```

ep->du_magic = DMP_MAGIC_UD;
strcpy(ep->du_name, "header");
ep->du_len = sizeof(Sample_hdr);
ep->du_ptr = &Sample_hdr
ep->du_segval = DUMP_GEN_SEGVAL;
ep->du_xmemp = NULL;
p->sdpst_table = hp;
/*
 * There is a work area in the sdp_start_t data area
 * for use in keeping the state across dump calls, generally
 * between RASCD_SDMP_START and RASCD_SDMP_AGAIN
 * commands.
 * In our case, we'll keep a pointer to the next cb to dump,
 * NULL initially.
 */
wkptr = (sample_cb_t**)p->sdpst_wk;
*wkptr = NULL;
break;
}
case RASCD_SDMP_AGAIN:{
/*
 * This is similar to RASCD_SDMP_START, but is received to dump
 * subsequent data for an unlimited dump table.
 */
int i;
sdp_start_t *p = (sdp_start_t*)arg;
struct cdt_nn_head_u *hp;
struct cdt_entry_ul *up;
sample_cb_t *cbp;
/* The dump table goes in the staging area */
hp = (struct cdt_nn_head_u*)p->sdpst_buffer;
up = (struct cdt_entry_ul*)((struct cdt_nn_u*)hp)->cdtnu_entry;
/* Point to the first/next cb */
wkptr = (sample_cb_t**)p->sdpst_wk;
/* For the first AGAIN call, cbp will be Sample_hdr.sh_cbp. */
cbp = (*wkptr)? *wkptr: Sample_hdr.sh_cbp;
/* Validate the pointer. */
cbp = get_scbp(cbp, 0, 0);
/* Set up cdt_nn_head_u */
hp->cdtnu_magic = DMP_MAGIC_NU;
hp->cdtnu_nentries = 1;
/* Set up cdt_entry_ul */
up->dul_magic = DMP_MAGIC_UL;
strcpy(up->dul_name, "cb");
up->dul_nentries = 0;
up->dul_len = sizeof(sample_cb_t);
/* Dump up to 3, NENTRIES, control blocks. */
for (i=0; iscb_next, 0, 0) {
    up->dul_entry[i].dle_vmhandle = DUMP_GEN_SEGVAL;
    up->dul_entry[i].dle_ptr = cbp;
    up->dul_nentries++;
}
/* Save address of next cb. */
*wkptr = cbp;
/* Set the table address to NULL when finished. */
p->sdpst_table = (!up->dul_nentries)? NULL: p->sdpst_buffer;
break;
}
case RASCD_SDMP_FINISHED:
/* Nothing to do here. */
break;

case RASCD_DMP_PASS_THROUGH:{
/* pass through */
printf("%s\n", arg);
break;
}
}

```

```

default: {
    printf("bad ras_control command.\n");
    rv = EINVAL_RAS_CONTROL_BADCMD;
}
}

return(rv);
}

/*
 * Allocate sample data
 *
 * Input:
 *   n - number of sample cbs to allocate.
 *
 * Returns:
 *   0 - success
 *   errno - errno from failure.
 */
int
alloc_sample(int n)
{
    sample_cb_t *cbp, *prev_cbp=NULL;

    /* Allocate n cbs */
    while (n--) {
        cbp = xmalloc(sizeof(*cbp), 3, kernel_heap);
        if (!cbp) {
            free_sample();
            return(ENOMEM);
        }
        if (!prev_cbp) Sample_hdr.sh_cbp = cbp;
        else prev_cbp->scb_next = cbp;
        cbp->scb_eyec = EYEC_SCB;
        cbp->scb_next = INVALID_SCB_PTR;
        prev_cbp = cbp;
    }

    return(0);
}

/*
 * Free sample data
 */
void
free_sample()
{
    sample_cb_t *cbp = Sample_hdr.sh_cbp;

    /* Validate cbp */
    cbp = get_scbp(cbp, 0, 0);

    while(cbp != INVALID_SCB_PTR) {
        sample_cb_t *save_cbp = cbp;
        cbp = get_scbp(cbp->scb_next, 0, 0);
        xfree(save_cbp, kernel_heap);
    }
}

/*
 * Validate the cb at the supplied address.
 * This ensures we won't get an exception for a bad scb_next ptr.
 *
 * Input:
 *   addr - - address to read from
 *   id - - dump ID
 *   ldmpflag - 1 if this is a live dump.

```

```

*
* Returns:
*   The pointer at the address.
*   INVALID_SCB_PTR if the memory at addr is bad.
*/
static sample_cb_t *
get_scbp(sample_cb_t *addr, dumpid_t id, int ldmpflag)
{
    sample_cb_t scb;
    char str[80];

    /* Just return if addr is the terminating value. */
    if (addr == INVALID_SCB_PTR) return(addr);

    /* Carefully get the storage at addr. */
    if (raschk_safe_read(addr, &scb, sizeof(scb), RAS_SR_NOFAULT)) {
        /* addr is bad. */
        if (ldmpflag) {
            sprintf(str, "The scb address 0x%lx is bad.\n", addr);
            (void)ldmp_errstr(id, Rascb, str);
        }
        addr = INVALID_SCB_PTR;
    }
    else {
        /* Validate the control block */
        if (scb.scb_eyec != EYEC_SCB) {
            /* cb appears bad. */
            if (ldmpflag) {
                sprintf(str, "cb at 0x%lx is invalid, eyec = 0x%lx.\n",
                    addr, scb.scb_eyec);
                (void)ldmp_errstr(id, Rascb, str);
            }
            addr = INVALID_SCB_PTR;
        }
    }
}

return(addr);
}

```

---

## Component Trace Facility

Component Trace (CT) is an important First Failure Data Capture (FFDC) and Second Failure Data Capture (SFDC) tool available to the kernel, kernel extensions, and device drivers. With the Component Trace facility, a component can capture trace events to aid both debugging and system analysis and can provide focused trace data on larger server systems.

Component Trace uses mechanisms similar to system trace. Existing **TRCHKx** and **TRCGEN** macros can be replaced with **CT** macros to trace into system trace buffers and private buffers of memory-trace mode.

If recorded, Component Trace events can be retrieved with the **ctctrl** command. Extraction with the **ctctrl** command is relevant only to in-memory tracing. Component Trace events can also be present in a system trace. You can use the **trcrpt** command for Component Trace and system trace to process the events.

## Component Trace Modes

Component Trace has two modes that can be used simultaneously:

### system trace mode

The system trace mode sends trace entries to the existing system trace. The following settings can be changed:

Setting	Default setting
on or off	The mode is on.

Setting	Default setting
trace level	The default level of system trace is CT_LVL_NORMAL (level 3).

### memory trace mode

The memory trace mode stores the trace entries in a memory buffer. The buffer is either private to the component or to a per-processor memory buffer that is dedicated to the lightweight memory trace. The following settings can be changed:

Setting	Default setting
on or off	The mode is off.
serialization policy	The recording of trace entries is serialized by the framework.
tracing status	Tracing is suspended.
the size of the private buffer	0
trace level	The default level of memory trace is CT_LVL_NORMAL (level 3).

Component trace entries can be traced to the private buffer of the component, the lightweight memory trace, the system trace, or any combination of these destinations. The destination is governed by flags specified in the **CT\_HOOKx** and **CT\_GEN** macros. The **MT\_COMMON** flag causes the entry to be traced into the common, lightweight-memory-trace buffer. The **MT\_RARE** flag causes the entry to go to the rare, lightweight-memory-trace buffer. Do not specify both the **MT\_COMMON** and **MT\_RARE** flags. The **MT\_PRIV** flag traces the entry into the private buffer of the component. The **MT\_SYSTEM** flag puts the entry into a system trace if system trace is active.

Generic trace entries, which are traced with the **CT\_GEN** macro, cannot be traced into the lightweight memory trace.

In the memory trace mode, you have the choice for each component, at initialization, to store their trace entries either in a component private buffer or in one of the memory buffers managed by the lightweight memory trace. When entries are stored in a lightweight-memory-trace buffer, the memory type (common or rare) is chosen for each trace entry.

The private buffer of the component is a pinned memory buffer that can be allocated by the framework at the component registration or later. The buffer is only attached to this component. You can dynamically change the buffer size with the **CT** API. Administrators can dynamically change the buffer size using the **ctctrl** command.

Private buffers and lightweight memory buffers are used in a circular mode, that is, if the buffer is full, the last trace entries overwrite the first one.

For each component, the serialization of the buffers can be managed either by the component (by the component owner) or by the Component Trace framework. This serialization policy is chosen at registration and cannot be changed during the life of the component.

The system trace mode is another function that is provided by Component Trace. When a component is traced using a system trace, each trace entry is sent to the current system trace. In this mode, Component Trace acts as a front-end filter for the existing system trace. By setting the level of the system trace, a component can control which trace hooks enter the system trace buffer.

## Using the Component Trace Facility

This section covers the life cycle of a component using Component Trace.

- “Initializing a Component Trace” on page 340

- “Unregistering a Component” on page 341
- “Tracing Events into the Component Trace Private Buffer” on page 341
- “Controlling Component Trace for Your Component” on page 341
- “Managing Trace Levels” on page 341

## Initializing a Component Trace

Component Trace is initialized in the process environment, preferably during the initialization of the driver or the kernel subsystem for which the component will be used. The initialization must be done for all components and subcomponents. Initialization is done in three steps:

1. “Registration”
2. “Changing Component Trace Properties” (optional)
3. “Customization” on page 341

**Registration:** A component is registered into the Component Trace framework by calling the **ras\_register** kernel service. The following code is an example of registering a parent (base) component.

```
ras_block_t rasb_eth;

kernno_t err;
...
err = ras_register(&rasb_eth, "ethernet", NULL, RAS_TYPE_NETWORK_ETHERNET, "All ethernet devices",
    RASF_TRACE_AWARE, eth_callback, NULL);
```

The example creates and registers a base component named `ethernet` of the network type and the `ethernet` subtype (the fourth argument). The *type/subtype* field is used to categorize a component. You can get a list of available types and subtypes in the **sys/ras\_base.h** file. The `rasb_eth` argument, which is the component identifier (or reference), is used for all actions to apply to the component. The third argument is set to `NULL`, indicating that this component is a base component with no parent. See “Callback Routine” on page 342 for details about the callback routine (the seventh argument). In the example, no *callback\_data* area (the last argument) is specified. If one is specified, it is passed through the callback routine. For more information, see the **ras\_register** kernel service.

**Changing Component Trace Properties:** Before customizing, you might want to change some of the default Component Trace properties. You can change the following settings:

- To activate the memory trace mode, use the following code:

```
/* set a buffer size (default size is 0) */
err = ras_control(rasb_eth, RASCT_SET_MEMBUFSIZE, size, 0);

/* allocate the private buffer */
err = ras_control(rasb_eth, RASCT_SET_ALLOC_BUFFER, 0, 0);

/* activate memory trace mode */
err = ras_control(rasb_eth, RASCT_SET_MEMTRC_RESUME, 0, 0);
```

If you activate memory trace mode without allocating a private buffer, leave out the `RASCT_SET_MEMBUFSIZE` and `RASCT_SET_ALLOC_BUFFER` calls. This is useful if you want to just send events to the system trace or to the lightweight memory trace. Failure in RAS routines should not cause a driver or kernel subsystem to abort initialization.

- To use component serialization rather than infrastructure serialization, use the following code:

```
err = ras_control(rasb_eth, RASCT_SET_CT_SERIALIZED, 0, 0);
```

You cannot change this setting after customization.

- To change the level of trace to the minimal value for system trace mode, use the following code:

```
err = ras_control(rasb_eth, RASCT_SET_SYSTRC_LVL, CT_LVL_MINIMAL, 0);
```

You can get the definitions of available levels in the **ras\_trace.h** file.

- To change the level of trace to a detailed value for memory trace mode, use the following code:

```
err = ras_control(rasb_eth, RASCT_SET_MEMTRC_LVL, CT_LVL_DETAIL, 0);
```

For more information, see the **ras\_control** kernel service.

**Customization:** The customization step enables component tracing for the specified component and retrieves any saved component settings. Customization is useful for keeping settings of components after a system restart.

The customization step is mandatory to put the component into a usable state. Before this call, no tracing is active even if either trace mode has been activated. You can customize the ethernet component as follows, continuing the examples in the “Changing Component Trace Properties” on page 340 section:

```
err = ras_customize(rasb_eth);
```

After this call, the initialization of your component with the Component Trace framework is complete, and tracing is active unless a persistent value has been set to indicate that the Component Trace is off.

For example, if the **ctctrl** command was used to set Component Tracing to be persistently off, any call to the **ras\_customize** kernel service results in the component trace of the calling component being turned off. In this case, unless the **ctctrl** command is used to set Component Trace to be persistently on, all requests to turn on or resume component tracing are denied.

## Unregistering a Component

If the registered component is a driver, a kernel extension or a subsystem that can be stopped. Unregister your component at stop or unload time. A component cannot be unregistered from the framework if it has subcomponents. Therefore, subcomponents must be unregistered first. Moreover, the memory trace mode must be stopped.

```
ras_control(rasb_eth, RASCT_SET_MEMTRC_SUSPEND, 0, 0); /* might need to be serialized */
ras_unregister(rasb_eth);
```

The **ras\_unregister** call must be done from the process environment. No other driver or subsystem operations should occur during **ras\_unregister** calls. For more information, see the **ras\_unregister** kernel service.

## Tracing Events into the Component Trace Private Buffer

For more information about how to trace events into the Component Trace private buffer, see the **CT\_HOOKx** and **CT\_GEN** macros.

## Controlling Component Trace for Your Component

For more information about how to control Component Trace and dump Component Trace buffers, see the **ctctrl** command .

## Managing Trace Levels

The trace levels are set to a value of -1 if the mode is suspended or off.

To access the memory trace mode level of a component, use the following code:

```
rasrb_trace_memlevel(rasb_eth)
```

To access the system trace mode level, use the following code:

```
rasrb_trace_syslevel(rasb_eth)
```

You can use the **CT\_TRCON(rasb\_eth, level)** macro to know if the trace entries traced at that level are recorded in one or both of the trace modes (system or memory).

The system mode is considered to be on if system trace is running and the system trace mode is on for the specified component.

## Callback Routine

For components that are registered for Component Trace, a callback routine is mandatory. The routine is called by the Component Trace framework to inform registered components about an event that might require action. The component can accept or reject a request by either performing or not performing the appropriate commands that the **ras\_control** kernel service passes to the callback.

The callback routine is called when a new setting is requested using the **ctctrl** user command (except for DR events). The following commands must be handled to ensure the correct operation of the **ctctrl** command:

Command	Description
<b>RASCT_MEMTRC_ON</b>	Sets the memory trace mode on for this component. This command can only be called from the process environment. The <b>RASCT_MEMTRC_ON</b> command is passed to the component through the callback routine. The component must perform the <b>RASCT_SET_ALLOC_BUFFER</b> command to allocate a buffer and perform the <b>RASCT_SET_MEMTRC_RESUME</b> command to enable memory trace mode.
<b>RASCT_MEMTRC_OFF</b>	Sets the memory trace mode off for this component. This command can only be called from the process environment. The <b>RASCT_MEMTRC_OFF</b> command is passed to the component through the callback routine. The component must perform the <b>RASCT_SET_MEMTRC_SUSPEND</b> command to disable memory trace mode followed by the <b>RASCT_SET_FREE_BUFFER</b> command to free the buffer. If the component is not framework-serialized, the component must serialize the <b>RASCT_SET_MEMTRC_SUSPEND</b> call.  When the <b>RASCT_MEMTRC_OFF</b> command is received, the component must free its private buffer. Otherwise, the framework will free the private buffer. The <b>RASCT_MEMTRC_SUSPEND</b> command is used instead of the <b>RASCT_MEMTRC_OFF</b> command if the buffer is needed again.
<b>RASCT_MEMTRC_SUSPEND</b> <b>RASCT_SET_MEMTRC_SUSPEND</b>	Suspends the memory trace mode for this component by turning this mode off without freeing any private buffer. This command can only be called from the process environment. The <b>RASCT_MEMTRC_SUSPEND</b> command is passed to the component through the callback routine. The component must perform the <b>RASCT_SET_MEMTRC_SUSPEND</b> command to apply the settings. If the component is not framework-serialized, the component must serialize the <b>RASCT_SET_MEMTRC_SUSPEND</b> call.
<b>RASCT_MEMTRC_RESUME</b> <b>RASCT_SET_MEMTRC_RESUME</b>	Starts the memory trace mode for this component to record trace events.  The command is usually used after tracing has suspended with a <b>RASCT_SET_MEMTRC_SUSPEND</b> command. This command can only be called from the process environment. The <b>RASCT_MEMTRC_RESUME</b> command is passed to the component through the callback routine. The component must perform a <b>RASCT_SET_MEMTRC_RESUME</b> control command to apply the settings. This command fails if a resize or save operation is in progress, or if a buffer is required but not allocated.
<b>RASCT_SYSTRC_ON</b> <b>RASCT_SET_SYSTRC_ON</b>	Sets the system trace mode on for this component. This command can be called in both interrupt and process environments. The <b>RASCT_SYSTRC_ON</b> command is passed to the component through the callback routine. The component must perform the <b>RASCT_SET_SYSTRC_ON</b> control command to apply the settings.
<b>RASCT_SYSTRC_OFF</b> <b>RASCT_SET_SYSTRC_OFF</b>	Sets the system trace mode off for this component. This command can be called in both interrupt and process environments. The <b>RASCT_SYSTRC_OFF</b> command is passed to the component through the callback routine. The component must perform the <b>RASCT_SET_SYSTRC_OFF</b> control command to apply the settings.



Command	Description
<b>RASCT_MEMTRC_LVL</b> <b>RASCT_SET_MEMTRC_LVL</b>	Changes the level of trace for the memory trace mode. The level of trace must be passed through the <i>arg</i> parameter of the <b>ras_control</b> kernel service. This command can be called in both interrupt and process environments. The <b>RASCT_MEMTRC_LVL</b> command is passed to the component through the callback routine. The component must perform the <b>RASCT_SET_MEMTRC_LVL</b> control command to apply the settings.
<b>RASCT_SYSTRC_LVL</b> <b>RASCT_SET_SYSTRC_LVL</b>	Changes the level of trace for the system trace mode. The level of trace must be passed through the <i>arg</i> parameter. This command can be called in both interrupt and process environments. The <b>RASCT_SYSTRC_LVL</b> command is passed to the component through the callback routine. The component must perform the <b>RASCT_SET_SYSTRC_LVL</b> control command to apply the settings.
<b>RASCT_MEMBUFSIZE</b>	Changes the size of the private buffer. The new size must be given in bytes through the <i>arg</i> parameter. This command can be called only in the process environment. The <b>RASCT_MEMBUFSIZE</b> command is passed to the component; the component must call the <b>RASCT_SET_MEMBUFSIZE</b> control command to apply the property. The <b>RASCT_SET_MEMBUFSIZE</b> command can be called regardless of whether a private buffer is allocated. But first, the memory trace mode must be suspended with the <b>RASCT_SET_MEMTRC_SUSPEND</b> command if the component is already on. Also, if a private buffer was allocated, the last traced entries are preserved. The exact amount depends on the new buffer size. Finally, the <b>RASCT_SET_ALLOC_BUFFER</b> command must be used. The command is necessary if no private buffer was allocated previously, but it can be used in either case. If a nonzero value is requested, a minimum value of 2 times the size of a full trace entry (with five data words) is required.  The <b>RASCT_MEMBUFSIZE</b> command is called if the <code>ctctl memtracebufsize=size</code> user command is performed. The component might want to keep track of the suspend or resume state in order not to resume the memory trace mode if this mode was not on before the call.
<b>RASCT_SET_ALLOC_BUFFER</b>	Allocates the private memory buffer. This command can be called only from the process environment. This command fails if the memory trace mode is not suspended.
<b>RASCT_SET_FREE_BUFFER</b>	Frees the private memory buffer. This command can be called only from the process environment. This command fails if the memory trace mode is not suspended, or if a resize or save operation is in progress.
<b>RASCT_SET_CT_SERIALIZE</b>	Sets the serialization policy for the component. If the value of the <i>arg</i> parameter is TRUE (nonzero), the framework serializes buffer access. Otherwise, it must be serialized by the component. The default value is TRUE. This property must be set at initialization phase, and it can not be changed after the customization phase. Moreover, only the component can set the property. This command can be called in both interrupt and process environments.
<b>RASCT_DR_MEM</b>	Informs the component of add-memory or remove-memory operation. The command is passed from the process environment. The total amount of system memory that is added or removed is passed to the component through the <i>arg</i> field of the <b>ras_callback</b> function. A component can resize its private buffer at this time.
<b>RASCT_PASS_THROUGH</b>	Passes string data to the callback function. The <b>ctctl</b> command can be used to pass a string (for example, <code>ctctl -c socket "passthrough string"</code> ). The string is passed with the <i>arg</i> parameter. The component can use this command to perform some specific actions. It returns the <b>EINVAL</b> error code for all unavailable commands. This command can be used only from the process environment.

Command	Description
<b>RASCT_GETBUFFER</b>	Retrieves a copy of the current contents of the private buffer. The memory trace mode must be disabled during the copy to preserve the integrity of the copy. If the memory trace mode is enabled, this command attempts to issue the <b>RASCT_MEMTRC_SUSPEND</b> command, retrieve the buffer contents, and then issue the <b>RASCT_MEMTRC_RESUME</b> command. The <i>arg</i> and <i>size</i> parameters must contain the pointer and the size of the provided buffer for the copy. This command can be called only from the process environment.

By convention, the callback must return a value of 0, if it accepts the setting and all **ras\_control** commands run by the callback return successfully. Otherwise the callback must return an error value.

All **RASCT\_SET\_ALLOC\_BUFFER**, **RASCT\_SET\_FREE\_BUFFER**, and **RASCT\_SET\_MEMBUFSIZE** commands must be called when memory trace mode is disabled (suspended) to ensure that the buffer is not used for tracing (by **CT\_HOOK** calls). The **RASCT\_SET\_MEMTRC\_SUSPEND** command must then be serialized the same way as the **CT\_HOOKx** and **CT\_GEN** calls.

Components can update some internal values or keep track of the context modification using the **callback\_data** area. The **callback\_data** area that is passed to the callback routine is the one given at the registration time.

---

## Error Logging

The error facility records device-driver entries in the system error log. These error log entries record any software or hardware failures that need to be available either for informational purposes or for fault detection and corrective action. The device driver, using the **errsave** kernel service, adds error records to the **/dev/error** special file.

The **errdemon** daemon picks up the error record and creates an error log entry. When you access the error log either through SMIT (System Management Interface Tool) or with the **errpt** command, the error record is formatted according to the error template in the error template repository and presented in either a summary or detailed report.

Before initiating the error logging process, determine what services are available to developers, and what services are available to the customer, service personnel, and defect personnel.

- **Determine the Importance of the Error:** Use system resources for logging only information that is important or helpful to the intended audience. Work with the hardware developer, if possible, to identify detectable errors and the information that should be relayed concerning those errors.
- **Determine the Text of the Message:** Use regular national language support (NLS) XPG/4 messages instead of the codepoints. For more information about NLS messages, see *Message Facility in AIX Version 6.1 National Language Support Guide and Reference*.
- **Determine the Correct Level of Thresholding:** Each software or hardware error to be logged, can be limited by thresholding to avoid filling the error log with duplicate information. Side effects of runaway error logging include overwriting existing error log entries and unduly alarming the end user. The error log is limited in size. When its size limit is reached, the log wraps. If a particular error is repeated needlessly, existing information is overwritten, which might cause inaccurate diagnostic analyses. The end user or service person can perceive a situation as more serious or pervasive than it is if they see hundreds of identical or nearly identical error entries.

You are responsible for implementing the proper level of thresholding in the device driver code.

The default size of the error log is 1 MB. As shipped, it cleans up any entries older than 30 days. To ensure that your error log entries are informative, noticed, and remain intact, *test your driver thoroughly*.

## Setting up Error Logging

To begin error logging, do the following:

1. Construct error templates.
2. Add error logging calls into your code.

### Constructing error templates

An error template is used to associate an error with message text which is output by the **errpt** command when the error log is viewed. Error templates are described in detail in the **errupdate** command article. You normally want to define your own error text. The error text can come from an XPG/4 message catalogue. There are also many predefined error text messages which you can view with the **errmsg -w ALL** command. For information on the use of error text in error templates, see the **errupdate** command.

### Adding Error Logging Calls into the Code

The **errsave** kernel service allows the kernel and kernel extensions to write to the error log. Typically, you define a routine in the device driver that can be called by other device driver routines when a loggable error is encountered. This function takes the data passed to it, puts it into the proper structure and calls the **errsave** kernel service. The syntax for the **errsave** kernel service is:

```
#include <sys/errids.h>
void errsave(buf, cnt)
char *buf;
unsigned int cnt;
```

where:

**buf** Specifies a pointer to a buffer that contains an error record as described in the **sys/errids.h** header file.  
**cnt** Specifies a number of bytes in the error record contained in the buffer pointed to by the *buf* parameter.

The following sample code is an example of a device driver error logging routine. This routine takes data passed to it from some part of the main body of the device driver. This code simply fills in the structure with the pertinent information, then passes it on using the **errsave** kernel service.

```
void
errsv_ex (int err_id, unsigned int port_num,
          int line, char *file, uint data1, uint data2)
{
    dderr    log;
    char     errbuf[255];
    ddex_dds *p_dds;

    p_dds = dds_dir[port_num];
    log.err.error_id = err_id;

    if (port_num = BAD_STATE) {
        sprintf(log.err.resource_name, "%s :%d",
            p_dds->dds_vpd.adpt_name, data1);
        data1 = 0;
    }

    else
        sprintf(log.err.resource_name, "%s", p_dds->dds_vpd.devname);

    sprintf(errbuf, "line: %d file: %s", line, file);
    strncpy(log.file, errbuf, (size_t)sizeof(log.file));

    log.data1 = data1;
    log.data2 = data2;

    errsave(&log, (uint)sizeof(dderr)); /* run actual logging */
} /* end errlog_ex */
```

The data to be passed to the **errsave** kernel service is defined in the **dderr** structure, which is defined in a local header file, **dderr.h**. The definition for **dderr** is:

```
typedef struct dderr {
    struct err_rec0 err;
    int data1; /* use data1 and data2 to show detail */
    int data2; /* data in the errlog report. Define */
              /* these fields in the errlog template */
              /* These fields may not be used in all */
              /* cases. */
} dderr;
```

The first field of the **dderr.h** header file is comprised of the **err\_rec0** structure, which is defined in the **sys/err\_rec.h** header file. This structure contains the ID (or label) and a field for the resource name. The two data fields hold the detail data for the error log report. As an alternative, you could simply list the fields within the function.

You can also log a message into the error log from the command line. To do this, use the **errlogger** command.

Errors can be logged from applications using the **errlog** subroutine.

**Note:** Care must be taken when logging a data structure, because error logging does not support padding done by the compiler.

Consider the following data structure:

```
struct {
    struct err_rec0 err;
    long data;
} myerr;
```

Because **err\_rec0** is 20 bytes in length, 0x14 bytes, the compiler normally inserts 4 bytes of padding before data, when compiling in 64-bit mode. The structure then looks like the following:

```
struct {
    struct err_rec0 err;
    int padding;
    long data; /* 64 bits of data, 64-bit aligned */
} myerr;
```

Thus the **Detail\_Data** item in the template begins formatting at the padding data item rather than at data.

This can be overcome, if you use the Xlc compiler as follows:

```
#pragma options align=packed
struct {
    struct err_rec0 err;
    long data;
} myerr;
#pragma options align=reset
```

After you add the templates using the **errupdate** command, compile the device driver code along with the new header file. Simulate the error and verify that it was written to the error log correctly. Some details to check for include:

- Is the error demon running? This can be verified by running the **ps -ef** command and checking for **/usr/lib/errdemon** as part of the output.
- Is the error part of the error template repository? Verify this by running the **errpt -at** command.
- Was the new header file, which was created by the **errupdate** command and which contains the error label and unique error identification number, included in the device driver code when it was compiled?

---

## Debug and Performance Tracing

The **trace** facility is useful for observing a running device driver and system. The **trace** facility captures a sequential flow of time-stamped system events, providing a fine level of detail on system activity. Events are shown in time sequence and in the context of other events. The **trace** facility is useful in expanding the trace event information to understand who, when, how, and even why the event happened.

### Introduction

The operating system is shipped with permanent trace event points. These events provide general visibility to system execution. You can extend the visibility into applications by inserting additional events and providing formatting rules.

The collection of **trace** data was designed so that system performance and flow would be minimally altered by activating **trace**. Because of this, the facility is extremely useful as a performance analysis tool and as a problem determination tool.

The **trace** facility is more flexible than traditional system monitor services that access and present statistics maintained by the system. With traditional monitor services, data reduction (conversion of system events to statistics) is largely coupled to the system instrumentation. For example, the system can maintain the minimum, maximum, and average elapsed time observed for runs of a task and permit this information to be extracted.

The **trace** facility does not strongly couple data reduction to instrumentation but provides a stream of system events. It is not required to presuppose what statistics are needed. The statistics or data reduction are to a large degree separated from the instrumentation.

You can choose to develop the minimum, maximum, and average time for task A from the flow of events. But it is also possible to extract the average time for task A when called by process B, extract the average time for task A when conditions XYZ are met, develop a standard deviation for task A, or even decide that some other task, recognized by a stream of events, is more meaningful to summarize. This flexibility is invaluable for diagnosing performance or functional problems.

The **trace** facility generates large volumes of data. This data cannot be captured for extended periods of time without overflowing the storage device. This allows two practical ways that the **trace** facility can be used natively.

First, the **trace** facility can be triggered in multiple ways to capture small increments of system activity. It is practical to capture seconds to minutes of system activity in this way for post-processing. This is sufficient time to characterize major application transactions or interesting sections of a long task.

Second, the **trace** facility can be configured to direct the event stream to standard output. This allows a real-time process to connect to the event stream and provide data reduction in real-time, thereby creating a long term monitoring capability. A logical extension for specialized instrumentation is to direct the data stream to an auxiliary device that can either store massive amounts of data or provide dynamic data reduction.

For AIX 5.3, tracing can be limited to a specified set of processes or threads. This can greatly reduce the amount of data generated and allow you to target the trace to report on specific tasks of interest.

You can start the system trace from:

- The command line
- SMIT
- Software

The trace facility causes predefined events to be written to a trace log. The tracing action is then stopped.

Tracing from a command line is discussed in “Controlling trace.” Tracing from a software application is discussed and an example is presented in “Examples of Coding Events and Formatting Events” on page 364.

After a trace is started and stopped, you must format it before viewing it.

To format the trace events that you have defined, you must provide a stanza that describes how the trace formatter is to interpret the data that has been collected. This is described in “Syntax for Stanzas in the trace Format File” on page 351.

The **trcrpt** command provides a general purpose report facility. The report facility provides little data reduction, but converts the raw binary event stream to a readable ASCII listing of the event stream. Data can be visually extracted by a reader, or tools can be developed to further reduce the data.

For an event to be traced, you must write an *event hook* (sometimes called a *trace hook*) into the code that you want to trace. Tracing can be done on either the system channel (channel 0) or on a generic channel (channels 1-7). All preshipped trace points are output to the system channel.

Usually, when you want to show interaction with other system routines, use the system channel. The generic channels are provided so that you can control how much data is written to the trace log. Only your data is written to one of the generic channels.

For more information on trace hooks, see “Macros for Recording trace Events” on page 349.

## Using the trace Facility

The following sections describe the use of the **trace** facility.

### Configuring and Starting trace Data Collection

The **trace** command configures the trace facility and starts data collection. You can start **trace** from the command line or with a **trcstart** subroutine call. The **trcstart** subroutine is in the **librts.a** library. The syntax of the **trcstart** subroutine is:

```
int trcstart(char *args)
```

where *args* is simply the options list desired that you would enter using the trace command if starting a system trace (channel 0). If starting a generic trace, include a **-g** option in the *args* string. On successful completion, **trcstart** returns the channel ID. For generic tracing this channel ID can be used to record to the private generic channel.

For an example of the **trcstart** routine, see “Examples of Coding Events and Formatting Events” on page 364.

When compiling a program using this subroutine, you must request the link to the **librts.a** library. Use **-l rts** as a compile option.

## Controlling trace

Basic controls for the **trace** facility exist as trace subcommands, standalone commands, and subroutines.

If you configure the **trace** routine to run asynchronously (the **-a** option), you can control the trace facility with the following commands:

<b>trcon</b>	Turns collection of trace data on.
<b>trcoff</b>	Turns collection of trace data off.
<b>trcstop</b>	Stops collection of trace data (like <b>trcoff</b> ) and terminates the <b>trace</b> routine.

## Producing a trace Report

The trace report facility formats and displays the collected event stream in readable form. This report facility displays text and data for each event according to rules provided in the trace format file. The default trace format file is `/etc/trcfmt` and contains a stanza for each event ID. The stanza for the event provides the report facility with formatting rules for that event. This technique allows you to add your own events to programs and insert corresponding event stanzas in the format file to have their new events formatted.

This report facility does not attempt to extract summary statistics (such as CPU utilization and disk utilization) from the event stream. This can be done in several other ways. To create simple summaries, consider using `awk` scripts to process the output obtained from the `trcrpt` command.

## Defining trace Events

The operating system is shipped with predefined trace hooks (events). You need only activate `trace` to capture the flow of events from the operating system. You might want to define trace events in your code during development for tuning purposes. This provides insight into how the program is interacting with the system. The following sections provide the information that is required to do this.

### Possible Forms of a trace Event Record

A trace event can take several forms. An event consists of the following:

- Hookword
- Data words (optional)
- A TID, or thread identifier
- Timestamp

An event record should be as short as possible. Many system events use only the hookword and timestamp. There is another event type you should seldom use because it is less efficient. It is a long format that allows you to record a variable length data. In this long form, the 16-bit data field of the hookword is converted to a *length* field that describes the length of the event record.

### Macros for Recording trace Events

The following macros should always be used to generate trace data. Do not call the tracing functions directly. There is a macro to record each possible type of event record. The macros are defined in the `sys/trcmacros.h` header file. Most event IDs are defined in the `sys/trchkid.h` header file. Include these two header files in any program that is recording `trace` events.

The macros to record system (channel 0) events with a time stamp are:

- `TRCHKL0T` (hw)
- `TRCHKL1T` (hw,D1)
- `TRCHKL2T` (hw,D1,D2)
- `TRCHKL3T` (hw,D1,D2,D3)
- `TRCHKL4T` (hw,D1,D2,D3,D4)
- `TRCHKL5T` (hw,D1,D2,D3,D4,D5)

Similarly, to record non-time stamped system events (channel 0) on versions of AIX prior to AIX 5L Version 5.3 with the 5300-05 Technology Level, use the following macros:

- `TRCHKL0` (hw)
- `TRCHKL1` (hw,D1)
- `TRCHKL2` (hw,D1,D2)
- `TRCHKL3` (hw,D1,D2,D3)
- `TRCHKL4` (hw,D1,D2,D3,D4)

- **TRCHKL5** (hw,D1,D2,D3,D4,D5)

In AIX 5L Version 5.3 with the 5300-05 Technology Level and above, a time stamp is recorded with each event regardless of the type of macro used.

There are only two macros to record events to one of the generic channels (channels 1-7). They are:

- **TRCGEN** (ch,hw,d1,len,buf)
- **TRCGENT** (ch,hw,d1,len,buf)

These macros record a hookword (hw), a data word (d1), and a length of data (len) specified in bytes from the user's data segment at the location specified (buf) to the event stream specified by the channel (ch). In AIX 5L Version 5.3 with the 5300-05 Technology Level and above, the time stamp is recorded with both macros.

## Use of Event IDs

Applications running on systems earlier than AIX 6.1 and 32-bit applications running on AIX 6.1 use 12-bit (or 3-hex-digit) event IDs (hook IDs), for a maximum of 4096 IDs. Beginning with AIX 6.1, 64-bit applications and kernel routines can use 16-bit (or 4-hex-digit) event IDs for a maximum of 65 536 IDs. Event IDs that are permanently left in and shipped with code must be permanently assigned.

If a 16-bit ID is less than 0x1000, its least significant digit must be 0 (in the form of 0x0hh0 where *h* is a hexadecimal digit).

**Tip:** A 12-bit hook ID in the form of 0xhhh is equivalent to a 16-bit ID in the form of 0xhhh0 where *h* is a hexadecimal digit.

To allow you to define events in your environments or during development, a range of event IDs exist for temporary use. Event IDs for temporary use ranges from hex 010 through hex 0FF for applications running on systems earlier than AIX 6.1 and 32-bit applications running on AIX 6.1 and later releases. For 64-bit applications and kernel routines beginning with AIX 6.1, event IDs for temporary use ranges from hex 0100 through hex 0FF0. No permanent (shipped) events are assigned in this range. You can freely use this range of IDs in your own environment. If you do use IDs in this range, do not let the code leave your environment.

Permanent events must have event IDs assigned by the current owner of the trace component. To obtain a trace event id, send a note with a subject of help to aixras@austin.ibm.com.

You should conserve event IDs because they are limited. Event IDs can be extended by the data field. The only reason to have a unique ID is that an ID is the level at which collection and report filtering is available in the trace facility. An ID can be collected or not collected by the trace collection process and reported or not reported by the trace report facility. Whole applications can be instrumented using only one event ID. The only restriction is that the granularity on choosing visibility is to choose whether events for that application are visible.

A new event can be formatted by the trace report facility (**trcrpt** command) if you create a stanza for the event in the trace format file. The trace format file is an editable ASCII file. The syntax for a format stanzas is shown in Syntax for Stanzas in the trace Format File. All permanently assigned event IDs should have an appropriate stanza in the default trace format file shipped with the base operating system.

## Suggested Locations and Data for Permanent Events

The intent of permanent events is to give an adequate level of visibility to determine execution, and data flow and have an adequate accounting for how CPU time is being consumed. During code development, it can be desirable to make very detailed use of trace for a component. For example, you can choose to trace the entry and exit of every subroutine in order to understand and tune path length. However, this would generally be an excessive level of instrumentation to ship for a component.



Consult a performance analyst for decisions regarding what events and data to capture as permanent events for a new component. The following paragraphs provide some guidelines for these decisions.

Events should capture execution flow and data flow between major components or major sections of a component. For example, there are existing events that capture the interface between the virtual memory manager and the logical volume manager. If work is being queued, data that identifies the queued item (a handle) should be recorded with the event. When a queue element is being processed, the "dequeue" event should provide this identifier as data also, so that the queue element being serviced is identified.

Data or requests that are identified by different handles at different levels of the system should have events and data that allow them to be uniquely identified at any level. For example, a read request to the physical file system is identified by a file descriptor and a current offset in the file. To a virtual memory manager, the same request is identified by a segment ID and a virtual page address. At the disk device driver level, this request is identified as a pointer to a structure that contains pertinent data for the request.

The file descriptor or segment information is not available at the device driver level. Events must provide the necessary data to link these identifiers so that, for example, when a disk interrupt occurs for incoming data, the identifier at that level (which can simply be the buffer address for where the data is to be copied) can be linked to the original user request for data at some offset into a file.

Events should provide visibility to major protocol events such as requests, responses, acknowledgements, errors, and retries. If a request at some level is fragmented into multiple requests, a trace event should indicate this and supply linkage data to allow the multiple requests to be tracked from that point. If multiple requests at some level are coalesced into a single request, a trace event should also indicate this and provide appropriate data to track the new request.

Use events to give visibility to resource consumption. Whenever resources are claimed, returned, created, or deleted an event should record the fact. For example, claiming or returning buffers to a buffer pool or growing or shrinking the number of buffers in the pool.

The following guidelines can help you determine where and when you should have trace hooks in your code:

- Tracing entry and exit points of every function is not necessary. Provide only key actions and data.
- Show linkage between major code blocks or processes.
- If work is queued, associate a name (handle) with it and output it as data.
- If a queue is being serviced, the trace event should indicate the unique element being serviced.
- If a work request or response is being referenced by different handles as it passes through different software components, trace the transactions so the action or receipt can be identified.
- Place trace hooks so that requests, responses, errors, and retries can be observed.
- Identify when resources are claimed, returned, created, or destroyed.

Also note that:

- A trace ID can be used for a group of events by "switching" on one of the data fields. This means that a particular data field can be used to identify from where the trace point was called. The trace format routine can be made to format the trace data for that unique trace point.
- The trace hook is the level at which a group of events can be enabled or disabled. Note that trace hooks can be grouped in SMIT. For more information, see "Trace Event Groups" on page 366.

## **Syntax for Stanzas in the trace Format File**

The intent of the trace format file is to provide rules for presentation and display of the expected data for each event ID. This allows new events to be formatted without changing the report facility. Rules for new events are simply added to the format file. The syntax of the rules provide flexibility in the presentation of the data.

A trace format stanza can be as long as required to describe the rules for any particular event. The stanza can be continued to the next line by terminating the present line with a backslash (\). The fields are:

#### **event\_id**

Each stanza begins with the 3-digit hexadecimal event ID that the stanza describes, followed by a space.

**Note:** Beginning with AIX 6.1, each stanza can also begin with a 4-digit hexadecimal event ID.

**V.R** This field describes the version (V) and release (R) that the event was first assigned. Any integers work for V and R, and you might want to keep your own tracking mechanism.

**L=** The text description of an event can begin at various indentation levels. This improves the readability of the report output. The indentation levels correspond to the level at which the system is running. The recognized levels are:

**APPL** Application level

**SVC** Transitioning system call

**KERN** Kernel level

**INT** Interrupt

#### **event\_label**

The *event\_label* is an ASCII text string that describes the overall use of the event ID. This is used by the **-j** option of the **trcrpt** command to provide a listing of events and their first level description. The event label also appears in the formatted output for the event unless the *event\_label* field starts with an @ character.

**\n** The event stanza describes how to parse, label, and present the data contained in an event record. You can insert a **\n** (newline) in the event stanza to continue data presentation of the data on a new line. This allows the presentation of the data for an event to be several lines long.

**\t** The **\t** (tab) function inserts a tab at the point it is encountered in parsing the description. This is similar to the way the **\n** function inserts new lines. Spacing can also be inserted by spaces in the *data\_label* or *match\_label* fields.

#### **starttimer(,#,#)**

The *starttimer* and *endtimer* fields work together. The (,#,#) field is a unique identifier that associates a particular *starttimer* value with an *endtimer* that has the same identifier. By convention, if possible, the identifiers should be the ID of starting event and the ID of the ending event.

When the report facility encounters a start timer directive while parsing an event, it associates the starting events time with the unique identifier. When an end timer with the same identifier is encountered, the report facility outputs the delta time (this appears in brackets) that elapsed between the starting event and ending event. The begin and end system call events make use of this capability. On the return from system call event, a delta time is shown that indicates how long the system call took.

#### **endtimer(,#,#)**

See the *starttimer* field in the preceding paragraph.

#### **data\_descriptor**

The *data\_descriptor* field is the fundamental field that describes how the report facility consumes, labels, and presents the data.

The various subfields of the *data\_descriptor* field are:

#### **data\_label**

The data label is an ASCII string that can optionally precede the output of data consumed by the following format field.

## format

You can think of the report facility as having a pointer into the data portion of an event. This data pointer is initialized to point to the beginning of the event data (the 16-bit data field in the hookword). The format field describes how much data the report facility consumes from this point and how the data is considered. For example, a value of **Bm.n** tells the report facility to consume *m* bytes and *n* bits of data and to consider it as binary data.

The possible format fields are described in the following section. If this field is not followed by a comma, the report facility outputs the consumed data in the format specified. If this field is followed by a comma, it signifies that the data is not to be displayed but instead compared against the following `match_vals` field. The data descriptor associated with the matching `match_val` field is then applied to the remainder of the data.

## match\_val

The match value is data of the same format described by the preceding format fields. Several match values typically follow a format field that is being matched. The successive match fields are separated by commas. The last match value is not followed by a comma. Use the character string `\*` as a pattern-matching character to match anything. A pattern-matching character is frequently used as the last element of the `match_val` field to specify default rules if the preceding `match_val` field did not occur.

## match\_label

The match label is an ASCII string that is output for the corresponding match.

Each of the possible format fields is described in the comments of the `/etc/trcfmt` file. The following shows several possibilities:

**Format field**                      descriptions

In most cases, the data length part of the specifier can also be the letter **"W"** which indicates that the word size of the trace hook is to be used. For example, **XW** will format 4 or 8 bytes into hexadecimal, depending upon whether the trace hook comes from a 32 or 64 bit environment.

<b>Am.n</b>	This value specifies that <i>m</i> bytes of data are consumed as ASCII text, and that it is displayed in an output field that is <i>n</i> characters wide. The data pointer is moved <i>m</i> bytes.
<b>S1, S2, S4</b>	Left justified string. The length of the field is defined as 1 byte (S1), 2 bytes (S2), or 4 bytes (S4) and so on. The data pointer is moved accordingly. <b>SW</b> indicates that the word size for the trace event is to be used.
<b>Bm.n</b>	Binary data of <i>m</i> bytes and <i>n</i> bits. The data pointer is moved accordingly.
<b>Xm</b>	Hexadecimal data of <i>m</i> bytes. The data pointer is moved accordingly.
<b>D2, D4</b>	Signed decimal format. Data length of 2 (D2) bytes or 4 (D4) bytes is consumed.
<b>U2, U4</b>	Unsigned decimal format. Data length of 2 or 4 bytes is consumed.
<b>F4, F8</b>	Floating point of 4 or 8 bytes.
<b>Gm.n</b>	Positions the data pointer. It specifies that the data pointer is positioned <i>m</i> bytes and <i>n</i> bits into the data.
<b>Om.n</b>	Skip or omit data. It omits <i>m</i> bytes and <i>n</i> bits.
<b>Rm</b>	Reverse the data pointer <i>m</i> bytes.
<b>Wm</b>	Position <code>DATA_POINTER</code> at word <i>m</i> . The word size is either 4 or 8 bytes, depending upon whether or not this is a 32 or 64 bit format trace. This bares no relation to the <code>%W</code> format specifier.
<b>m8</b>	Output the next 8 bytes as time in milliseconds from the beginning of the trace. <b>mW</b> will format only 8 bytes of data. The <code>DATA_POINTER</code> is advanced by 8 bytes.
<b>u4, u8</b>	Output the next 4 or 8 bytes as time in microseconds. <b>mW</b> will format either 4 or 8 bytes of data depending on whether the current hook is 32 or 64 bits. The <code>DATA_POINTER</code> is advanced by 4 or 8 bytes.

Some macros are provided that can be used as format fields to quickly access data. For example:

\$D1, \$D2, \$D3, \$D4, \$D5	These macros access data words 1 through 5 of the event record without moving the data pointer. The data accessed by a macro is hexadecimal by default. A macro can be cast to a different data type (X, D, U, B) by using a % character followed by the new format code. For example, the following macro causes data word one to be accessed, but to be considered as 2 bytes and 3 bits of binary data: \$D1%B2.3
\$HD	This macro accesses the first 16 bits of data contained in the hookword, in a similar manner as the \$D1 through \$D5 macros access the various data words. It is also considered as hexadecimal data, and also can be cast.

You can define other macros and use other formatting techniques in the trace format file. This is shown in the following trace format file example.

### Example Trace Format File

```
# @(#)65 1.142 src/bos/usr/bin/trcrpt/trcfmt, cmdtrace, bos43N, 9909A_43N 2/12/99 13:15:34
# COMPONENT_NAME: CMDTRACE system trace logging and reporting facility
#
# FUNCTIONS: template file for trcrpt
#
# ORIGINS: 27, 83
#
# (C) COPYRIGHT International Business Machines Corp. 1988, 1993
# All Rights Reserved
# Licensed Materials - Property of IBM
#
# US Government Users Restricted Rights - Use, duplication or
# disclosure restricted by GSA ADP Schedule Contract with IBM Corp.
#
# LEVEL 1, 5 Years Bull Confidential Information
#
# I. General Information
#
# The formats shown below apply to the data placed into the
# trcrpt format buffer. These formats in general mirror the binary
# format of the data in the trace stream. The exceptions are
# hooks from a 32-bit application on a 64-bit kernel, and hooks from a
# 64-bit application on a 32-bit kernel. These exceptions are noted
# below as appropriate.
#
# Trace formatting templates should not use the thread id or time
# stamp from the buffer. The thread id should be obtained with the
# $TID macro. The time stamp is a raw timer value used by trcrpt to
# calculate the elapsed and delta times. These values are either
# 4 or 8 bytes depending upon the system the trace was run on, not upon
# the environment from which the hook was generated.
# The system environment, 32 or 64 bit, and the hook's
# environment, 32 or 64 bit, are obtained from the $TRACEENV and $HOOKENV
# macros discussed below.
#
# To interpret the time stamp, it is necessary to get the values from
# hook 0x00a, subhook 0x25c, used to convert it to nanoseconds.
# The 3 data words of interest are all 8 bytes in length and are in
# the generic buffer, see the template for hook 00A.
# The first data word gives the multiplier, m, and the second word
```

```

# is the divisor, d. These values should be set to 1 if the
# third word doesn't contain a 2. The nanosecond time is then
# calculated with  $nt = t * m / d$  where t is the time from the trace.
#
# Also, on a 64-bit system, there will be a header on the trace stream.
# This header serves to identify the stream as coming from a
# 64-bit system. There is no such header on the data stream on a
# 32-bit system. This data stream, on both systems, is produced with
# the "-o -" option of the trace command.
# This header consists only of a 4-byte magic number, 0xEFDF1114.
#
# A. Binary format for the 32-bit trace data
# TRCHKL0      MMTDDDDiiiiiii
# TRCHKL0T    MMTDDDDiiiiiiiittttttt
# TRCHKL1      MMTDDDD11111111iiiiiii
# TRCHKL1T    MMTDDDD11111111iiiiiiiittttttt
# Note that trchkg covers TRCHKL2-TRCHKL5.
# trchkg      MMTDDDD1111111122222222333333334444444455555555iiiiiii
# trchkgT     MMTDDDD1111111122222222333333334444444455555555 i... t...
# trcgent     MMTLLLL11111111vvvvvvvvvvvvvvvvvvvvvvvvvvvvvxxxxx i... t...
#
# legend:
# MMM = hook id
# T = hooktype
# D = hookdata
# i = thread id, 4 bytes on a 32 byte system and 8 bytes on a 64-bit
# system. The thread id starts on a 4 or 8 byte boundary.
# t = timestamp, 4 bytes on a 32-bit system or 8 on a
# 64-bit system.
# 1 = d1 (see trchkid.h for calling syntax for the tracehook routines)
# 2 = d2, etc.
# v = trcgen variable length buffer
# L = length of variable length data in bytes.
#
# The DATA_POINTER starts at the third byte in the event, ie.,
# at the 16 bit hookdata DDDD.
# The trcgen() is an exception. The DATA_POINTER starts at
# the fifth byte, ie., at the 'd1' parameter 11111111.
#
# Note that a generic trace hook with a hookid of 0x00b is
# produced for 64-bit data traced from a 64-bit app running on
# a 32-bit kernel. Since this is produced on a 32-bit system, the
# thread id and time stamp will be 4 bytes in the data stream.
#
# B. 64-bit trace hook format
#
# TRCHK64L0 ffff1111hhhhssss iiiiiiiiiiiiii
# TRCHK64L0T ffff1111hhhhssss iiiiiiiiiiiiii tttttttttttttt
# TRCHK64L1 ffff1111hhhhssss 1111111111111111 i...
# ...
# TRCGEN ffff1111hhhhssss ddddddddddddd "string" i...
# TRCGENT ffff1111hhhhssss ddddddddddddd "string" i... t...
#
# Legend
# f - flags
# tgbuuuuuuuuuuuuuuu: t - time stamped, g - generic (trcgen),
# b - 32-bit data, u - unused.
# l - length, number of bytes traced.
# For TRCHKL0 l111 = 0,
# for TRCHKL5T l111 = 40, 0x28 (5 8-byte words)

```

```

#     h - hook id
#     s - subhook id
#     l - data word 1, ...
#     d - generic trace data word.
#     i - thread id, 8 bytes on a 64-bit system, 4 on a 32-bit system.
#         The thread id starts on an 8-byte boundary.
#     t - time stamp, 8 bytes on a 64-bit system, 4 on a 32-bit system.
#
# For non-generic entries, the data pointer starts at the
# subhook id, offset 6. This is compatible with the 32-bit
# hook format shown above.
# For generic (trcgen) hooks, the g flag above is on. The
# length shows the number of variable bytes traced and does not include
# the data word.
# The data pointer starts at the 64-bit data word.
# Note that the data word is 64 bits here.
#
# C. Trace environments
# The trcrpt, trace report, utility must be able to tell whether
# the trace it's formatting came from a 32 or a 64 bit system.
# This is accomplished by the log file header's magic number.
# In addition, we need to know whether 32 or 64 bit data was traced.
# It is possible to run a 32-bit application on a 64-bit kernel,
# and a 64-bit application on a 32-bit kernel.
# In the case of a 32-bit app on a 64-bit kernel, the "b" flag
# shown under item B above is set on. The trcrpt program will
# then present the data as if it came from a 32-bit kernel.
# In the second case, if the reserved hook id 00b is seen, the data
# traced by the 32-bit kernel is made to look as if it came
# from a 64-bit trace. Thus the templates need not be kernel aware.
#
# For example, if a 32-bit app uses
# TRCHKL5T(0x50000005, 1, 2, 3, 4, 5)
# and is running on a 64-bit kernel, the data actually traced
# will look like:
#     ffff1111hhhhssss 1111111111111111 2222222222222222 3333333333333333
#     a000001450000005 0000000100000002 0000000300000004 0000000500000000 i t
# Here, the flags have the T and B bits set (a000) which says
# the hook is timestamped and from a 32-bit app.
# The length is 0x14 bytes, 5 4-byte registers 00000001 through
# 00000005.
# The hook id is 0x5000.
# The subhook id is 0x0005.
# i and t refer to the 8-byte thread id and time stamp.
#
# This would be reformatted as follows before being processed
# by the corresponding template:
#     500e0005 00000001 00000002 00000003 00000004 00000005
# Note this is how the data would look if traced on a 32-bit kernel.
# Note also that the data would be followed by an 8-byte thread id and
# time stamp.
#
# Similarly, consider the following hook traced by a 64-bit app
# on a 32-bit kernel:
#     TRCHKL5T(0x50000005, 1, 2, 3, 4, 5)
# The data traced would be:
#     00b8002c 80000028 50000005 0000000000000001 ... 0000000000000005 i t
# Note that this is a generic trace entry, T = 8.
# In the generic entry, we're using the 32-bit data word for the flags
# and length.

```

```

# The trcrpt utility would reformat this before processing by
# the template as follows:
# 8000002850000005 0000000000000001 ... 0000000000000005 i8 t8
#
# The thread id and time stamp in the data stream will be 4 bytes,
# because the hook came from a 32-bit system.
#
# If a 32-bit app traces generic data on a 64-bit kernel, the b
# bit will be set on in the data stream, and the entry will be formatted
# like it came from a 32-bit environment, i.e. with a 32-bit data word.
# For the case of a 64-bit app on a 32-bit kernel, generic trace
# data is handled in the same manner, with the flags placed
# into the data word.
# For example, if the app issues
# TRCGEN(1, 0x50000005, 1, 6, "hello")
# The 32-bit kernel trace will generate
# 00b00012 40000006 50000005 0000000000000001 "hello"
# This will be reformatted by trcrpt into
# 4000000650000005 0000000000000001 "hello"
# with the data pointer starting at the data word.
#
# Note that the string "hello" could have been 4096 bytes. Therefore
# this generic entry must be able to violate the 4096 byte length
# restriction.
#
# D. Indentation levels
# The left margin is set per template using the 'L=XXXX' command.
# The default is L=KERN, the second column.
# L=APPL moves the left margin to the first column.
# L=SVC moves the left margin to the second column.
# L=KERN moves the left margin to the third column.
# L=INT moves the left margin to the fourth column.
# The command if used must go just after the version code.
#
# Example usage:
#113 1.7 L=INT "stray interrupt" ... \
#
# E. Continuation code and delimiters.
# A '\' at the end of the line must be used to continue the template
# on the next line.
# Individual strings (labels) can be separated by one or more blanks
# or tabs. However, all whitespace is squeezed down to 1 blank on
# the report. Use '\t' for skipping to the next tabstop, or use
# A0.X format (see below) for variable space.
#
# II. FORMAT codes
#
# A. Codes that manipulate the DATA_POINTER
# Gm.n
# "Goto" Set DATA_POINTER to byte.bit location m.n
#
# Om.n
# "Omit" Advance DATA_POINTER by m.n byte.bits
#
# Rm
# "Reverse" Decrement DATA_POINTER by m bytes. R0 byte aligns.
#
# Wm
# Position DATA_POINTER at word m. The word size is either 4 or 8

```

```

# bytes, depending upon whether or not this is a 32 or 64 bit format
# trace. This bares no relation to the %W format specifier.
#
# B. Codes that cause data to be output.
# Am.n
# Left justified ascii.
# m=length in bytes of the binary data.
# n=width of the displayed field.
# The data pointer is rounded up to the next byte boundary.
# Example
# DATA_POINTER|
#           V
#          xxxxxhello world\0xxxxxx
#
# i. A8.16 results in:                |hello wo      |
# DATA_POINTER-----|
#                   V
#          xxxxxhello world\0xxxxxx
#
# ii. A16.16 results in:              |hello world   |
# DATA_POINTER-----|
#                   V
#          xxxxxhello world\0xxxxxx
#
# iii. A16 results in:                |hello world|
# DATA_POINTER-----|
#                   V
#          xxxxxhello world\0xxxxxx
#
# iv. A0.16 results in:              |                |
# DATA_POINTER|
#           V
#          xxxxxhello world\0xxxxxx
#
# Sm (m = 1, 2, 4, or 8)
# Left justified ascii string.
# The length of the string is in the first m bytes of
# the data. This length of the string does not include these bytes.
# The data pointer is advanced by the length value.
# SW specifies the length to be 4 or 8 bytes, depending upon whether
# this is a 32 or 64 bit hook.
# Example
# DATA_POINTER|
#           V
#          xxxxxBhello worldxxxxxx (B = hex 0x0b)
#
# i. S1 results in:                   |hello world|
# DATA_POINTER-----|
#                   V
#          xxxxBhello worldxxxxxx
#
# $reg%S1
# A register with the format code of 'Sx' works "backwards" from
# a register with a different type. The format is Sx, but the length
# of the string comes from $reg instead of the next n bytes.
#
# Bm.n
# Binary format.
# m = length in bytes
# n = length in bits

```



```

# The length in bits of the data is  $m * 8 + n$ . B2.3 and B0.19 are the same.
# Unlike the other printing FORMAT codes, the DATA_POINTER
# can be bit aligned and is not rounded up to the next byte boundary.
#
# Xm
# Hex format.
# m = length in bytes. m=0 thru 16
# X0 is the same as X1, except that no trailing space is output after
# the data. Therefore X0 can be used with a LOOP to output an
# unbroken string of data.
# The DATA_POINTER is advanced by m (1 if m = 0).
# XW will format either 4 or 8 bytes of data depending upon whether
# this is a 32 or 64 bit hook. The DATA_POINTER is advanced
# by 4 or 8 bytes.
#
# Dm (m = 2, 4, or 8)
# Signed decimal format.
# The length of the data is m bytes.
# The DATA_POINTER is advanced by m.
# DW will format either 4 or 8 bytes of data depending upon whether
# this is a 32 or 64 bit hook. The DATA_POINTER is advanced
# by 4 or 8 bytes.
#
# Um (m = 2, 4, or 8)
# Unsigned decimal format.
# The length of the data is m bytes.
# The DATA_POINTER is advanced by m.
# UW will format either 4 or 8 bytes of data depending upon whether
# this is a 32 or 64 bit hook. The DATA_POINTER is advanced
# by 4 or 8 bytes.
#
# om (m = 2, 4, or 8)
# Octal format.
# The length of the data is m bytes.
# The DATA_POINTER is advanced by m.
# ow will format either 4 or 8 bytes of data depending upon whether
# this is a 32 or 64 bit hook. The DATA_POINTER is advanced
# by 4 or 8 bytes.
#
# F4
# Floating point format. (like %0.4E)
# The length of the data is 4 bytes.
# The format of the data is that of C type 'float'.
# The DATA_POINTER is advanced by 4.
#
# F8
# Floating point format. (like %0.4E)
# The length of the data is 8 bytes.
# The format of the data is that of C type 'double'.
# The DATA_POINTER is advanced by 8.
#
# HB
# Number of bytes in trcgen() variable length buffer.
# The DATA_POINTER is not changed.
#
# HT
# 32-bit hooks:
# The hooktype. (0 - E)
# trcgen = 0, trchk = 1, trchl = 2, trchkg = 6
# trcgent = 8, trchkt = 9, trchlt = A, trchkgt = E

```

```

# HT & 0x07 masks off the timestamp bit
# This is used for allowing multiple, different trchhook() calls with
# the same template.
# The DATA_POINTER is not changed.
# 64-bit hooks
# This is the flags field.
# 0x8000 - hook is time stamped.
# 0x4000 - This is a generic trace.
#
# Note that if the hook was reformatted as discussed under item
# I.C above, HT is set to reflect the flags in the new format.
#
# C. Codes that interpret the data in some way before output.
# Tm (m = 4, or 8)
# Output the next m bytes as a data and time string,
# in GMT timezone format. (as in ctime(&seconds))
# The DATA_POINTER is advanced by m bytes.
# Only the low-order 32-bits of the time are actually used.
# TW will format either 4 or 8 bytes of data depending upon whether
# this is a 32 or 64 bit hook. The DATA_POINTER is advanced
# by 4 or 8 bytes.
#
# Em (m = 1, 2, 4, or 8)
# Output the next m bytes as an 'errno' value, replacing
# the numeric code with the corresponding #define name in
# /usr/include/sys/errno.h
# The DATA_POINTER is advanced by 1, 2, 4, or 8.
# EW will format either 4 or 8 bytes of data depending upon whether
# this is a 32 or 64 bit hook. The DATA_POINTER is advanced
# by 4 or 8 bytes.
#
# Pm (m = 4, or 8)
# Use the next m bytes as a process id (pid), and
# output the pathname of the executable with that process id.
# Process ids and their pathnames are acquired by the trace command
# at the start of a trace and by trcrpt via a special EXEC tracehook.
# The DATA_POINTER is advanced by 4 or 8 bytes.
# PW will format either 4 or 8 bytes of data depending upon whether
# this is a 32 or 64 bit hook.
#
# \t
# Output a tab. \t\t\t outputs 3 tabs. Tabs are expanded to spaces,
# using a fixed tabstop separation of 8. If L=0 indentation is used,
# the first tabstop is at 3.
#
# \n
# Output a newline. \n\n\n outputs 3 newlines.
# The newline is left-justified according to the INDENTATION LEVEL.
#
# $macro
# Undefined macros have the value of 0.
# The DATA_POINTER is not changed.
# An optional format can be used with macros:
# $v1%X8 will output the value $v1 in X8 format.
# $zz%B0.8 will output the value $v1 in 8 bits of binary.
# Understood formats are: X, D, U, B and W. Others default to X2.
#
# The W format is used to mask the register.
# Wm.n masks off all bits except bits m through n, then shifts the
# result right m bits. For example, if $ZZ = 0x12345678, then

```

```

#   $zz%W24.27 yields 2. Note the bit numbering starts at the right,
#   with 0 being the least significant bit.
#
# "string"      'string' data type
#   Output the characters inside the double quotes exactly. A string
#   is treated as a descriptor. Use "" as a NULL string.
#
# `string format $macro` If a string is backquoted, it is expanded
#   as a quoted string, except that FORMAT codes and $registers are
#   expanded as registers.
#
# III. SWITCH statement
#   A format code followed by a comma is a SWITCH statement.
#   Each CASE entry of the SWITCH statement consists of
#     1. a 'matchvalue' with a type (usually numeric) corresponding to
#        the format code.
#     2. a simple 'string' or a new 'descriptor' bounded by braces.
#        A descriptor is a sequence of format codes, strings, switches,
#        and loops.
#     3. and a comma delimiter.
#   The switch is terminated by a CASE entry without a comma delimiter.
#   The CASE entry selected is the first entry whose matchvalue
#   is equal to the expansion of the format code.
#   The special matchvalue '*' is a wildcard and matches anything.
#   The DATA_POINTER is advanced by the format code.
#
# IV. LOOP statement
#   The syntax of a 'loop' is
#   LOOP format_code { descriptor }
#   The descriptor is executed N times, where N is the numeric value
#   of the format code.
#   The DATA_POINTER is advanced by the format code plus whatever the
#   descriptor does.
#   Loops are used to output binary buffers of data, so descriptor is
#   usually simply X1 or X0. Note that X0 is like X1 but does not
#   supply a space separator ' ' between each byte.
#
# V. macro assignment and expressions
#   'macros' are temporary (for the duration of that event) variables
#   that work like shell variables.
#   They are assigned a value with the syntax:
#   {{ $xxx = EXPR }}
#   where EXPR is a combination of format codes, macros, and constants.
#   Allowed operators are + - / *
#   For example:
#   {{{ $dog = 7 + 6 }} {{ $cat = $dog * 2 }} $dog $cat
#
#   will output:
#   #000D 001A
#
#   Macros are useful in loops where the loop count is not always
#   just before the data:
#   #G1.5 {{ $count = B0.5 }} G11 LOOP $count {X0}
#
#   Up to 255 macros can be defined per template.
#
# VI. Special macros:

```

```

# $HOOKENV      This is either "32" or "64" depending upon
#               whether this is a 32 or 64 bit trace hook.
#               This can be used to interpret the HT value.
# $TRACEENV    This is either "32" or "64" depending upon
#               whether this is a 32 or 64 bit trace, i.e., whether the
#
#               trace was generated by a 32 or 64 bit kernel.
#               Since hooks will be formatted according to the environment
#               they came from, $HOOKENV should normally be used.
# $RELLINENO   line number for this event. The first line starts at 1.
# $D1 - $D5   dataword 1 through dataword 5. No change to datapointer.
#             The data word is either 4 or 8 bytes.
# $L1 - $L5   Long dataword 1,5(64 bits). No change to datapointer.
# $HD         hookdata (lower 16 bits)
#             For a 32-bit generic hook, $HD is the length of the
#             generic data traced.
#             For 32 or 64 bit generic hooks, use $HL.
# $HL         Hook data length. This is the length in bytes of the hook
#             data. For generic entries it is the length of the
#             variable length buffer and doesn't include the data word.
# $WORDSIZE   Contains the word size, 4 or 8 bytes, of the current
#             entry, (i.e.) $HOOKENV / 8.
# $GENERIC    specifies whether the entry is a generic entry. The
#             value is 1 for a generic entry, and 0 if not generic.
#             $GENERIC is especially useful if the hook can come from
#             either a 32 or 64 bit environment, since the types (HT)
#             have different formats.
# $TOTALCPUS  Output the number of CPUs in the system.
# $TRACEDCPUS Output the number of CPUs that were traced.
# $REPORTEDCPUS Output the number of CPUs active in this report.
#             This can decrease as CPUs stop tracing when, for example,
#             the single-buffer trace, -f, was used and the buffers for
#             each CPU fill up.
# $LARGEDATATYPES This is set to 1 if the kernel is supporting large data
#             types for 64-bit applications.
# $SVC        Output the name of the current SVC
# $EXECPATH   Output the pathname of the executable for current process.
# $PID        Output the current process id.
# $TID        Output the current thread id.
# $CPUID      Output the current processor id.
# $PRI        Output the current process priority
# $ERROR      Output an error message to the report and exit from the
#             template after the current descriptor is processed.
#             The error message supplies the logfile, logfile offset of the
#             start of that event, and the traceid.
# $LOGIDX     Current logfile offset into this event.
# $LOGIDX0    Like $LOGIDX, but is the start of the event.
# $LOGFILE    Name of the logfile being processed.
# $TRACEID    Traceid of this event.
# $DEFAULT    Use the DEFAULT template 008
# $STOP       End the trace report right away
# $BREAK      End the current trace event
# $SKIP       Like break, but don't print anything out.
# $DATAPOINTER The DATA_POINTER. It can be set and manipulated
#             like other user-macros.
#             {{ $DATAPOINTER = 5 }} is equivalent to G5
#
# Note: For generic trace hooks, $DATAPOINTER points to the
# data word. This means it is 0x4 for 32-bit hooks, and 0x8 for
# 64-bit hooks.

```

```

# For non-generic hooks, $DATAPOINTER is set to 2 for 32-bit hooks
# and to 6 for 64 bit trace hooks. This means it always
# points to the subhook id.
#
# $BASEPOINTER Usually 0. It is the starting offset into an event. The actual
# offset is the DATA_POINTER + BASE_POINTER. It is used with
# template subroutines, where the parts on an event have the
# same structure, and can be printed by the same template, but
# might have different starting points into an event.
# $IPADDR IP address of this machine, 4 bytes.
# $BUFF Buffer allocation scheme used, 1=kernel heap, 2=separate segment.
#
# VII. Template subroutines
# If a macro name consists of 3 hex digits, it is a "template subroutine".
# The template whose traceid equals the macro name is inserted in place
# of the macro.
#
# The data pointer is where it was when the template
# substitution was encountered. Any change made to the data pointer
# by the template subroutine remains in affect when the template ends.
#
# Macros used within the template subroutine correspond to those in the
# calling template. The first definition of a macro in the called template
# is the same variable as the first in the called. The names are not
# related.
#
# NOTE: Nesting of template subroutines is supported to 10 levels.
#
# Example:
# Output the trace label ESDI STRATEGY.
# The macro '$stat' is set to bytes 2 and 3 of the trace event.
# Then call template 90F to interpret a buf header. The macro '$return'
# corresponds to the macro '$rv', because they were declared in the same
# order. A macro definition with no '=' assignment just declares the name
# like a place holder. When the template returns, the saved special
# status word is output and the returned minor device number.
#
#900 1.0 "ESDI STRATEGY" {{ $rv = 0 }} {{ $stat = X2 }} \
# $90F \n\
#special_esdi_status=$stat for minor device $rv
#
#90F 1.0 "" G4 {{ $return }} \
# block number X4 \n\
# byte count X4 \n\
# B0.1, 1 B_FLAG0 \
# B0.1, 1 B_FLAG1 \
# B0.1, 1 B_FLAG2 \
# G16 {{ $return = X2 }}
#
#
# Note: The $DEFAULT reserved macro is the same as $008
#
# VIII. BITFLAGS statement
# The syntax of a 'bitflags' is
# BITFLAGS [format_code|register],
# flag_value string {optional string if false}, or
# '&' mask field_value string,
# ...
#
# This statement simplifies expanding state flags, because it looks

```

```

#      a lot like a series of #defines.
#      The '&' mask is used for interpreting bit fields.
#      The mask is anded to the register and the result is compared to
#      the field_value. If a match, the string is printed.
#      The base is 16 for flag_values and masks.
#      The DATA_POINTER is advanced if a format code is used.
#      Note: the default base for BITFLAGS is 16. If the mask or field value
#      has a leading "o", the number is octal. 0x or 0X makes the number hexadecimal.

```

## Examples of Coding Events and Formatting Events

There are five basic steps involved in generating a trace from your software program.

**Step 1: Enable the trace:** Enable and disable the trace from your software that has the trace hooks defined. The following code shows the use of trace events to time the running of a program loop.

```

#include <sys/trcctl.h>
#include <sys/trcmacros.h>
#include <sys/trchkid.h>

char *ctl_file = "/dev/systrctl";
int   ctldfd;
int   i;

main()
{
    printf("configuring trace collection \n");
    if (trcstart("-ad")){
        perror("trcstart");
        exit(1);
    }
    printf("turning trace on \n");
    if(trcon(0)){
        perror("TRCON");
        exit(1);
    }
    /* here is the code that is being traced */
    for(i=1;i<11;i++){
        TRCHKLIT(HKWD_USER1,i);
        /* sleep(1) */
        /* you can uncomment sleep to make the loop
        /* take longer. If you do, you will want to
        /* filter the output or you will be */
        /* overwhelmed with 11 seconds of data */
    }
    /* stop tracing code */
    printf("turning trace off\n");
    if(trcstop(0)){
        perror("TRCOFF");
        exit(1);
    }
}

```

**Step 2: Compile your program:** When you compile the sample program, you need to link to the **librts.a** library:

```
cc -o sample sample.c -l rts
```

**Step 3: Run the program:** Run the program. In this case, it can be done with the following command:

```
./sample
```

**Step 4: Add a stanza to the format file:** This provides the report generator with the information to correctly format your file. The report facility does not know how to format the **HKWD\_USER1** event, unless you provide rules in the trace format file.

The following is an example of a stanza for the **HKWD\_USER1** event. The **HKWD\_USER1** event is event ID 010 hexadecimal. You can verify this by looking at the **sys/trchkid.h** header file.

```
# User event HKWD_USER1 Formatting Rules Stanza
# An example that will format the event usage of the sample program
010 1.0 L=APPL "USER EVENT - HKWD_USER1" 02.0 \n\
    "The # of loop iterations =" U4\n\
    "The elapsed time of the last loop = "\
    endtimer(0x010,0x010) starttimer(0x010,0x010)
```

**Note:** When entering the example stanza, do not modify the master format file **/etc/trcfmt**. Instead, make a copy and keep it in your own directory. This allows you to always have the original trace format file available. If you are going to ship your formatting stanzas, the **trcupdate** command is used to add your stanzas to the default trace format file. See the **trcupdate** command in *AIX Version 6.1 Commands Reference, Volume 5* for information about how to code the input stanzas.

**Step 5: Run the format/filter program:** Filter the output report to get only your events. To do this, run the **trcrpt** command:

```
trcrpt -d 010 -t mytrcfmt -0 exec-on -o sample.rpt
```

The formatted trace results are:

ID	PROC	NAME	I	ELAPSED_SEC	DELTA_MSEC	APPL	SYSCALL	KERNEL	INTERRUPT
010	sample			0.000105984	0.105984	USER HOOK 1			
						The data field for the user hook = 1			
010	sample			0.000113920	0.007936	USER HOOK 1			
						The data field for the user hook = 2 [7 usec]			
010	sample			0.000119296	0.005376	USER HOOK 1			
						The data field for the user hook = 3 [5 usec]			
010	sample			0.000124672	0.005376	USER HOOK 1			
						The data field for the user hook = 4 [5 usec]			
010	sample			0.000129792	0.005120	USER HOOK 1			
						The data field for the user hook = 5 [5 usec]			
010	sample			0.000135168	0.005376	USER HOOK 1			
						The data field for the user hook = 6 [5 usec]			
010	sample			0.000140288	0.005120	USER HOOK 1			
						The data field for the user hook = 7 [5 usec]			
010	sample			0.000145408	0.005120	USER HOOK 1			
						The data field for the user hook = 8 [5 usec]			
010	sample			0.000151040	0.005632	USER HOOK 1			
						The data field for the user hook = 9 [5 usec]			
010	sample			0.000156160	0.005120	USER HOOK 1			
						The data field for the user hook = 10 [5 usec]			

## Usage Hints

The following sections provide some examples and suggestions for use of the trace facility.

### Viewing trace Data

Including several optional columns of data in the trace output can cause the output to exceed 80 columns. It is best to view the report on an output device that supports 132 columns. You can also use the **-O 2line=on** option to produce a more narrow report.

### Bracketing Data Collection

Trace data accumulates rapidly. Bracket the data collection as closely around the area of interest as possible. One technique for doing this is to issue several commands on the same command line. For example, the command

```
trace -a; cp /etc/trcfmt /tmp/junk; trcstop
```

The following command captures the total copy command run, and stops the trace when the command finishes:

```
trace -ax "cp /etc/trcfmt /tmp/junk"
```

captures the total execution of the copy command.

**Note:** This example is more educational if the source file is not already cached in system memory. The **trcfmt** file can be in memory if you have been modifying it or producing trace reports. In that case, choose as the source file some other file that is 50 to 100 KB and has not been touched.

## Reading a trace Report

The trace facility displays system activity. It is a useful learning tool to observe how the system actually performs. The previous output is an interesting example to browse. To produce a report of the copy, use the following:

```
trcrpt -0 "exec=on,pid=on" > cp.rpt
```

In the **cp.rpt** file you can see the following activities:

- The exec and page fault activities of the **cp** process.
- The fork, exec, and page fault activities of the **cp** process.
- The opening of the **/etc/trcfmt** file for reading and the creation of the **/tmp/junk** file.
- The successive **read** and **write** subroutines to accomplish the copy.
- The **cp** process becoming blocked while waiting for I/O completion, and the wait process being dispatched.
- How logical volume requests are translated to physical volume requests.
- The files are mapped rather than buffered in traditional kernel buffers. The read accesses cause page faults that must be resolved by the virtual memory manager.
- The virtual memory manager senses sequential access and begins to prefetch the file pages.
- The size of the prefetch becomes larger as sequential access continues.
- The writes are delayed until the file is closed (unless you captured execution of the **sync** daemon that periodically forces out modified pages).
- The disk device driver coalesces multiple file requests into one I/O request to the drive when possible.

## Effective Filtering of the trace Report

The full detail of the trace data might not be required. You can choose specific events of interest to be shown. For example, it is sometimes useful to find the number of times a certain event occurred. Answer the question, "How many opens occurred in the copy example?" First, find the event ID for the **open** subroutine:

```
trcrpt -j | pg
```

You can see that event ID 15b is the open event. Now, process the data from the copy example (the data is probably still in the log file) as follows:

```
trcrpt -d 15b -0 "exec=on"
```

The report is written to standard output and you can determine the number of opens that occurred. If you want to see only the opens that were performed by the **cp** process, run the report command again using:

```
trcrpt -d 15b -p cp -0 "exec=on"
```

This command shows only the opens performed by the **cp** process.

## Trace Event Groups

Combining multiple trace hooks into a trace event group allows all hooks to be turned on or off at once when starting a trace.



Trace event groups should only be manipulated using either the **trcevgrp** command, or SMIT. The **trcevgrp** command allows groups to be created, modified, removed, and listed.

Reserved event groups may not be changed or removed by the **trcevgrp** command. These are generally groups used to perform system support. A reserved event group must be created using the ODM facilities. Such a group will have three attributes as shown below:

```
SWservAt:
    attribute = "(name)_trcgrp"
    default = " "
    value = "(list-of-hooks)"

SWservAt:
    attribute = "(name)_trcgrpdesc"
    default = " "
    value = "description"

SWservAt:
    attribute = "(name)_trcgrptype"
    default = " "
    value = "reserved"
```

The hook IDs must be enclosed in double quotation marks (") and separated by commas.

---

## Memory Overlay Detection System (MODS)

Some of the most difficult types of problems to debug are what are generally called "memory overlays." Memory overlays include the following:

- Writing to memory that is owned by another program or routine
- Writing past the end (or before the beginning) of declared variables or arrays
- Writing past the end (or before the beginning) of dynamically allocated memory
- Writing to or reading from freed memory
- Freeing memory twice
- Calling memory allocation routines with incorrect parameters or under incorrect conditions.

In the kernel environment (including the kernel, kernel extensions, and device drivers), memory overlay problems have been especially difficult to debug because tools for finding them have not been available. Starting with AIX 4.2.1, however, the Memory Overlay Detection System (MODS) helps detect memory overlay problems in the kernel, kernel extensions, and device drivers.

**Note:** This feature does not detect problems in application code; it only monitors kernel and kernel extension code.

### bosdebug command

The **bosdebug** command turns the MODS facility on and off. Only the root user can run the **bosdebug** command.

To turn on the base MODS support, type:

```
bosdebug -M
```

For a description of all the available options, type:

```
bosdebug -?
```

After you have run **bosdebug** with the options you want, run the **bosboot -a** command, then shut down and reboot your system (using the **shutdown -r** command). If you need to make any changes to your **bosdebug** settings, you must run **bosboot -a** and **shutdown -r** again.

## When to use the MODS feature

This feature is useful in the following circumstances:

- When developing your own kernel extensions or device drivers and you want to test them thoroughly.
- When asked to turn this feature on by IBM technical support service to help in further diagnosing a problem that you are experiencing.

## How MODS works

The primary goal of the MODS feature is to produce a dump file that accurately identifies the problem.

MODS works by turning on additional checking to help detect the conditions listed above. When any of these conditions is detected, your system crashes immediately and produces a dump file that points directly at the offending code. (In previous versions, a system dump might point to unrelated code that happened to be running later when the invalid situation was finally detected.)

If your system crashes while the MODS is turned on, then MODS has most likely done its job.

The **xmalloc** subcommand provides details on exactly what memory address (if any) was involved in the situation, and displays mini-tracebacks for the allocation or free records of this memory.

Similarly, the **netm** command displays allocation and free records for memory allocated using the **net\_malloc** kernel service (for example, **mbufs**, **mclusters**, etc.).

You can use these commands, as well as standard crash techniques, to determine exactly what went wrong.

## MODS limitations

There are limitations to the Memory Overlay Detection System. Although it significantly improves your chances, MODS cannot detect all memory overlays. Also, turning MODS on has a variable negative impact (depending on how frequently **xmalloc** or **xmfree** is called) on overall system performance and causes somewhat more memory to be used in the kernel and the network memory heaps. If your system is running at full processor utilization, or if you are already near the maximums for kernel memory usage, turning on the MODS might cause performance degradation and/or system hangs.

Practical experience with the MODS, however, suggests that the great majority of customers can use it with minimal impact to their systems.

## MODS benefits

The following benefits are gained from using the MODS:

- You can more easily test and debug your own kernel extensions and device drivers.
- Difficult problems that previously required multiple attempts to recreate and debug them generally require many fewer such attempts.

---

## Related Information

Software Product Packaging in the *AIX Version 6.1 General Programming Concepts: Writing and Debugging Programs*

Changing or Removing a Paging Space in the *Operating system and device management*

## Commands References

The **dumpctrl** command, **errinstall** command, **errlogger** command, **errmsg** command, **errupdate** command, **extendlv** command in the *AIX Version 6.1 Commands Reference, Volume 2*.

The **livedumpstart** command in the *AIX Version 6.1 Commands Reference, Volume 3*.

The **sysdumpdev** command, **sysdumpstart** command, **trace** command, **trcrpt** command in the *AIX Version 6.1 Commands Reference, Volume 5*.

## Technical References

The **errsave**, **ras\_register**, **ras\_control**, **ldmp\_setupparms**, **livedump** and **dmp\_compspec** kernel services in the *AIX Version 6.1 Technical Reference: Kernel and Subsystems Volume 1*.



---

# Chapter 17. Loadable Authentication Module Programming Interface

---

## Overview

The loadable authentication module interface provides a means for extending identification and authentication (I&A) for new technologies. The interface implements a set of well-defined functions for performing user and group account access and management.

The degree of integration with the system administrative commands is limited by the amount of functionality provided by the module. When all of the functionality is present, the administrative commands are able to create, delete, modify and view user and group accounts.

The security library and loadable authentication module communicate through the `secmethod_table` interface. The `secmethod_table` structure contains a list of subroutine pointers. Each subroutine pointer performs a well-defined operation. These subroutines are used by the security library to perform the operations which would have been performed using the local security database files.

---

## Load Module Interfaces

Each loadable module defines a number of interface subroutines. The interface subroutines which must be present are determined by how the loadable module is to be used by the system. A loadable module may be used to provide identification (account name and attribute information), authentication (password storage and verification) or both. All modules may have additional support interfaces for initializing and configuring the loadable module, creating new user and group accounts, and serializing access to information. This table describes the purpose of each interface. Interfaces may not be required if the loadable module is not used for the purpose of the interface. For example, a loadable module which only performs authentication functions is not required to have interfaces which are only used for identification operations.

Method Interface Types		
Name	Type	Required
<code>method_attrlist</code>	Support	No
<code>method_authenticate</code>	Authentication	No [ 3]
<code>method_chpass</code>	Authentication	Yes
<code>method_close</code>	Support	No
<code>method_commit</code>	Support	No
<code>method_delgroup</code>	Support	No
<code>method_deluser</code>	Support	No
<code>method_getentry</code>	Identification [ 1]	No
<code>method_getgracct</code>	Identification	No
<code>method_getgrgid</code>	Identification	Yes
<code>method_getgrnam</code>	Identification	Yes
<code>method_getgrset</code>	Identification	Yes
<code>method_getgrusers</code>	Identification	No
<code>method_getpasswd</code>	Authentication	No
<code>method_getpwnam</code>	Identification	Yes

Method Interface Types		
Name	Type	Required
method_getpwuid	Identification	Yes
method_lock	Support	No
method_newgroup	Support	No
method_newuser	Support	No
method_normalize	Authentication	No
method_open	Support	No
method_passwdexpired	Authentication [ 2]	No
method_passwdrestrictions	Authentication [ 2]	No
method_putentry	Identification [ 1]	No
method_putgrent	Identification	No
method_putgrusers	Identification	No
method_putpwent	Identification	No
method_unlock	Support	No

**Notes:**

1. Any module which provides a *method\_attrlist()* interface must also provide this interface.
2. Attributes which are related to password expiration or restrictions should be reported by the *method\_attrlist()* interface.
3. If this interface is not provided the *method\_getpasswd()* interface must be provided.

Several of the functions make use of a *table* parameter to select between user, group and system identification information. The *table* parameter has one of the following values:

Identification Table Names	
Value	Description
"user"	The table containing user account information, such as user ID, full name, home directory and login shell.
"group"	The table containing group account information, such as group ID and group membership list.
"system"	The table containing system information, such as user or group account default values.

When a *table* parameter is used by an authentication interface, "user" is the only valid value.

---

## Authentication Interfaces

Authentication interfaces perform password validation and modification. The authentication interfaces verify that a user is allowed access to the system. The authentication interfaces also maintain the authentication information, typically passwords, which are used to authorize user access.

### The *method\_authenticate* Interface

```
int method_authenticate (char *user, char *response,
                        int **reenter, char **message);
```

The *user* parameter points to the requested user. The *response* parameter points to the user response to the previous message or password prompt. The *reenter* parameter points to a flag. It is set to a non-zero value when the contents of the *message* parameter must be used as a prompt and the user's response used as the *response* parameter when this method is re-invoked. The initial value of the reenter flag is zero. The *message* parameter points to a character pointer. It is set to a message which is output to the user when an error occurs or an additional prompt is required.

*method\_authenticate* verifies that a named user has the correct authentication information, typically a password, for a user account.

*method\_authenticate* is called indirectly as a result of calling the authenticate subroutine. The grammar given in the **SYSTEM** attribute normally specifies the name of the loadable authentication module, but it is not required to do so.

*method\_authenticate* returns **AUTH\_SUCCESS** with a *reenter* value of zero on success. On failure a value of **AUTH\_FAILURE**, **AUTH\_UNAVAIL** or **AUTH\_NOTFOUND** is returned. An **AUTH\_PWDMUSTCHANGE** value with a reenter value of zero indicates success, with the requirement that the user must change the password.

## The method\_chpass Interface

```
int method_chpass (char *user, char *oldpassword,  
                  char *newpassword, char **message);
```

The *user* parameter points to the requested user. The *oldpassword* parameter points to the user's current password. The *newpassword* parameter points to the user's new password. The *message* parameter points to a character pointer. It will be set to a message which is output to the user.

*method\_chpass* changes the authentication information for a user account.

*method\_chpass* is called indirectly as a result of calling the chpass subroutine. The security library will examine the **registry** attribute for the user and invoke the *method\_chpass* interface for the named loadable authentication module.

*method\_chpass* returns zero for success or -1 for failure. On failure the *message* parameter should be initialized with a user message.

## The method\_getpasswd Interface

```
char *method_getpasswd (char *user);
```

The *user* parameter points to the requested user.

*method\_getpasswd* provides the encrypted password string for a user account. The encrypted password string consists of two *salt* characters and 11 encrypted password characters. The *crypt* subroutine is used to create this string and encrypt the user-supplied password for comparison.

*method\_getpasswd* is called when *method\_authenticate* would have been called, but is undefined. The result of this call is compared to the result of a call to the *crypt* subroutine using the response to the password prompt. See the description of the *method\_authenticate* interface for a description of the *response* parameter.

*method\_getpasswd* returns a pointer to an encrypted password on success. On failure a **NULL** pointer is returned and the global variable **errno** is set to indicate the error. A value of **ENOSYS** is used when the module cannot return an encrypted password. A value of **EPERM** is used when the caller does not have the required permissions to retrieve the encrypted password. A value of **ENOENT** is used when the requested user does not exist.

## The `method_normalize` Interface

```
int method_normalize (char *longname, char *shortname);
```

The *longname* parameter points to a fully-qualified user name for modules which include domain or registry information in a user name. The *shortname* parameter points to the shortened name of the user, without the domain or registry information.

*method\_normalize* determines the shortened user name which corresponds to a fully-qualified user name. The shortened user name is used for user account queries by the security library. The fully-qualified user name is only used to perform initial authentication.

If the fully-qualified user name is successfully converted to a shortened user name, a non-zero value is returned. If an error occurs a zero value is returned.

## The `method_passwdexpired` Interface

```
int method_passwdexpired (char *user, char **message);
```

The *user* parameter points to the requested user. The *message* parameter points to a character pointer. It will be set to a message which is output to the user.

*method\_passwdexpired* determines if the authentication information for a user account is expired. This method distinguishes between conditions which allow the user to change their information and those which require administrator intervention. A message is returned which provides more information to the user.

*method\_passwdexpired* is called as a result of calling the `passwdexpired` subroutine.

*method\_passwdexpired* returns 0 when the password has not expired, 1 when the password is expired and the user is permitted to change their password and 2 when the password has expired and the user is not permitted to change their password. A value of -1 is returned when an error has occurred, such as the user does not exist.

## The `method_passwdrestrictions` Interface

```
int method_passwdrestrictions (char *user, char *newpassword,  
                               char *oldpassword, char **message);
```

The *user* parameter points to the requested user. The *newpassword* parameter points to the user's new password. The *oldpassword* parameter points to the user's current password. The *message* parameter points to a character pointer. It will be set to a message which is output to the user.

*method\_passwdrestrictions* determines if new password meets the system requirements. This method distinguishes between conditions which allow the user to change their password by selecting a different password and those which prevent the user from changing their password at the present time. A message is returned which provides more information to the user.

*method\_passwdrestrictions* is called as a result of calling the security library subroutine `passwdrestrictions`.

*method\_passwdrestrictions* returns a value of 0 when *newpassword* meets all of the requirements, 1 when the password does not meet one or more requirements and 2 when the password may not be changed. A value of -1 is returned when an error has occurred, such as the user does not exist.

---

## Identification Interfaces

Identification interfaces perform user and group identity functions. The identification interfaces store and retrieve user and group identifiers and account information.



The identification interfaces divide information into three different categories: user, group and system. User information consists of the user name, user and primary group identifiers, home directory, login shell and other attributes specific to each user account. Group information consists of the group identifier, group member list, and other attributes specific to each group account. System information consists of default values for user and group accounts, and other attributes about the security state of the current system.

## The method\_getentry Interface

```
int method_getentry (char *key, char *table, char *attributes[],
                    attrval_t results[], int size);
```

The *key* parameter refers to an entry in the named table. The *table* parameter refers to one of the three tables. The *attributes* parameter refers to an array of pointers to attribute names. The *results* parameter refers to an array of value return data structures. Each value return structure contains either the value of the corresponding attribute or a flag indicating a cause of failure. The *size* parameter is the number of array elements.

*method\_getentry* retrieves user, group and system attributes. One or more attributes may be retrieved for each call. Success or failure is reported for each attribute.

*method\_getentry* is called as a result of calling the *getuserattr*, *getgroupattr* and *getconfattr* subroutines.

*method\_getentry* returns a value of 0 if the *key* entry was found in the named *table*. When the entry does not exist in the table, the global variable **errno** must be set to **ENOENT**. If an error in the value of *table* or *size* is detected, the **errno** variable must be set to **EINVAL**. Individual attribute values have additional information about the success or failure for each attribute. On failure a value of -1 is returned.

## The method\_getgracct Interface

```
struct group *method_getgracct (void *id, int type);
```

The *id* parameter refers to a group name or GID value, depending upon the value of the *type* parameter. The *type* parameters indicates whether the *id* parameter is to be interpreted as a (char \*) which references the group name, or (gid\_t) for the group.

*method\_getgracct* retrieves basic group account information. The *id* parameter may be a group name or identifier, as indicated by the *type* parameter. The basic group information is the group name and identifier. The group member list is not returned by this interface.

*method\_getgracct* may be called as a result of calling the *IDtogroup* subroutine.

*method\_getgracct* returns a pointer to the group's group file entry on success. The group file entry may not include the list of members. On failure a **NULL** pointer is returned.

## The method\_getgrgid Interface

```
struct group *method_getgrgid (gid_t gid);
```

The *gid* parameter is the group identifier for the requested group.

*method\_getgrgid* retrieves group account information given the group identifier. The group account information consists of the group name, identifier and complete member list.

*method\_getgrgid* is called as a result of calling the *getgrgid* subroutine.

*method\_getgrgid* returns a pointer to the group's group file structure on success. On failure a **NULL** pointer is returned.

## The method\_getgrnam Interface

```
struct group *method_getgrnam (char *group);
```

The *group* parameter points to the requested group.

*method\_getgrnam* retrieves group account information given the group name. The group account information consists of the group name, identifier and complete member list.

*method\_getgrnam* is called as a result of calling the *getgrnam* subroutine. This interface may also be called if *method\_getentry* is not defined.

*method\_getgrnam* returns a pointer to the group's group file structure on success. On failure a **NULL** pointer is returned.

## The method\_getgrset Interface

```
char *method_getgrset (char *user);
```

The *user* parameter points to the requested user.

*method\_getgrset* retrieves supplemental group information given a user name. The supplemental group information consists of a comma separated list of group identifiers. The named user is a member of each listed group.

*method\_getgrset* is called as a result of calling the *getgrset* subroutine.

*method\_getgrset* returns a pointer to the user's concurrent group set on success. On failure a **NULL** pointer is returned.

## The method\_getgrusers Interface

```
int method_getgrusers (char *group, void *result,  
                      int type, int *size);
```

The *group* parameter points to the requested group. The *result* parameter points to a storage area which will be filled with the group members. The *type* parameters indicates whether the *result* parameter is to be interpreted as a (*char \*\**) which references a user name array, or (*uid\_t*) array. The *size* parameter is a pointer to the number of users in the named group. On input it is the size of the *result* field.

*method\_getgrusers* retrieves group membership information given a group name. The return value may be an array of user names or identifiers.

*method\_getgrusers* may be called by the security library to obtain the group membership information for a group.

*method\_getgrusers* returns 0 on success. On failure a value of -1 is returned and the global variable **errno** is set. The value **ENOENT** must be used when the requested group does not exist. The value **ENOSPC** must be used when the list of group members does not fit in the provided array. When **ENOSPC** is returned the *size* parameter is modified to give the size of the required *result* array.

## The method\_getpwnam Interface

```
struct passwd *method_getpwnam (char *user);
```

The *user* parameter points to the requested user.

*method\_getpwnam* retrieves user account information given the user name. The user account information consists of the user name, identifier, primary group identifier, full name, login directory and login shell.

*method\_getpwnam* is called as a result of calling the *getpwnam* subroutine. This interface may also be called if *method\_getentry* is not defined.

*method\_getpwnam* returns a pointer to the user's password structure on success. On failure a **NULL** pointer is returned.

## The *method\_getpwuid* Interface

```
struct passwd *method_getpwuid (uid_t uid);
```

The *uid* parameter points to the user ID of the requested user.

*method\_getpwuid* retrieves user account information given the user identifier. The user account information consists of the user name, identifier, primary group identifier, full name, login directory and login shell.

*method\_getpwuid* is called as a result of calling the *getpwuid* subroutine.

*method\_getpwuid* returns a pointer to the user's password structure on success. On failure a **NULL** pointer is returned.

## The *method\_putentry* Interface

```
int method_putentry (char *key, char *table, char *attributes,  
                    attrval_t values[], int size);
```

The *key* parameter refers to an entry in the named table. The *table* parameter refers to one of the three tables. The *attributes* parameter refers to an array of pointers to attribute names. The *values* parameter refers to an array of value structures which correspond to the attributes. Each value structure contains a flag indicating if the attribute was output. The *size* parameter is the number of array elements.

*method\_putentry* stores user, group and system attributes. One or more attributes may be retrieved for each call. Success or failure is reported for each attribute. Values will be saved until *method\_commit* is invoked.

*method\_putentry* is called as a result of calling the *putuserattr*, *putgroupattr* and *putconfattr* subroutines.

*method\_putentry* returns 0 when the attributes have been updated. On failure a value of -1 is returned and the global variable **errno** is set to indicate the cause. A value of **ENOSYS** is used when updating information is not supported by the module. A value of **EPERM** is used when the invoker does not have permission to create the group. A value of **ENOENT** is used when the entry does not exist. A value of **EROFS** is used when the module was not opened for updates.

## The *method\_putgrent* Interface

```
int method_putgrent (struct group *entry);
```

The *entry* parameter points to the structure to be output. The account name is contained in the structure.

*method\_putgrent* stores group account information given a group entry. The group account information consists of the group name, identifier and complete member list. Values will be saved until *method\_commit* is invoked.

*method\_putgrent* may be called as a result of calling the *putgroupattr* subroutine.

*method\_putgrent* returns 0 when the group has been successfully updated. On failure a value of -1 is returned and the global variable **errno** is set to indicate the cause. A value of **ENOSYS** is used when updating groups is not supported by the module. A value of **EPERM** is used when the invoker does not have permission to update the group. A value of **ENOENT** is used when the group does not exist. A value of **EROFS** is used when the module was not opened for updates.

## The *method\_putgrusers* Interface

```
int method_putgrusers (char *group, char *users);
```

The *group* parameter points to the requested group. The *users* parameter points to a **NUL** character separated, double **NUL** character terminated, list of group members.

*method\_putgrusers* stores group membership information given a group name. Values will be saved until *method\_commit* is invoked.

*method\_putgrusers* may be called as a result of calling the *putgroupattr* subroutine.

*method\_putgrusers* returns 0 when the group has been successfully updated. On failure a value of -1 is returned and the global variable **errno** is set to indicate the cause. A value of **ENOSYS** is used when updating groups is not supported by the module. A value of **EPERM** is used when the invoker does not have permission to update the group. A value of **ENOENT** is used when the group does not exist. A value of **EROFS** is used when the module was not opened for updates.

## The *method\_putpwent* Interface

```
int method_putpwent (struct passwd *entry);
```

The *entry* parameter points to the structure to be output. The account name is contained in the structure.

*method\_putpwent* stores user account information given a user entry. The user account information consists of the user name, identifier, primary group identifier, full name, login directory and login shell. Values will be saved until *method\_commit* is invoked.

*method\_putpwent* may be called as a result of calling the *putuserattr* subroutine.

*method\_putpwent* returns 0 when the user has been successfully updated. On failure a value of -1 is returned and the global variable **errno** is set to indicate the cause. A value of **ENOSYS** is used when updating users is not supported by the module. A value of **EPERM** is used when the invoker does not have permission to update the user. A value of **ENOENT** is used when the user does not exist. A value of **EROFS** is used when the module was not opened for updates.

---

## Support Interfaces

Support interfaces perform functions such as initiating and terminating access to the module, creating and deleting accounts, and serializing access to information.

## The *method\_attrlist* Interface

```
attrtab **method_attrlist (void);
```

This interface does not require any parameters.

*method\_attrlist* provides a means of defining additional attributes for a loadable module. Authentication-only modules may use this interface to override attributes which would normally come from the identification module half of a compound load module.

*method\_attrlist* is called when a loadable module is first initialized. The return value will be saved for use by later calls to various identification and authentication functions.

## The `method_close` Interface

```
void method_close (void *token);
```

The *token* parameter is the value of the corresponding *method\_open* call.

*method\_close* indicates that access to the loadable module has ended and all system resources may be freed. The loadable module must not assume this interface will be invoked as a process may terminate without calling this interface.

*method\_close* is called when the session count maintained by *enduserdb* reaches zero.

There are no defined error return values. It is expected that the *method\_close* interface handle common programming errors, such as being invoked with an invalid token, or repeatedly being invoked with the same token.

## The `method_commit` Interface

```
int method_commit (char *key, char *table);
```

The *key* parameter refers to an entry in the named table. If it is **NULL** it refers to all entries in the table. The *table* parameter refers to one of the three tables.

*method\_commit* indicates that the specified pending modifications are to be made permanent. An entire table or a single entry within a table may be specified. *method\_lock* will be called prior to calling *method\_commit*. *method\_unlock* will be called after *method\_commit* returns.

*method\_commit* is called when *putgroupattr* or *putuserattr* are invoked with a *Type* parameter of **SEC\_COMMIT**. The value of the *Group* or *User* parameter will be passed directly to *method\_commit*.

*method\_commit* returns a value of 0 for success. A value of -1 is returned to indicate an error and the global variable **errno** is set to indicate the cause. A value of **ENOSYS** is used when the load module does not support modification requests for any users. A value of **EROFS** is used when the module is not currently opened for updates. A value of **EINVAL** is used when the *table* parameter refers to an invalid table. A value of **EIO** is used when a potentially temporary input-output error has occurred.

## The `method_delgroup` Interface

```
int method_delgroup (char *group);
```

The *group* parameter points to the requested group.

*method\_delgroup* removes a group account and all associated information. A call to *method\_commit* is not required. The group will be removed immediately.

*method\_delgroup* is called when *putgroupattr* is invoked with a *Type* parameter of **SEC\_DELETE**. The value of the *Group* and *Attribute* parameters will be passed directly to *method\_delgroup*.

*method\_delgroup* returns 0 when the group has been successfully removed. On failure a value of -1 is returned and the global variable **errno** is set to indicate the cause. A value of **ENOSYS** is used when deleting groups is not supported by the module. A value of **EPERM** is used when the invoker does not have permission to delete the group. A value of **ENOENT** is used when the group does not exist. A value of **EROFS** is used when the module was not opened for updates. A value of **EBUSY** is used when the group has defined members.

## The method\_deluser Interface

```
int method_deluser (char *user);
```

The *user* parameter points to the requested user.

*method\_delgroup* removes a user account and all associated information. A call to *method\_commit* is not required. The user will be removed immediately.

*method\_deluser* is called when *putuserattr* is invoked with a *Type* parameter of **SEC\_DELETE**. The value of the *User* and *Attribute* parameters will be passed directly to *method\_deluser*.

*method\_deluser* returns 0 when the user has been successfully removed. On failure a value of -1 is returned and the global variable **errno** is set to indicate the cause. A value of **ENOSYS** is used when deleting users is not supported by the module. A value of **EPERM** is used when the invoker does not have permission to delete the user. A value of **ENOENT** is used when the user does not exist. A value of **EROFS** is used when the module was not opened for updates.

## The method\_lock Interface

```
void *method_lock (char *key, char *table, int wait);
```

The *key* parameter refers to an entry in the named table. If it is **NULL** it refers to all entries in the table. The *table* parameter refers to one of the three tables. The *wait* parameter is the number of second to wait for the lock to be acquired. If the *wait* parameter is zero the call returns without waiting if the entry cannot be locked immediately.

*method\_lock* informs the loadable modules that access to the underlying mechanisms should be serialized for a specific table or table entry.

*method\_lock* is called by the security library when serialization is required. The return value will be saved and used by a later call to *method\_unlock* when serialization is no longer required.

## The method\_newgroup Interface

```
int method_newgroup (char *group);
```

The *group* parameter points to the requested group.

*method\_newgroup* creates a group account. The basic group account information must be provided with calls to *method\_putgrent* or *method\_putentry*. The group account information will not be made permanent until *method\_commit* is invoked.

*method\_newgroup* is called when *putgroupattr* is invoked with a *Type* parameter of **SEC\_NEW**. The value of the *Group* parameter will be passed directly to *method\_newgroup*.

*method\_newgroup* returns 0 when the group has been successfully created. On failure a value of -1 is returned and the global variable **errno** is set to indicate the cause. A value of **ENOSYS** is used when creating group is not supported by the module. A value of **EPERM** is used when the invoker does not have permission to create the group. A value of **EEXIST** is used when the group already exists. A value of **EROFS** is used when the module was not opened for updates. A value of **EINVAL** is used when the group has an invalid format, length or composition.

## The method\_newuser Interface

```
int method_newuser (char *user);
```

The *user* parameter points to the requested user.

*method\_newuser* creates a user account. The basic user account information must be provided with calls to *method\_putpwent* or *method\_putentry*. The user account information will not be made permanent until *method\_commit* is invoked.

*method\_newuser* is called when *putuserattr* is invoked with a *Type* parameter of **SEC\_NEW**. The value of the *User* parameter will be passed directly to *method\_newuser*.

*method\_newuser* returns 0 when the user has been successfully created. On failure a value of -1 is returned and the global variable **errno** is set to indicate the cause. A value of **ENOSYS** is used when creating users is not supported by the module. A value of **EPERM** is used when the invoker does not have permission to create the user. A value of **EEXIST** is used when the user already exists. A value of **EROFS** is used when the module was not opened for updates. A value of **EINVAL** is used when the user has an invalid format, length or composition.

## The *method\_open* Interface

```
void *method_open (char *name, char *domain,  
                  int mode, char *options);
```

The *name* parameter is a pointer to the stanza name in the configuration file. The *domain* parameter is the value of the **domain=** attribute in the configuration file. The *mode* parameter is either **O\_RDONLY** or **O\_RDWR**. The *options* parameter is a pointer to the **options=** attribute in the configuration file.

*method\_open* prepares a loadable module for use. The domain and options attributes are passed to *method\_open*.

*method\_open* is called by the security library when the loadable module is first initialized and when *setuserdb* is first called after *method\_close* has been called due to an earlier call to *enduserdb*. The return value will be saved for a future call to *method\_close*.

## The *method\_unlock* Interface

```
void method_unlock (void *token);
```

The *token* parameter is the value of the corresponding *method\_lock* call.

*method\_unlock* informs the loadable modules that an earlier need for access serialization has ended.

*method\_unlock* is called by the security library when serialization is no longer required. The return value from the earlier call to *method\_lock* be used.

---

## Configuration Files

The security library uses the `/usr/lib/security/methods.cfg` file to control which modules are used by the system. A stanza exists for each loadable module which is to be used by the system. Each stanza contains a number of attributes used to load and initialize the module. The loadable module may use this information to configure its operation when the *method\_open()* interface is invoked immediately after the module is loaded.

## The options Attribute

The options attribute will be passed to the loadable module when it is initialized. This string is a comma-separated list of *Flag* and *Flag=Value* entries. The entire value of the *options* attribute is passed to the *method\_open()* subroutine when the module is first initialized. Five pre-defined flags control how the library uses the loadable module.

<b>auth=module</b>	<i>Module</i> will be used to perform authentication functions for the current loadable authentication module. Subroutine entry points dealing with authentication-related operations will use method table pointers from the named module instead of the module named in the <i>program=</i> or <i>program_64=</i> attribute.
<b>authonly</b>	The loadable authentication module only performs authentication operations. Subroutine entry points which are not required for authentication operations, or general support of the loadable module, will be ignored.
<b>db=module</b>	<i>Module</i> will be used to perform identification functions for the current loadable authentication module. Subroutine entry points dealing with identification related operations will use method table pointers from the name module instead of the module named in the <i>program=</i> or <i>program_64=</i> attribute.
<b>dbonly</b>	The loadable authentication module only provides user and group identification information. Subroutine entry points which are not required for identification operations, or general support of the loadable module, will be ignored.
<b>noprompt</b>	The initial password prompt for authentication operations is suppressed. Password prompts are normally performed prior to a call to <i>method_authenticate()</i> . <i>method_authenticate()</i> must be prepared to receive a <b>NULL</b> pointer for the <i>response</i> parameter and set the <i>reenter</i> parameter to <b>TRUE</b> to indicate that the user must be prompted with the contents of the <i>message</i> parameter prior to <i>method_authenticate()</i> being re-invoked. See the description of <i>method_authenticate</i> for more information on these parameters.

---

## Compound Load Modules

Compound load modules are created with the *auth=* and *db=* attributes. The security library is responsible for constructing a new method table to perform the compound function.

Interfaces are divided into three categories: identification, authentication and support. Identification interfaces are used when a compound module is performing an identification operation, such as the *getpwnam()* subroutine. Authentication interfaces are used when a compound module is performing an authentication operation, such as the *authenticate()* subroutine. Support subroutines are used when initializing the loadable module, creating or deleting entries, and performing other non-data operations. The table Method Interface Types describes the purpose of each interface. The table below describes which support interfaces are called in a compound module and their order of invocation.

Support Interface Invocation	
Name	Invocation Order
<i>method_attrlist</i>	Identification, Authentication
<i>method_close</i>	Identification, Authentication
<i>method_commit</i>	Identification, Authentication
<i>method_deluser</i>	Authentication, Identification
<i>method_lock</i>	Identification, Authentication
<i>method_newuser</i>	Identification, Authentication
<i>method_open</i>	Identification, Authentication
<i>method_unlock</i>	Authentication, Identification



## Related Information

Identification and Authentication Subroutines

`/usr/lib/security/methods.cfg` File



---

## Chapter 18. Kernel Storage-Protection Keys

Kernel storage-protection keys provide a storage-protection mechanism for the kernel and kernel extensions.

Storage protection keys have the following elements:

- Each virtual memory page is assigned a small integer protection key, which represents a required access authority to make the page readable or writable in the current context.
- A register, the authority mask register, is used to encode the access authority of the current context.

The hardware can provide up to 32 protection keys. The authority mask register contains a pair of bits representing the write and read access authorities for each of the defined protection keys, making it possible for the kernel to control access to as many as 32 classes of memory, where each class consists of virtual memory pages sharing the same specific protection key.

In all cases, a running thread can modify its authority mask register, and hence its authority to access specific memory classes. In AIX environment, the storage-protection mechanism can catch programming errors caused by inadvertent storage overlays. When you write a program using storage protection keys, the program can voluntarily subdivide its memory for protection purposes, decreasing the likelihood that an accidental reference to memory goes undetected. Such undetected errors can be difficult to debug; however, by coding to detect them as they occur, a program can improve its reliability, availability, and serviceability characteristics.

When you use storage protection keys in the kernel, the following programming errors can be detected and easily repaired:

- Incorrect reference to application memory by the kernel
- Incorrect reference to kernel private memory by key-safe kernel extensions
- Incorrect reference to kernel extension private memory by the kernel
- Incorrect reference to memory controlled by one kernel subsystem by another

Kernel extensions can be classified into the following categories from the storage protection aspect:

### **Key-unsafe kernel extensions**

Key-unsafe kernel extensions are traditional extensions that were written without regard to storage key protection. When loaded in a key-safe environment, such extensions are automatically run with wide access authority to provide binary compatibility. With the ability to access essentially all kernel memory without restrictions, key-unsafe kernel extensions can continue to function, even though they might violate the protection domain of the kernel.

In some other cases key-unsafe kernel extensions can cause system crashes therefore it is highly recommended that kernel extensions should be made key-safe before running. See “Making a Kernel Extension Key Safe” on page 388.

### **Key-safe kernel extensions**

Key-safe kernel extensions do not directly refer to either the internal data structures of the kernel or user space addresses. See “Making a Kernel Extension Key Safe” on page 388.

### **Key-protected kernel extensions**

Key-protected kernel extensions can coexist with the key-protected kernel environment, and also become key-protected by identifying and protecting their own private data. To place kernel extension data in key-protected memory, see “Designing the Key Protection in a Key-protected Kernel Extension” on page 388.

---

## Kernel Keys and Kernel Key Sets

A kernel key is a virtual key that is assigned by the kernel. Kernel programs can use more virtual keys than exist in the hardware, because many kernel keys can share a single hardware key.

The kernel data is classified into keys according to function. You can use kernel keys to define a large number of virtual storage protection keys. Most kernel keys are used only within the kernel. The **sys/skeys.h** file contains a list of all keys. The following kernel keys can be useful in your kernel extensions:

### **KKEY\_PUBLIC**

Used to access stack, bss and data.

### **KKEY\_BLOCK\_DEV**

Required for block device drivers. The **buf** structures must be either public or in the key.

### **KKEY\_COMMO**

Required for communication drivers.

### **KKEY\_NETM**

Required for drivers to reference memory that is allocated by the **net\_malloc** kernel service.

### **KKEY\_USB**

Required for USB device drivers.

### **KKEY\_GRAPHICS**

Required for graphics device drivers.

### **KKEY\_DMA**

Required for direct memory access (DMA) information.

### **KKEY\_TRB**

Required for timer services (the **trb** structure).

### **KKEY\_IOMAP**

Used to access I/O mapped segments.

### **KKEY\_FILE\_SYSTEM**

Used to access vnodes and gnodes (**vnop** callers).

You can use a set of keys for a typical kernel extension using one of the predefined kernel key sets. The **sys/skeys.h** file contains a list of all key sets. The following predefined kernel key sets are useful with read and write accesses to the member kernel keys. Similar sets whose names are appended with **\_READ** grant only read access.

### **KKEYSET\_GLOBAL**

All known kernel keys

### **KKEYSET\_LEGACY**

The large set of kernel keys that is granted to kernel extensions that are not key-safe

### **KKEYSET\_KERNEXT**

The kernel keys common to most kernel extensions

### **KKEYSET\_COMMO**

The common keys plus those that are used by communication drivers

### **KKEYSET\_BLOCK**

The common keys plus those that are used by block I/O drivers

### **KKEYSET\_GRAPHICS**

The common keys plus those that are used by graphics drivers

## **KKEYSET\_USB**

The common keys plus those that are used by USB drivers

## **KKEYSET\_USERDATA**

All known user-space kernel keys

---

## **Protection Gates**

To make a kernel extension key-safe or key-protected, you must add protection gates, typically at all entry and exit points of your module. With proper protection gates, your module has access to the data that it requires, and does not have access to the data that it does not require.

You can specify the following gates at an entry point:

### **An add gate**

With an add gate, you can augment the callers key set with your own. You can specify an add gate for a service where your caller passes pointers to data. The data of the caller might be protected with a key you are unaware of, therefore you must retain these keys to refer to the data of the caller. You can add additional keys so that you can also refer to any private data that you need.

### **A replace gate**

With a replace gate, you can switch to your own key set. You can specify a replace gate at an entry point for a callback that you have registered with, for example, the device switch table. You can also specify the replace gate to relinquish the keys of the caller, when the kernel is your caller. (Do not retain access to internal data of the kernel.) You can use the predefined key sets as described in “Kernel Keys and Kernel Key Sets” on page 386 to form the basis of a typical replace gate.

In both cases, the protection gate service returns the original authority mask register value to you, so that you can restore keys at your exit points.

You can decide where to place these gates in your program flow. You can identify and place gates at all entry points that are externally visible. One common exception is to defer the gate when you take advantage of keys of the caller to pick up potentially private data being passed in. Then switch to your own key set with a replace gate. This technique provides better storage protection than adding an add gate right at the entry point.

**Important:** Remember to restore the necessary keys at any exit point, where data might be stored back through a parameter pointer.

You can use the following kernel services to implement your protection gates:

### **hkeyset\_add**

Adds access rights to the currently addressable hardware key set, and returns the old hardware key set.

### **hkeyset\_replace**

Loads a new set of access rights, and returns the old hardware key set.

### **hkeyset\_restore**

Restores a saved set of access rights.

### **hkeyset\_get**

Returns the current access rights.

To identify the entry points of your kernel extension, look for the following typical ones:

- Device-switch table callbacks for open, close, read, write, ioctl, strategy, select, print, dump, mpx, and revoke entry points

**Tip:** Typically, the config entry point does not need a protection gate, but includes the initialization necessary for protection gates, heaps and so on, for subsequent uses by other entry points.

- The timer event handler
- The watchdog handler
- The enhanced I/O error handling handlers
- The interrupt handler (INTCLASSx, iodone, offlevel)
- The environmental and power warning handler
- Exported system calls
- The dynamic reconfiguration and high availability handlers
- The shutdown notification handler
- The RAS callback
- Streams callback functions
- Process state change notification handlers
- Function pointers that are passed outside of your module

---

## Making a Kernel Extension Key Safe

To make a kernel extension key-safe, follow these steps:

1. Decide which kernel key set, if any, can be the basis for the key set of your module.
2. Optionally, remove any unnecessary keys from your copy of the kernel key set.
3. Convert the kernel key set to a hardware key set.
4. Place add and replace protection gates at or near all entry points (except initialization). See “Protection Gates” on page 387.
5. Place restore gates at or near exit points.
6. Link your extension with the new **-b ras** flag to identify the extension to the system as reliability, availability, and serviceability aware.

**Restriction:** You must specify **-q noinlgue** to ensure that the compiler does not generate inline pointer glue.

---

## Designing the Key Protection in a Key-protected Kernel Extension

To design key protection into a kernel extension, perform the following steps:

1. Define memory classes. A memory class is a set of virtual pages that are associated with the same storage protection key. Ideally, each memory class is represented by a unique kernel key, although there are limits on both the actual number of kernel keys that can be used, and the actual number of hardware keys that are available for the kernel and kernel extensions.
2. Allocate a kernel key for each memory class. You can use the **kkey\_assign\_private** kernel service to allocate unique kernel keys.
3. Assign kernel keys to the (pages of) memory to be included in each memory class. You can use the **heap\_create** kernel service to create a memory heap whose pages are protected with a specified kernel key. Your new heap can then be used by the **xmalloc** and **xmfree** kernel services as usual. When your program ends, use the **heap\_destroy** kernel service to clean up any heaps you have created.
4. Determine which kernel keys are required by the kernel extension to satisfy the requirements of other kernel extension data that is referred to. For the best result, add the access rights from one of these predefined key sets to your own, rather than build up your key set from scratch. See “Kernel Keys and Kernel Key Sets” on page 386.
5. Create the kernel key sets defining the access authorities that are needed by your extension to refer to its private memory classes. You can use the following kernel services:

**kkeyset\_create**

Allocates an empty key set.

**kkeyset\_add\_key**

Adds access rights for a kernel key to your key set.

**kkeyset\_remove\_key**

Removes access rights to a given kernel key from your key set.

**kkeyset\_add\_set**

Adds the access rights that are defined in one key set to another.

**kkeyset\_remove\_set**

Removes the access rights that are defined in one key set from another.

**kkeyset\_delete**

Frees a key set after it has been used to create a hardware key set by 6.

6. Create hardware key sets that you need to control access to the memory classes during execution. These hardware key sets are needed to implement your protection gates. The **kkeyset\_to\_hkeyset** service performs this translation.

**Attention:** The **HKEYSET\_GLOBAL** key set is defined to access any memory no matter which key protects it. Therefore, this hardware key set should be used with caution.

7. Implement protection gates as needed to grant access to the memory classes, and to restore the prior authority mask register state. Typically, you can perform these steps at the entry and exit points of the major functions of the kernel extension. You can implement a replace gate at an entry point in your driver to grant access appropriately. Without one, the driver has the inappropriate access rights of its caller, typically the kernel. You can implement a restore gate at exit points to restore the access rights back to those of the calling program. See “Protection Gates” on page 387 for more information about protection gates.
8. Link the kernel extension using the new **-b ras** flag, so that the kernel can recognize it as a key-safe kernel extension.

If the extension must refer to user space memory, use the **copyin**, **copyout**, or one of the cross-memory kernel services. These services acquire the appropriate additional hardware keys necessary to access and protect user space. (The user-space application might use storage protect keys itself.) It is possible, but not suggested, to add user keys to the current authority mask register with the **hkeyset\_update\_userkeys** service, and to restore the previous state with the **hkeyset\_restore\_userkeys** service.

Any kernel keys can be used, with appropriate care, by any kernel extension. You can also remove keys that you do not need. For example, you can remove the **KKEY\_TRB** key if your extension does not use services such as the **tstart** kernel service, which uses the **trb** structure. Though not required, the action can increase the benefit of using storage protect keys by minimizing the access rights of your kernel extension. Start with a predefined key set, and remove specific keys; rather than start with an empty key set, and figure out which ones you need to add.

In general, you do not need to test your code for the presence, absence, or number of available hardware keys. All key-protection kernel services work efficiently in all cases.

## Examples

1. The following example creates a pinned, shared heap that is assigned to a specified kernel key.

```
#include <sys/malloc.h>
#include <sys/skeys.h>
kkey_t my_kkey;
heapattr_t heapattr;
heapaddr_t my_heap = HPA_INVALID_HEAP

rc = kkey_assign_private("my string", 0, 0, &my_kkey);
bzero(&heapattr, sizeof(heapattr));
```

```

heapattr.hpa_eyec = EYEC_HEAPATTR;
heapattr.hpa_version = HPA_VERSION;
heapattr.hpa_flags = HPA_PINNED | HPA_SHARED;
heapattr.hpa_debug_level = HPA_DEFAULT_DEBUG;
heapattr.hpa_kkey = my_kkey;
rc = heap_create(&heapattr, &my_heap);

```

## 2. The following examples create a hardware key set.

- Initialize a hardware key set (errors ignored for brevity)

```

#include <sys/skeys.h>
hkeyset_t my_hkeyset; /* globally visible hardware keyset */
kkeyset_t my_kkeyset; /* temporary kernel keyset */

```

```

my_kkeyset = KKEYSET_INVALID;
rc = kkeyset_create(&my_kkey);
rc = kkeyset_add_set(my_kkeyset, KKEYSET_BLOCK);
rc = kkeyset_add_key(my_kkeyset, my_kkey, KA_RW);
rc = kkeyset_to_hkeyset(my_kkeyset, &my_hkeyset);
kkeyset_delete(my_kkeyset);

```

- Implement a protection gate

```

hkeyset_t old_hkeyset;
old_hkeyset = hkeyset_replace(my_hkeyset);
... ..
hkeyset_restore(old_hkeyset);

```

---

## Related Information

The **kkey\_assign\_private**, **kkeyset\_create**, **kkeyset\_add\_key**, **kkeyset\_remove\_key**, **kkeyset\_add\_set**, **kkeyset\_remove\_set**, **kkeyset\_delete**, and **kkeyset\_to\_hkeyset** kernel services.

The **hkeyset\_add**, **hkeyset\_replace**, **hkeyset\_restore**, or **hkeyset\_get** kernel service.



---

## Appendix. Notices

This information was developed for products and services offered in the U.S.A.

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 give you any license to these patents. You can send license inquiries, in writing, to:

IBM Director of Licensing  
IBM Corporation  
North Castle Drive  
Armonk, NY 10504-1785  
U.S.A.

**The following paragraph does not apply to the United Kingdom or any other country where such provisions are inconsistent with local law:** 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 states 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.

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 Corporation  
Dept. LRAS/Bldg. 003  
11400 Burnet Road  
Austin, TX 78758-3498  
U.S.A.

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.

For license inquiries regarding double-byte (DBCS) information, contact the IBM Intellectual Property Department in your country or send inquiries, in writing, to:

IBM World Trade Asia Corporation  
Licensing  
2-31 Roppongi 3-chome, Minato-ku  
Tokyo 106, Japan

IBM may use or distribute any of the information you supply in any way it believes appropriate without incurring any obligation to you.

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.

Any references in this information to non-IBM Web sites are provided for convenience only and do not in any manner serve as an endorsement of those Web sites. The materials at those Web sites are not part of the materials for this IBM product and use of those Web sites is at your own risk.

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 the names and addresses used by an actual business enterprise is entirely coincidental.

#### COPYRIGHT LICENSE:

This information contains sample application programs in source language, which illustrates 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. 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 IBM's application programming interfaces.

Each copy or any portion of these sample programs or any derivative work, must include a copyright notice as follows:

(c) (your company name) (year). Portions of this code are derived from IBM Corp. Sample Programs. (c) Copyright IBM Corp. \_enter the year or years\_. All rights reserved.

---

## Trademarks

IBM, the IBM logo, and [ibm.com](http://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 [www.ibm.com/legal/copytrade.shtml](http://www.ibm.com/legal/copytrade.shtml)

Other company, product, or service names may be the trademarks or service marks of others.

---

# Index

## Numerics

- 64-bit
  - kernel extension 18, 19

## A

- accented characters 206
- asynchronous I/O subsystem
  - subroutines 87
  - subroutines affected by 88
- ataide\_buf structure (IDE) 306
  - fields 307
- ATM LAN Emulation device driver 114
  - close 120
  - configuration parameters 116
  - data reception 120
  - data transmission 120
  - entry points 119
  - open 119
  - trace and error logging 125
- ATM LANE
  - clients
    - adding 115
- ATM MPOA client
  - tracing and error logging 127
- atmle\_ctl 121
- ATMLE\_MIB\_GET 121
- ATMLE\_MIB\_QUERY 121
- atomic operations 62
- attributes 102

## B

- block (physical volumes) 209
- block device drivers
  - I/O kernel services 41
- block I/O buffer cache
  - managing 48
  - supporting user access to device drivers 48
  - using write routines 49
- block I/O buffer cache kernel services 42
- bootlist command
  - altering list of boot devices 105

## C

- callback
  - function 47
- cfgmgr command
  - configuring devices 99, 105
- character I/O kernel services 42
- chdev command
  - changing device characteristics 105
  - configuring devices 99
- child devices 101
- CIO\_ASYNC\_STATUS 111
- CIO\_HALT\_DONE 110

- CIO\_LOST\_STATUS 110
- CIO\_NULL\_BLK 110
- CIO\_START\_DONE 110
- CIO\_TX\_DONE 110
- clients
  - ATM LANE
    - adding 115
- commands
  - errlogger 346
  - errmsg 344
  - errpt 344, 346
  - errupdate 346
  - trcrpt 348, 349
- communications device handlers
  - common entry points 108
  - common status and exception codes 109
  - common status blocks 109
  - interface kernel services 74
  - kernel-mode interface 107
  - mbuf structures 108
  - types
    - Ethernet 155
    - Fiber Distributed Data Interface (FDDI) 127
    - Forum Compliant ATM LAN Emulation 114
    - Multiprotocol (MPQP) 111
    - PCI Token-Ring device drivers 146
    - SOL (serial optical link) 112
    - Token-Ring (8fa2) 139
    - Token-Ring (8fc8) 131
    - user-mode interface 107
- communications I/O subsystem
  - physical device handler model 108
- compiling
  - when using trace 364
- complex locks 61
- component trace 315
- Component Trace 338
  - callback routine 342
  - controlling 339
  - initializing 339
  - managing trace levels 339
  - tracing events 339
  - unregistering a component 339
- configuration
  - low function terminal interface 203
- cross-memory kernel services 67

## D

- DASD subsystem
  - device block level description 313
  - device block operation
    - cylinder 314
    - head 314
    - sector 313
    - track 313
- data flushing 69
- dataless workstations, copying a system dump on 320

- DDS 103
- debug 347
- debugger 315
- device attributes
  - accessing 102
  - modifying 103
- device configuration database
  - configuring 95
  - customized database 95
  - predefined database 95, 100
- device configuration manager
  - configuration hierarchy 96
  - configuration rules 96
  - device dependencies graph 96
  - device methods 98
  - invoking 97
- device configuration subroutines 105
- device configuration subsystem 95, 96
  - adding unsupported devices 100
  - configuration commands 105
  - configuration database structure 94
  - configuration subroutines 105
  - database configuration procedures 95
  - device classifications 93
  - device dependencies 101
  - device method level 94
  - device types 97
  - high-level perspective 94
  - low-level perspective 95
  - object classes in 97
  - run-time configuration commands 99
  - scope of support 93
  - writing device methods for 98
- Device control operations 184
  - NDD\_CLEAR\_STATS 185
  - NDD\_DISABLE\_ADAPTER 187
  - NDD\_DISABLE\_ADDRESS 185
  - NDD\_DISABLE\_MULTICAST 186
  - NDD\_ENABLE\_ADAPTER 187
  - NDD\_ENABLE\_ADDRESS 184
  - NDD\_ENABLE\_MULTICAST 186
  - NDD\_GET\_ALL\_STATS 185
  - NDD\_GET\_STATS 184
  - NDD\_MIB\_ADDR 185
  - NDD\_MIB\_GET 184
  - NDD\_MIB\_QUERY 184
  - NDD\_PROMISCUOUS\_OFF 186
  - NDD\_PROMISCUOUS\_ON 186
  - NDD\_SET\_LINK\_STATUS 187
  - NDD\_SET\_MAC\_ADDR 187
- Device Control Operations
  - NDD\_CLEAR\_STATS 137
  - NDD\_DISABLE\_ADDRESS 137
  - NDD\_ENABLE\_ADDRESS 136
  - NDD\_GET\_ALL\_STATS 137
  - NDD\_GET\_STATS 136
  - NDD\_MIB\_ADDR 137
  - NDD\_MIB\_GET 136
  - NDD\_MIB\_QUERY 136
- device dependent structure
  - format 104
- device dependent structure (*continued*)
  - updating
    - using the Change method 103
- device driver
  - including in a system dump 317
- device driver management kernel services 58
- device drivers
  - adding 101
  - device dependent structure 103
  - display 205
  - entry points 204
  - interface 204
  - pseudo
    - low function terminal 204
- device methods
  - adding devices 101
  - Change method and device dependent structure 103
  - changing device states 99
  - Configure method and device dependent structure 103
  - for changing the database and not device state 100
  - interfaces 98
  - interfaces to
    - run-time commands 99
  - invoking 98
  - method types 98
  - source code examples of 98
  - writing 98
- device states 99
- devices
  - child 101
  - dependencies 101
  - SCSI 223
- diacritics 206
- diagnostics
  - low function terminal interface 205
- direct access storage device subsystem 209
- diskless systems
  - configuring dump device 316
  - dump device for 316
- display device driver 205
  - interface 205
- DMA management
  - setting up transfers 50
- DMA management kernel services 43
- dump 315
  - configuring dump devices 316
  - copying from dataless machines 320
  - copying to other media 320
  - starting 317
  - system dump facility 315
- dump device
  - determining the size of 323
  - determining the type of logical volume 323
  - increasing the size of 322, 323
- dump devices 316

## E

- eeh callback
  - function 47
- EEH error handling
  - kernel services
    - table 47
- encapsulation 75
- entry points
  - communications physical device handler 108
  - device driver 204
  - IDE adapter driver 309
  - IDE device driver 309
  - logical volume device driver 213
  - MPQP device handler 111
  - SCSI adapter device driver 241
  - SCSI device driver 241
  - SOL device handler 112
- errlogger command 346
- errmsg command 344
- error conditions
  - SCSI\_ADAPTER\_HDW\_FAILURE 292
  - SCSI\_ADAPTER\_SFW\_FAILURE 292
  - SCSI\_CMD\_TIMEOUT 292
  - SCSI\_FUSE\_OR\_TERMINAL\_PWR 292
  - SCSI\_HOST\_IO\_BUS\_ERR 291
  - SCSI\_NO\_DEVICE\_RESPONSE 292
  - SCSI\_TRANSPORT\_BUSY 292
  - SCSI\_TRANSPORT\_DEAD 292
  - SCSI\_TRANSPORT\_FAULT 291
  - SCSI\_TRANSPORT\_RESET 292
  - SCSI\_WW\_NAME\_CHANGE 292
- error logging 344
  - adding logging calls 345
  - coding steps 345
  - determining the importance 344
  - determining the text of the error message 344
  - thresholding level 344
- error messages
  - determining the text of 344
- errpt command 344, 346
- errsave kernel service 344, 345
- errupdate command 346
- Ethernet device driver 155
  - asynchronous status 183
  - configuration parameters 157
  - device control operations 184
  - entry points 179
  - NDD\_CLEAR\_STATS 185
  - NDD\_DISABLE\_ADAPTER 187
  - NDD\_DISABLE\_ADDRESS 185
  - NDD\_DISABLE\_MULTICAST 186
  - NDD\_ENABLE\_ADAPTER 187
  - NDD\_ENABLE\_ADDRESS 184
  - NDD\_ENABLE\_MULTICAST 186
  - NDD\_GET\_ALL\_STATS 185
  - NDD\_GET\_STATS 184
  - NDD\_MIB\_ADDR 185
  - NDD\_MIB\_GET 184
  - NDD\_MIB\_QUERY 184
  - NDD\_PROMISCUOUS\_OFF 186
  - NDD\_PROMISCUOUS\_ON 186

- Ethernet device driver (*continued*)
  - NDD\_SET\_LINK\_STATUS 187
  - NDD\_SET\_MAC\_ADDR 187
- events
  - management of 75
- exception codes
  - communications device handlers 109
- exception handlers
  - implementing
    - in kernel-mode 15, 16, 17, 18
    - in user-mode 18
  - registering 75
- exception handling
  - interrupts and exceptions 14
- modes
  - kernel 14
  - user 18
- processing exceptions
  - basic requirements 15
  - default mechanism 14
  - kernel-mode 14
- exception management kernel services 74
- execution environments
  - interrupt 6
  - process 5

## F

- FCP 256
  - asynchronous event handling 284
  - autosense data 285
  - consolidated commands 288
  - fragmented commands 289
  - initiator-mode recovery 286
  - NACA=1 error 285
  - openx subroutine options 296
  - recovery from failure 284
  - returned status 286
  - SC\_CHECK\_CONDITION 287
  - scsi\_buf structure 290
  - spanned commands 288
- FCP adapters
  - IOCINFO 264
- FDDI device driver 127
  - configuration parameters 127
  - entry points 128
  - trace and error logging 129
- Fiber Distributed Data Interface device driver 127
- file descriptor 63
- file systems
  - logical file system 35
  - virtual file system 36
- files
  - /dev/error 344
  - /dev/systrctl 348
  - /etc/trcfmt 349, 365
  - sys/err\_rec.h 346
  - sys/errids.h 345
  - sys/trchkid.h 349, 350, 365
  - sys/trcmacros.h 349
- filesystem 35

- fine granularity timer services 79
- Forum Compliant ATM LAN Emulation device driver 114
- function
  - callback 47

## G

- g-nodes 37
- getattr subroutine
  - modifying attributes 103
- graphic input device 197

## H

- hardware interrupt kernel services 42

## I

- I/O kernel services
  - block I/O 41
  - buffer cache 42
  - character I/O 42
  - DMA management 43
  - interrupt management 42
  - memory buffer (mbuf) 43
- IDE subsystem
  - adapter driver
    - entry points 309
    - ioctl commands 310, 311
    - performing dumps 309
  - consolidated commands 306
  - device communication
    - initiator-mode support 303
  - error processing 309
  - error recovery
    - analyzing returned status 304
    - initiator mode 304
  - fragmented commands 306
  - IDE device driver
    - design requirements 309
    - entry points 309
    - internal commands 305
    - responsibilities relative to adapter device driver 303
  - IDEIOIDENT 312
  - IDEIOINQU 311
  - IDEIOREAD 311
  - IDEIORESET 311
  - IDEIOSTART 311
  - IDEIOSTOP 311
  - IDEIOSTUNIT 311
  - IDEIOTUR 311
  - initiator I/O request execution 305
  - spanned commands 306
  - structures
    - ataide\_buf structure 306
  - typical adapter transaction sequence 304
- input device, subsystem 197
- input ring mechanism 204

- interface
  - low function terminal subsystem 203
- interrupt execution environment 6
- interrupt management
  - defining levels 49
  - setting priorities 50
- interrupt management kernel services 49
- interrupts
  - management services 42
- INTSTOLLONG macro 26
- IOCINFO
  - FCP adapters 264
  - iSCSI adapters 264
  - Virtual SCSI 264
- ioctl commands
  - SCIOCMD 248
- iSCSI 256, 290
  - autosense data 285
  - consolidated commands 288
  - fragmented commands 289
  - initiator-mode recovery 286
  - NACA=1 error 285
  - openx subroutine options 296
  - returned status 286
  - SC\_CHECK\_CONDITION 287
  - scsi\_buf structure 290
  - spanned commands 288
- iSCSI adapters
  - IOCINFO 264

## K

- kernel data
  - accessing in a system call 24
- kernel environment 1
  - base kernel services 2
  - creation of kernel processes 8
  - exception handling 14
  - execution environments
    - interrupt 6
    - process 5
  - libraries
    - libcsys 4
    - libsys 5
  - loading kernel extensions 3
  - private routines 3
  - programming
    - kernel threads 6
- kernel environment, runtime 41
- kernel extension binding
  - adding symbols to the /unix name space 2
  - using existing libraries 4
- kernel extension development
  - 64-bit 18
- kernel extension libraries
  - libcsys 4
  - libsys 5
- kernel extension programming environment
  - 64-bit 19

- kernel extensions
  - accessing user-mode data
    - using cross-memory services 12
    - using data transfer services 12
  - interrupt priority
    - service times 50
  - loading 3
  - loading and binding services 58
  - management services 59
  - serializing access to data structures 12
  - unloading 3
  - using with system calls 2
- kernel key sets 386
- kernel keys 386
- kernel processes
  - accessing data from 9
  - comparison to user processes 8
  - creating 10, 74
  - executing 10
  - handling exceptions 11
  - handling signals 11
  - obtaining cross-memory descriptors 9
  - preempting 10
  - terminating 10
  - using system calls 11
- kernel protection domain 8, 9, 23
- kernel services 41
  - address family domain 72
  - atomic operations 62
  - categories
    - EEH 44
    - I/O 41, 42, 43
    - I/O, enhanced error handling 44
    - memory 65, 66, 67
  - communications device handler interface 74
  - complex locks 61
  - device driver management 58, 59
  - errsave 344, 345
  - exception management 74
  - fine granularity 79
  - interface address 73
  - loading 3
  - lock allocation 60
  - locking 60
  - logical file system 63
  - loopback 73
  - management 58, 59
  - memory 64
  - message queue 72
  - multiprocessor-safe timer service 81
  - network 72
  - network interface device driver 72
  - process level locks 62
  - process management 74
  - protocol 73
  - Reliability Availability Serviceability (RAS) 77
  - routing 73
  - security 77
  - simple locks 60
  - time-of-day 78
  - timer 79

- kernel services (*continued*)
  - unloading kernel extensions 3
  - virtual file system 82
- kernel structures
  - encapsulation 75
- kernel symbol resolution
  - using private routines 3
- kernel threads
  - creating 7, 74
  - executing 7
  - terminating 7
- key-protected 388
- key-safe 388

**L**

- ldata kernel services
  - ldata 65
- lft 203
- LFT
  - accented characters 206
- libraries
  - libcsys 4
  - libsys 5
- live dump 326, 330
  - callback commands 328
  - detail levels 330
  - initiating 327
- locking
  - conventional locks 12
  - kernel-mode strategy 13
  - serializing access to a predefined data structure and 12
- locking kernel services 60
- lockl locks 62
- locks
  - allocation 60
  - atomic operations 62
  - complex 61
  - lockl 62
  - simple 60
- logical file system 63
  - component structure 36
  - file routines 36
  - v-nodes 36
  - file system role 35
- logical volume device driver
  - bottom half 213
  - data structures 213
  - physical device driver interface 215
  - pseudo-device driver role 212
  - top half 213
- logical volume manager
  - DASD support 209
- logical volume subsystem
  - bad block processing 215
  - logical volume device driver 212
  - physical volumes
    - comparison with logical volumes 209
    - reserved sectors 210
- LONG32TOLONG64 macro 26

- loopback kernel services 73
- low function terminal
  - configuration commands 204
  - functional description 203
  - interface 203
    - components 204
    - configuration 203
    - device driver entry points 204
    - ioctl's 204
    - terminal emulation 203
    - to display device drivers 204
    - to system keyboard 204
- low function terminal interface
  - AIXwindows support 204
- low function terminal subsystem 203
  - accented characters supported 206
- lsattr command
  - displaying attribute characteristics of devices 105
- lscfg command
  - displaying device diagnostic information 105
- lsconn command
  - displaying device connections 105
- lsdev command
  - displaying device information 105
- lsparent command
  - displaying information about parent devices 105

## M

- macros
  - INTSTOLLONG 26
  - LONG32TOLONG64 26
  - memory buffer (mbuf) 43
- management kernel services 58
- management services
  - file descriptor 63
- mbuf structures
  - communications device handlers 108
- memory buffer (mbuf) kernel services 43
- memory buffer (mbuf) macros 43
- memory kernel services
  - memory management 65
  - memory pinning 65
  - user memory access 66
- message queue kernel services 72
- mkdev command
  - adding devices to the system 105
  - configuring devices 99
- MODS 315, 367
- MPQP device handlers
  - binary synchronous communication
    - message types 111
    - receive errors 112
  - entry points 111
- multiprocessor-safe timer services 81
- Multiprotocol device handlers 111

## N

- NACA=1 error 285
- NDD\_ADAP\_CHECK 134

- NDD\_AUTO\_RMV 134
- NDD\_BUS\_ERR 134
- NDD\_CLEAR\_STATS 122, 137, 185
- NDD\_CMD\_FAIL 134
- NDD\_DEBUG\_TRACE 123
- NDD\_DISABLE\_ADAPTER 187
- NDD\_DISABLE\_ADDRESS 122, 137, 185
- NDD\_DISABLE\_MULTICAST 122, 186
- NDD\_ENABLE\_ADAPTER 187
- NDD\_ENABLE\_ADDRESS 122, 136, 184
- NDD\_ENABLE\_MULTICAST 123, 186
- NDD\_GET\_ALL\_STATS 123, 137, 185
- NDD\_GET\_STATS 124, 136, 184
- NDD\_MIB\_ADDR 124, 137, 185
- NDD\_MIB\_GET 124, 136, 184
- NDD\_MIB\_QUERY 124, 136, 184
- NDD\_PIO\_FAIL 133
- NDD\_PROMISCUOUS\_OFF 186
- NDD\_PROMISCUOUS\_ON 186
- NDD\_SET\_LINK\_STATUS 187
- NDD\_SET\_MAC\_ADDR 187
- NDD\_TX\_ERROR 134
- NDD\_TX\_TIMEOUT 134
- network kernel services
  - address family domain 72
  - communications device handler interface 74
  - interface address 73
  - loopback 73
  - network interface device driver 72
  - protocol 73
  - routing 73

## O

- object data manager 100
- ODM 100
- odmadd command
  - adding devices to predefined database 100
- openx subroutine 296
  - SC\_DIAGNOSTIC 296
  - SC\_FORCED\_OPEN 296
  - SC\_NO\_RESERVE 296
  - SC\_RETAIN\_RESERVATION 296
  - SC\_SINGLE 296
- optical link device handlers 112

## P

- parameters
  - long 26
  - long long 26
  - scalar 26
  - signed long 26
  - uintptr\_t 26
- partition (physical volumes) 210
- PCI Token-Ring Device Driver
  - trace and error logging 151
- PCI Token-Ring High Device Driver
  - entry points 147
- PCI Token-Ring High Performance
  - configuration parameters 146



- performance tracing 315
- physical volumes
  - block 209
  - comparison with logical volumes 209
  - limitations 210
  - partition 210
  - reserved sectors 210
  - sector layout 210
- pinning
  - memory 65
- predefined attributes object class
  - accessing 102
  - modifying 103
- printer addition management subsystem
  - adding a printer definition 220
  - adding a printer formatter 221
  - adding a printer type 219
  - defining embedded references in attribute strings 221
  - modifying printer attributes 220
- printer formatter
  - defining embedded references 221
- printers
  - unsupported types 219
- private routines 3
- process execution environment 5
- process management kernel services 74
- processes
  - creating 74
- protection domains
  - kernel 23
  - understanding 23
  - user 23
- protection gates 387
- pseudo components 327
- pseudo device driver
  - low function terminal 204
- putattr subroutine
  - modifying attributes 103

## R

- RCM 205
- referenced routines
  - for memory pinning 71
  - to support address space operations 71
  - to support cross-memory operations 71
  - to support pager back ends 71
- Reliability Availability Serviceability (RAS) kernel services 77
- rendering context manager 204, 205
- restbase command
  - restoring customized information to configuration database 105
- rmdev command
  - configuring devices 99
  - removing devices from the system 105
- routine
  - callback 47
- runtime kernel environment 41

## S

- SAM 255, 277
  - adapter device driver interfaces 300
  - asynchronous event handling 283
  - closing the device 299
  - command tag queuing 289
  - device driver interfaces 300
  - driver transaction sequence 287
  - dumps 300
  - error processing 299
  - error recovery 285
  - FCP 255
  - initiator I/O requests 288
  - initiator-mode recovery 285
  - interfaces 300
  - iSCSI 255
  - SAS 255
- SAM Adapter device driver
  - ioctl commands, required 301
- SAM device driver
  - responsibilities 296
  - SC\_DIAGNOSTIC 297
  - SC\_FORCED\_OPEN 297
  - SC\_FORCED\_OPEN\_LUN 297
  - SC\_NO\_RESERVE 298
  - SC\_RETAIN\_RESERVATION 297
  - SC\_SINGLE 298
- sample code
  - trace format file 354
- savebase command
  - saving customized information to configuration database 105
- sc\_buf structure (SCSI) 232
- scalar parameters 26
- SCIOCMD 248
- SCIOLCHBA 277
- SCIOLCMD 274
- SCIOLEVENT 267
- SCIOLHALT 273
- SCIOLINQU 268
- SCIOLNMSRV 275
- SCIOLPASSTHRUHBA 277
- SCIOLPAYLD 276
- SCIOLQWWN 275
- SCIOLREAD 271
- SCIOLRESET 272
- SCIOLSTART 265
- SCIOLSTOP 267
- SCIOLSTUNIT 269
- SCIOLTUR 270
- SCSI
  - virtual 256, 290
- SCSI subsystem
  - adapter device driver
    - entry points 241
    - initiator-mode ioctl commands 247
    - ioctl operations 244, 247, 248, 249, 250, 251, 252
    - performing dumps 241
    - responsibilities relative to SCSI device driver 223
    - target-mode ioctl commands 250

- SCSI subsystem (*continued*)
  - asynchronous event handling 224
  - command tag queuing 232
  - device communication
    - initiator-mode support 224
    - target-mode support 224
  - error processing 240
  - error recovery
    - initiator mode 226
    - target mode 229
  - initiator I/O request execution
    - fragmented commands 231
    - gathered write commands 231
    - spanned or consolidated commands 230
  - initiator-mode adapter transaction sequence 229
  - SCSI device driver
    - asynchronous event-handling routine 226
    - closing a device 240
    - design requirements 237
    - entry points 241
    - internal commands 229
    - responsibilities relative to adapter device driver 223
    - using openx subroutine options 237
  - structures
    - sc\_buf structure 232
    - tm\_buf structure 240, 244
  - target-mode interface 242, 243, 245
    - interaction with initiator-mode interface 242
- SCSI\_ADAPTER\_HDW\_FAILURE 292
- SCSI\_ADAPTER\_SFW\_FAILURE 292
- scsi\_buf structure 290
  - fields 290
- SCSI\_CMD\_TIMEOUT 292
- SCSI\_FUSE\_OR\_TERMINAL\_PWR 292
- SCSI\_HOST\_IO\_BUS\_ERR 291
- SCSI\_NO\_DEVICE\_RESPONSE 292
- SCSI\_TRANSPORT\_BUSY 292
- SCSI\_TRANSPORT\_DEAD 292
- SCSI\_TRANSPORT\_FAULT 291
- SCSI\_TRANSPORT\_RESET 292
- SCSI\_WW\_NAME\_CHANGE 292
- security kernel services 77
- serial optical link device handlers 112
- signal management 75
- Small Computer Systems Interface subsystem 223
- SOL device handlers
  - changing device attributes 114
  - configuring physical and logical devices 113
  - entry points 112, 113
  - special files interfaces 113
- status and exception codes 109
- status blocks
  - communications device handler
    - CIO\_ASYNC\_STATUS 111
    - CIO\_HALT\_DONE 110
    - CIO\_LOST\_STATUS 110
    - CIO\_NULL\_BLK 110
    - CIO\_START\_DONE 110
    - CIO\_TX\_DONE 110
  - communications device handlers and 109
  - status codes
    - communications device handlers and 109
  - status codes, system dump 319
  - storage 209
  - storage protection keys
    - kernel storage-protection keys 385
  - stream-based tty subsystem 203
  - structures
    - scsi\_buf 290
  - subroutines
    - close 197
    - ioctl 197
    - open 197
    - read 197
    - write 197
  - subsystem
    - graphic input device 197
    - low function terminal 203
    - streams-based tty 203
  - system calls
    - accessing kernel data in 24
    - asynchronous signals 30
    - error information 32
    - exception handling 31
    - execution 24
    - in kernel protection domain 23
    - in user protection domain 23
    - nesting for kernel-mode use 31
    - page faulting 31
    - passing parameters 25
    - preempting 29
    - services for all kernel extensions 32
    - services for kernel processes only 32
    - setjmpx kernel service 30
    - signal handling in 30
    - stacking saved contexts 30
    - using with kernel extensions 2
    - wait termination 30
  - system dump 330
    - callback commands 323
    - checking status 319
    - configuring dump devices 316
    - copy from server 321
    - copying from dataless machines 320
    - copying on a non-dataless machine 321
    - copying to other media 320
    - firmware-assisted system dump 315
    - including device driver data 317
    - locating 321
    - reboot in normal mode 321
    - starting 317
    - traditional system dump 315
  - system dump facility 315

## T

- terminal emulation
  - low function terminal 203
- threads
  - creating 74
- time-of-day kernel services 78

- timer kernel services
  - coding the timer function 81
  - compatibility 79
  - determining the timer service to use 79
  - fine granularity 79
  - reading time into time structure 80
  - watchdog 79
- timer service
  - multiprocessor-safe 81
- tm\_buf structure (SCSI) 240
- TOK\_ADAP\_INIT 134
- TOK\_ADAP\_OPEN 135
- TOK\_DMA\_FAIL 135
- TOK\_RECOVERY\_THRESH 133
- TOK\_RING\_SPEED 135
- TOK\_RMV\_ADAP 135
- TOK\_WIRE\_FAULT 135
- Token-Ring (8fa2) device driver 139
  - asynchronous state 141
  - configuration parameters 139
  - data reception 140
  - data transmission 140
  - device driver close 140
  - device driver open 140
  - trace and error logging 144
- Token-Ring (8fc8) device 131
- Token-Ring (8fc8) device driver
  - configuration parameters 131
  - trace and error logging 138
- trace
  - controlling 348
- trace events
  - defining 349
  - event IDs 350
    - determining location of 350
  - format file example 354
  - format file stanzas 351
  - forms of 349
  - macros 349
- trace facility 347
  - configuring 348
  - controlling 348
  - controlling using commands 348
  - defining events 349
  - event IDs 350
  - events, forms of 349
  - hookids 350
  - reports 349
  - starting 347, 348
  - using 348
- trace report
  - filtering 366
  - producing 349
  - reading 366
- tracing 347
  - configuring 348
  - starting 347, 348
- trcrpt command 348, 349

## U

- user commands
  - configuration 204
- user protection domain 23

## V

- v-nodes 36
- virtual file system 35
  - configuring 39
  - data structures 38
  - file system role 36
  - generic nodes (g-nodes) 37
  - header files 38
  - interface requirements 37
  - mount points 36
  - virtual nodes (v-nodes) 36
- virtual file system kernel services 82
- virtual memory management
  - addressing data 68
  - data flushing 69
  - discarding data 69
  - executable data 69
  - installing pager backends 70
  - moving data 69
  - objects 68
    - protecting data 69
    - referenced routines
      - for manipulating objects 70
- virtual memory management kernel services 66
- virtual memory manager 68
- Virtual SCSI 256, 290
  - IOCINFO 264
- vm\_uiomove 67, 69, 71



---

# Readers' Comments — We'd Like to Hear from You

**AIX Version 6.1**  
**Kernel Extensions and Device Support Programming Concepts**

**Publication No. SC23-5261-02**

We appreciate your comments about this publication. Please comment on specific errors or omissions, accuracy, organization, subject matter, or completeness of this book. The comments you send should pertain to only the information in this manual or product and the way in which the information is presented.

For technical questions and information about products and prices, please contact your IBM branch office, your IBM business partner, or your authorized remarketer.

When you send comments to IBM, you grant IBM a nonexclusive right to use or distribute your comments in any way it believes appropriate without incurring any obligation to you. IBM or any other organizations will only use the personal information that you supply to contact you about the issues that you state on this form.

Comments:

Thank you for your support.

Submit your comments using one of these channels:

- Send your comments to the address on the reverse side of this form.
- Send your comments via e-mail to: [pserinfo@us.ibm.com](mailto:pserinfo@us.ibm.com)

If you would like a response from IBM, please fill in the following information:

\_\_\_\_\_

Name

\_\_\_\_\_

Address

\_\_\_\_\_

Company or Organization

\_\_\_\_\_

Phone No.

\_\_\_\_\_

E-mail address



Fold and Tape

Please do not staple

Fold and Tape



NO POSTAGE  
NECESSARY  
IF MAILED IN THE  
UNITED STATES

# BUSINESS REPLY MAIL

FIRST-CLASS MAIL PERMIT NO. 40 ARMONK, NEW YORK

POSTAGE WILL BE PAID BY ADDRESSEE

IBM Corporation  
Information Development  
Department 04XA-905-6B013  
11501 Burnet Road  
Austin, TX 78758-3400



Fold and Tape

Please do not staple

Fold and Tape





Printed in U.S.A.

SC23-5261-02

