

AIX Version 4.1

Writing a Device Driver



Print Graphic between registration marks



AIX Version 4.1

Writing a Device Driver

Spine Copy

Note to Printer:

Center crop marks are given.
Please align with center of book
spine than delete crop marks.

Anything over 220 pages will
have spine copy

AIX Version 4.1

Writing a Device Driver

SC23-2593-00

Printed in the U.S.A.

Place Form Number barcode
lower left corner here.

Place Part Number bar-code
lower left corner here.

First Edition (August 1994)

This edition of *AIX Version 4.1 Writing a Device Drive* applies to AIX Version 4.1 and to all subsequent releases of AIX until otherwise indicated in new releases or technical newsletters.

The following paragraph does not apply to the United Kingdom or any country where such provisions are inconsistent with local law: THIS MANUAL IS PROVIDED "AS IS" WITHOUT WARRANTY OF ANY KIND, EITHER EXPRESSED OR IMPLIED, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND 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.

It is not warranted that the contents of this publication or the accompanying source code examples, whether individually or as one or more groups, will meet your requirements or that the publication or the accompanying source code examples are error-free.

This publication 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.

It is possible that this publication may contain references to, or information about, products (machines and programs), programming, or services that are not announced in your country. Such references or information must not be construed to mean that such products, programming, or services will be offered in your country. Any reference to a licensed program in this publication is not intended to state or imply that you can use only that licensed program. You can use any functionally equivalent program instead.

The information provided regarding publications by other vendors does not constitute an expressed or implied recommendation or endorsement of any particular product, service, company or technology, but is intended simply as an information guide that will give a better understanding of the options available to you. The fact that a publication or company does not appear in this book does not imply that it is inferior to those listed. The providers of this book take no responsibility whatsoever with regard to the selection, performance, or use of the publications listed herein.

NO WARRANTIES OF ANY KIND ARE MADE WITH RESPECT TO THE CONTENTS, COMPLETENESS, OR ACCURACY OF THE PUBLICATIONS LISTED HEREIN. ALL WARRANTIES, EXPRESSED OR IMPLIED, INCLUDING BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE ARE SPECIFICALLY DISCLAIMED. This disclaimer does not apply to the United Kingdom or elsewhere if inconsistent with local law.

A reader's comment form is provided at the back of this publication. If the form has been removed, address comments to Publications Department, Internal Zip 9630, 11400 Burnet Road, Austin, Texas 78758-3493. To send comments electronically, use this commercial internet address: aix6kpub@austin.ibm.com. Any information that you supply may be used without incurring any obligation to you.

© Copyright International Business Machines Corporation 1994. All rights reserved.

Notice to U.S. Government Users — Documentation Related to Restricted Rights — Use, duplication or disclosure is subject to restrictions set forth in GSA ADP Schedule Contract.

Trademarks and Acknowledgements

The following trademarks and acknowledgements apply to this book:

AIX is a trademark of International Business Machines Corporation.

AIXwindows is a trademark of International Business Machines Corporation.

Display PostScript is a trademark of Adobe Corporation.

GL is a trademark of Silicon Graphics, Inc.

graPHIGS is a trademark of International Business Machines Corporation.

IBM is a registered trademark of International Business Machines Corporation.

Intel is a trademark of Intel Corporation.

InfoExplorer is a trademark of International Business Machines Corporation.

Micro Channel is a trademark of International Business Machines Corporation.

OpenGL is a trademark of Silicon Graphics, Inc.

POWER Architecture is a trademark of International Business Machines Corporation.

PowerPC is a trademark of International Business Machines Corporation.

PowerPC 601 is a trademark of International Business Machines Corporation.

POWERserver is a trademark of International Business Machines Corporation.

POWERstation is a trademark of International Business Machines Corporation.

RISC System/6000 is a trademark of International Business Machines Corporation.

UNIX is a registered trademark licensed exclusively by X/Open Company.

X Window System is a trademark of Massachusetts Institute of Technology.

X/Open is a trademark of X/Open Company Limited.

X11 is a trademark of Massachusetts Institute of Technology.

Contents

About This Book	xiii
Chapter 1. Device Driver Overview	1-1
Aspects of the Kernel that Affect Device Drivers	1-2
How Device Drivers Are Accessed	1-3
Types of Device Drivers	1-5
Block Device Drivers	1-6
STREAMS Device Drivers	1-7
Character Device Drivers	1-7
Device Driver Configuration	1-8
Object Data Manager (ODM) Database	1-9
Device Driver Entry Points	1-10
xyzconfig Entry Point	1-10
xyzopen and xyzclose Entry Points	1-11
xyzread Entry Point	1-11
xyzwrite Entry Point	1-12
xyzstrategy Entry Point	1-12
xyzioctl Entry Point	1-12
STREAMS Entry Points	1-13
xyzwput Entry Point	1-13
xyzwsrv, xyzrsrv Entry Points	1-13
Sample Device Driver	1-14
Files for Sample XYZ Device Driver	1-14
makefile for Sample XYZ Device Driver	1-15
Configuration Program for Sample XYZ Device Driver	1-15
Source Code for Sample XYZ Device Driver	1-17
User Program to Invoke Sample XYZ Device Driver	1-19
Running the Sample XYZ Device Driver	1-20
Trace Output for Sample XYZ Device Driver	1-20
Routines on the Interrupt Side	1-21
Pinning Device Driver Object Files	1-22
Driving a SCSI Attached Device	1-22
Other Topics	1-23
Chapter 2. Device I/O	2-1
Address Translation	2-1
Block Address Translation	2-2
Segment Address Translation	2-2
I/O Controller Types	2-4
I/O Space on PCI and ISA Systems	2-5
Programmed I/O to PCI and ISA Adapters	2-6
Direct Memory Access	2-7
DMA on POWER and POWER2 Architectures	2-7
DMA on RSC (Single-Chip) Architectures	2-7
DMA on PowerPC Architectures	2-7

DMA Routines for PCI and ISA Adapters	2-8
Page Protection	2-9
Peer-To-Peer DMA Support	2-9
DMA Master I/O for an ISA Adapter	2-10
DMA Slave Transfers on an ISA Adapter	2-12
DMA Master Transfers on a PCI Adapter	2-13
I/O Controller Interface Translation on Micro Channel Systems	2-13
I/O Address Spaces on Micro Channel Systems	2-15
Programmed I/O to Micro Channel Adapters	2-19
DMA Master I/O on a Micro Channel Adapter	2-20
DMA Slave I/O on a Micro Channel Adapter	2-20
Alignment Issues for DMA on Micro Channel	2-21
Chapter 3. Interrupts	3-1
Overview	3-1
Interrupt Hardware Support	3-2
Interrupt Levels	3-2
Interrupt Priorities	3-4
Interrupt Level Mapping	3-6
Interrupt Handling	3-9
Interrupt Management Kernel Services	3-9
Multiprocessor Interrupt Concerns	3-10
Chapter 4. Memory Management	4-1
Memory Allocation Services	4-1
Memory Pinning Services	4-3
Memory Access Services	4-5
Virtual Memory Management Services	4-6
Cross-Memory Services	4-12
Chapter 5. Synchronization and Serialization	5-1
Timer Services	5-2
Watchdog Timers	5-2
Real-Time Timers	5-4
Event Notification	5-6
Serialization Services	5-8
Uniprocessor (UP) Serialization	5-8
Multiprocessing (MP) Serialization	5-8
Lock Overview	5-9
Device Driver Lock Models	5-9
MP-Safe Coding Sample	5-10
MP-Efficient Coding Sample	5-13
Chapter 6. Device Configuration Methods	6-1
Device States	6-1
ODM Configuration Databases	6-2
Define Methods	6-3
Configure Methods	6-4
Change Methods	6-5

Unconfigure Methods	6-5
Undefine Methods	6-6
Configuring Devices with No Parent	6-6
Adapter Device Attributes and busresolve	6-7
Configuration of Devices on PCI and ISA Bus Systems	6-8
Chapter 7. Block Device Drivers	7-1
Introduction	7-1
Block I/O Device Driver Entry Points	7-1
ddconfig Entry Point	7-2
ddopen and ddclose Entry Points	7-3
ddstrategy Entry Point	7-3
dddump Entry Point	7-6
Character Access to Block Device Drivers	7-6
Block I/O Device Summary	7-7
Chapter 8. SCSI Device Drivers	8-1
Device Driver Overview	8-1
Adapter Device Driver Overview	8-1
SCSI Adapter/Device Interface	8-2
SCSI Adapter Device Driver Routines	8-4
SCSI Adapter ioctl Operations	8-6
SCSI Device Driver Routines	8-13
Top-Half Routines	8-14
Bottom-Half Routines	8-16
PVIDs	8-17
SCSI Device Attributes	8-18
SCSI Configuration Methods	8-18
Chapter 9. Writing a Virtual File System	9-1
Multiple File System Types within the Kernel	9-2
Data Structures within a Virtual File System	9-3
gfs Structure	9-5
vfs structure	9-5
vnode structure	9-6
gnode structure	9-7
File-Over-File Mounts	9-7
Components of a Third-Party Virtual File System	9-8
Creating the Virtual File System Kernel Extension	9-9
Entry Points within the File System Kernel Extension	9-9
VFS Operations within the File System Kernel Extension	9-10
Vnode Operations within the File System Kernel Extension	9-12
Virtual Memory Operations	9-15
File System Helper	9-16
Mount Helper	9-17
Virtual File System Configuration Program	9-18
Software Installation Package	9-19
Glossary	9-21

Chapter 10. STREAMS-Based TTY Subsystem Interface	10-1
Overview	10-1
Stream Head	10-3
TIOC Module	10-4
Open Routine	10-4
Copy in Data for an IOCTL	10-5
Copy out Data for an IOCTL	10-6
LDTERM Module	10-6
Open Routine	10-6
Close Routine	10-6
Read-Side Put Routine	10-7
Write-Side Put Routine: Immediate Processing	10-9
Write-Side Service Routine: Delayed Processing	10-10
Multibyte Processing	10-10
Messages Summary	10-10
SPTR Module	10-11
Open Routine	10-11
Read-Side Put Routine	10-11
Write-Side Put Routine	10-12
Messages Summary	10-12
SLIP Module	10-13
SLIP Applications	10-13
SLIP Routines	10-13
TTY Drivers	10-13
Drivers Configuration Routine	10-14
Open Disciplines	10-15
Pacing Disciplines	10-15
Open and Close Routines	10-15
Write-Side Put Routine	10-16
Read-Side Processing	10-17
Interface with the TIOC Module	10-18
Interface with the LDTERM Module	10-18
Interface with the SPTR Module	10-19
The TTY Subsystem in a Multiprocessor Environment	10-19
TTY Modules Other Than Driver	10-19
Drivers	10-20
Special Cases	10-20
IOCTL Support and Origin	10-21
TTY Data Structures	10-24

Chapter 11. Implementing Graphical Input and 2D Graphics Device Drivers .	11-1
Porting to the AIXwindows X Server: Overview	11-1
Porting 2D Graphics Adapters	11-2
Graphics Adapter Interface (GAI) Display Subsystem	11-3
Display Subsystem Definitions	11-4
Application Programming Interface (API)	11-9
X Server	11-10
GAI Load Modules	11-11
Kernel Components of the Display Subsystem	11-13
Display Device Driver	11-14
LFT Overview	11-14
Configuration and ODM Object Classes	11-14
Display Device Driver Subroutines	11-16
Configure the Device (vddconfig)	11-16
Open a Device (vddopen)	11-17
Close a Device (vddclose)	11-18
Device Control (vddioctl)	11-19
LFT Interface Routines	11-20
Activate (vttact)	11-20
Copy Full Lines (vttcfl)	11-21
Clear Rectangle (vttclr)	11-22
Copy Line Segment (vttcpl)	11-23
Deactivate (vttidact)	11-24
Define Cursor (vttdefc)	11-24
Initialize (vttinit)	11-25
Move Cursor (vttmovc)	11-27
Scroll (vttscr)	11-27
Terminate (vttterm)	11-28
Draw Text (vtttext)	11-29
Display Driver Structure Descriptions	11-31
vtt_rc_parms	11-31
vtt_box_rc_parms	11-31
vtt_cp_parms	11-31
font_data	11-32
phys_displays	11-33
Device Dependent Structure (DDS)	11-33
Graphics Adapter Interface (GAI) 2D Adapter Load Modules	11-35
Loadable DDX Interface	11-35
Selection of Adapters	11-35
X Server Initialization Subroutines	11-37
ddxProcessArgument	11-38
FindAllAvailableDisplays	11-39
InitOutput Subroutine	11-40
Device-Dependent Initialization Subroutines	11-41
xxxentryFunc Subroutine	11-41
xxxScrInit	11-42
xxxCloseScreen	11-42
Server Termination	11-43
Adapter Access and the aixgsc System Call	11-43
Implementation Details	11-44

Minimum Resource Management Subsystem (RMS) for 2D Adapters	11-45
Implementation	11-45
Include Files	11-46
Configuring the 2D Adapter into the ODM Database	11-46
Porting Input Devices	11-48
Input Device Driver Overview	11-48
Device Driver	11-48
X Server Input Ring	11-50
SIGMSG Signal	11-51
Block and Wakeup Handling	11-51
xxxBlockHandler Subroutine	11-52
xxxWakeupHandler Subroutine	11-52
Event Processing	11-53
AddInputCheck Subroutine	11-53
RemoveInputCheck Subroutine	11-54
Input Load Module	11-54
InputDevPrivate structure	11-54
ExtInItInput Subroutine	11-55
deviceProc Subroutine	11-56
setDeviceMode Subroutine	11-57
setDeviceValuators Subroutine	11-58
getDeviceControl Subroutine	11-58
changeDeviceControl Subroutine	11-59
processRawInputEvents Subroutine	11-59
ODM Database Entry for Input Devices	11-60
ODM Input Device Record Example	11-61
Sample Input Device Load Module	11-61
Building a Dynamically Loadable Module	11-62
Debugging Load Modules	11-63
List of X Server Porting Subroutines	11-64
X Server Initialization	11-64
Device-Dependent Initialization	11-64
Block and Wakeup Handling (Input Devices)	11-64
Event Processing (Input Devices)	11-64
Input Load Module (Input Devices)	11-64
Chapter 12. Implementing a Network Device Driver	12-1
Writing a Network Device Driver	12-2
Overview of Network Device Driver Changes in AIX Version 4.1	12-2
Network Device Driver Initialization and Termination	12-2
CDLI – Device Driver Interface	12-5
Device Driver – CDLI Interface	12-11
Writing a Network Demuxer	12-13
Demuxer Initialization	12-13
nd_add_filter Function	12-14
nd_del_filter Function	12-15
nd_add_status Function	12-15
nd_del_status Function	12-16
nd_receive Function	12-16
nd_status Function	12-17
nd_response Function	12-17

DLPI/Socket – Network Demuxer Interface	12-18
Device Driver – Network Demuxer Interface	12-20
Sample Code – DLPI Call to ns_add_filter	12-21
Writing a Network Interface Driver	12-22
Basic Functions of a Network Interface Driver	12-22
Summary of NID Changes in AIX Version 4.1	12-22
Network Interface Driver Functions	12-22
Loading and Initialization	12-22
Communicating with the IP	12-25
Outgoing Packets	12-26
Communicating with the Device Handler	12-29
Output Data	12-29
Translating Network Addresses to Hardware Addresses	12-29
Handling NID Specific ioctl Calls	12-31
Terminating	12-34
NID and ARP Data Structures	12-34
xx_softc	12-35
arpcom	12-35
ifnet	12-35
ifaddr	12-37
ifreq	12-37
arptab	12-37
arpreq	12-38
Include Files	12-38
Tracing and Debugging for NIDs	12-39
Configuration Method for NID	12-39
Chapter 13. Network Interfaces and Protocols	13-1
Detailed Network Interfaces	13-1
STREAMS User Interfaces	13-1
Protocol Interfaces via DLPI	13-2
Writing or Porting STREAMS Network Protocols	13-3
DLPI Interfaces Supported by AIX	13-3
AIX Interpretations of Source and Destination Addresses	13-4
Protocol Address Resolution	13-4
AIX STREAMS Loading Convention	13-4
MP Serialization and Locking Options for STREAMS Modules and Drivers	13-4
TLI and XTI Interface Protocols	13-5
Obtaining Copies of the DLPI and TPI Specifications	13-5
Writing or Porting Socket Network Protocols	13-7
Initialization	13-7
Loading	13-8
Socket-Protocol Interface	13-8
Protocol-Socket Interface	13-11
Protocol-Network Interface	13-12
Network – Protocol Interface	13-14
IP Encapsulation/Adding Protocols to the System IP Protocol Switch	13-14
Sample Socket Protocol	13-15
Sample Code for Direct Access to Device Driver via STREAMS	13-18

Chapter 14. Debugging Tools	14-1
System Dump	14-1
Initiating a System Dump	14-1
Including Device Driver Information in a System Dump	14-2
Formatting a System Dump	14-4
The crash Command	14-5
crash Subcommands	14-5
Low-Level Kernel Debugger	14-26
Entering the Debugger	14-26
Debugger Subcommands and Concepts	14-27
Maps and Listings	14-40
Using the Debugger	14-45
Error Logging	14-52
Precoding Steps to Consider	14-53
Coding Steps	14-54
Writing to the /dev/error Special File	14-60
Performance Tracing	14-61
Introduction	14-61
Using the trace Facility	14-63
Controlling trace	14-65
Producing a trace Report	14-68
Defining trace Events	14-70
Usage Hints	14-85
Appendix A. New Interfaces	A-1
d_map_clear Kernel Service	A-2
d_map_disable Kernel Service	A-3
d_map_enable Kernel Service	A-4
d_map_init Kernel Service	A-5
d_map_list Kernel Service	A-6
d_map_page Kernel Service	A-8
d_map_slave Kernel Service	A-10
d_unmap_list Kernel Service	A-12
d_unmap_page Kernel Service	A-13
d_unmap_slave Kernel Service	A-14
iomem_att Kernel Service	A-15
iomem_det Kernel Service	A-17
ns_alloc Network Service	A-18
ns_free Network Service	A-19
rmalloc Kernel Service	A-20
rmfree Kernel Service	A-21
Index	X-1

About This Book

AIX Version 4.1 Writing a Device Driver contains an overview of block and character device drivers and describes how to write a device driver for AIX Version 4.1. Also included is information on debugging and packaging device drivers.

Who Should Use This Book

This book is intended for programmers and software support personnel who need detailed information on writing device drivers. Readers of this book are expected to be familiar with the C programming language, AIX commands, subroutines, and special files.

How to Use This Book

Overview of Contents

Chapters 1 through 6 are intended for all readers of this book and discuss the following topics:

- Chapter 1 is a device driver overview.
- Chapter 2 discusses device input/output.
- Chapter 3 discusses interrupts.
- Chapter 4 is about memory management.
- Chapter 5 discusses synchronization and serialization.
- Chapter 6 is about device configuration methods.

Chapters 7 through 13 each discuss a particular type of device driver programming.

- Chapter 7 is about block device drivers.
- Chapter 8 is about SCSI device drivers.
- Chapter 9 is about Virtual File Systems.
- Chapter 10 describes the STREAMS-based tty interface.
- Chapter 11 discusses implementing graphical input and 2D graphics device drivers.
- Chapter 12 discusses implementing a network device driver.
- Chapter 13 is about network interfaces and protocols.

Chapter 14 contains information on debugging device drivers.

The appendix contains reference information about some new kernel services and network services.

Highlighting

The following highlighting conventions are used in this book:

Bold	Identifies commands, keywords, files, directories, and other items whose names are predefined by the system.
<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.

Related Publications

The following books contain information about or related to writing device drivers:

- *AIX Version 4.1 Commands Reference*, Order Number SBOF-1851.
- *AIX Version 4.1 General Programming Concepts, Volume 1: Writing Programs*, Order Number SC23-2533.
- *AIX Version 4.1 General Programming Concepts, Volume 2: Debugging Programs*, Order Number SC23-2490.
- *AIX Version 4.1 Communications Programming Concepts*, Order Number SC23-2610.
- *AIX Version 4.1 Kernel Extensions and Device Support Programming Concepts*, Order Number SC23-2611.
- *AIX Version 4.1 Files Reference*, Order Number SC23-2512.
- *AIX Version 4.1 Problem Solving Guide and Reference*, Order Number SC23-2606.
- *AIX Version 4.1 Technical Reference, Volume 5: Kernel and Subsystems*, Order Number SC23-2618.
- *AIX Version 4.1 Technical Reference, Volume 6: Kernel and Subsystems*, Order Number SC23-2619.
- *PowerPC Architecture*, Order Number SR28-5124.
- *POWERstation and POWERserver Hardware Technical Information-General Architectures*, Order Number SA23-2643.
- *UNIX System V Release 4, Programmer's Guide: STREAMS*, Englewood Cliffs, N.J.: Prentice-Hall, 1990.
- Angebrannt, Susan, Drewry, Raymond, Karlton, Philip, Newman, Todd, Packard, Keith and Scheifler, Robert W. *Strategies for Porting the X Window Server*, Massachusetts Institute of Technology. 1991.
- Fortune, Erik and Israel, Elias. *The X-Window Server*, Digital Press.
- Gettys, James, Newman, Ron and Scheifler, Robert W. *Xlib—C Language X Interface, MIT X Consortium Standard, X Version 11, Release 5*, MIT X Consortium 1991.
- Leffler, Samuel J., and others. *The Design and Implementation of the 4.3 BSD UNIX Operating System*, Addison-Wesley. 1990.
- Patrick, Mark, and Sachs, George. *X11 Input Extension Library Specification. MIT X Consortium Standard. X Version 11, Release 5*, Hewlett-Packard Company, Ardent Computer, and the Massachusetts Institute of Technology. 1989, 1990, 1991.
- Patrick, Mark, and Sachs, George. *X11 Input Extension Protocol Specification. MIT X Consortium Standard. X Version 11, Release 5*, Hewlett-Packard Company, Ardent Computer, and the Massachusetts Institute of Technology. 1989, 1990, 1991.
- Sachs, George. *X11 Input Extension Porting Document. MIT X Consortium Standard. X Version 11, Release 5*, Hewlett-Packard Company and the Massachusetts Institute of Technology. 1989, 1990, 1991.
- Scheifler, Robert W. *X Window System Protocol, MIT X Consortium Standard, X Version 11, Release 5*, MIT X Consortium 1991.
- Womack, et. al. *PEX Protocol Specification, X Version 5.1, MIT X Consortium Standard*, Massachusetts Institute of Technology 1988, 1989, 1990, 1991, 1992.

- Womack, et. al. *PEX Protocol Encoding & Version 5.1, MIT X Consortium Standard*
Massachusetts Institute of Technology 1988, 1989, 1990, 1991, 1992.

Ordering Additional Copies of This Book

You can order publications from your sales representative or from your point of sale.

If you received a printed copy of *Documentation Overview* with your system, use that book for information on related publications and for instructions on ordering them.

To order additional copies of this book, use Order Number SC23-2593.

Chapter 1. Device Driver Overview

Many computer programs are dedicated to working with attached devices in some way. For example, there are programs to send control characters to a printer, programs to receive characters from a terminal, and programs to read data from a tape. In a broad sense, each of these programs is a *device driver* because the program is dedicated to handling input from or output to a device. Such programs are usually regarded as being part of, or an extension of, the computer's operating system.

Any operating system that supports multitasking (such as AIX) needs some way to prevent one program from writing to, or changing the state of, some device that is already being accessed by another program. So, a multitasking operating system relies on the computer's processors to distinguish between privileged and non-privileged execution of instructions. Therefore, one must distinguish between programs that execute in privileged mode (kernel mode) and those that execute in user mode. The AIX kernel consists of all software that executes in kernel mode.

Even though AIX programs that execute in user mode can drive devices, such as a printer or some device attached to a serial port, they can only do so by invoking software that is part of the kernel. Because kernel device drivers are considerably more complex than drivers that execute in user mode, from here on, the term *device driver* will only refer to software that handles a device while executing in kernel mode.

Device drivers are more complex than user software for several reasons:

- Device drivers output data to a device or demand data from a device.

This means that the driver may have to read or write to registers on a card attached to an I/O bus, or the driver may have to set up the means for the data to be transferred in some other way. A device driver is intimately interconnected with processor memory design, how the processor performs I/O, and with the architecture of the I/O bus attached to the system. So, device drivers are not portable; migrating the driver routines from one system to another often requires the routines to be rewritten.

- Device drivers may have to process interrupts generated by a card attached to the system I/O bus.

When a terminal sends a character to the computer, or a printer runs out of paper, or a tape drive has completed writing a block of data, the card serving as an adapter between the device and the I/O bus on the computer generates an interrupt. The software routines, within a device driver, that process interrupts (called *interrupt handlers*) take some sort of action like buffering incoming data, or signaling a process. Because such interrupts occur *asynchronously*, meaning that they occur without regard to what instructions the computer's processors are executing, the interrupt may occur while a processor on the computer is in the middle of handling another interrupt. Therefore, interrupt handlers must be reentrant; in other words, they must be able to access shared resources (such as non-private data) and exclude concurrent access by any software including another instance of itself.

Device driver routines that (asynchronously) execute in the context of handling an interrupt are said to be, *on the interrupt side* and are occasionally referred to as the *device handler* but this is not the same thing as a *network device handler*.

Device driver routines that (synchronously) execute in the context of a calling process are said to be *on the call side*. For more information on concurrent access of shared data, see Chapter 5, "Synchronization and Serialization."

- Device drivers may have to execute *in real time*
Device drivers may have to respond to an interrupt, or perform some other function within a certain fixed period of time.
- A device driver is a collection of routines. There is no *main* routine.
The routines are usually written in C and compiled to produce one or two Extended Object File Format (XCOFF) object files. The object files are linked to enable the kernel loader to resolve kernel symbols. As a result of linking, the loader section is filled out with a list of symbols to import from the kernel. The symbols are in the file `/lib/kernex.exp`. The linking also establishes the driver's configuration routine as the default entry point for beginning execution. A simple example is shown in "Sample Device Driver" on page 1-14.

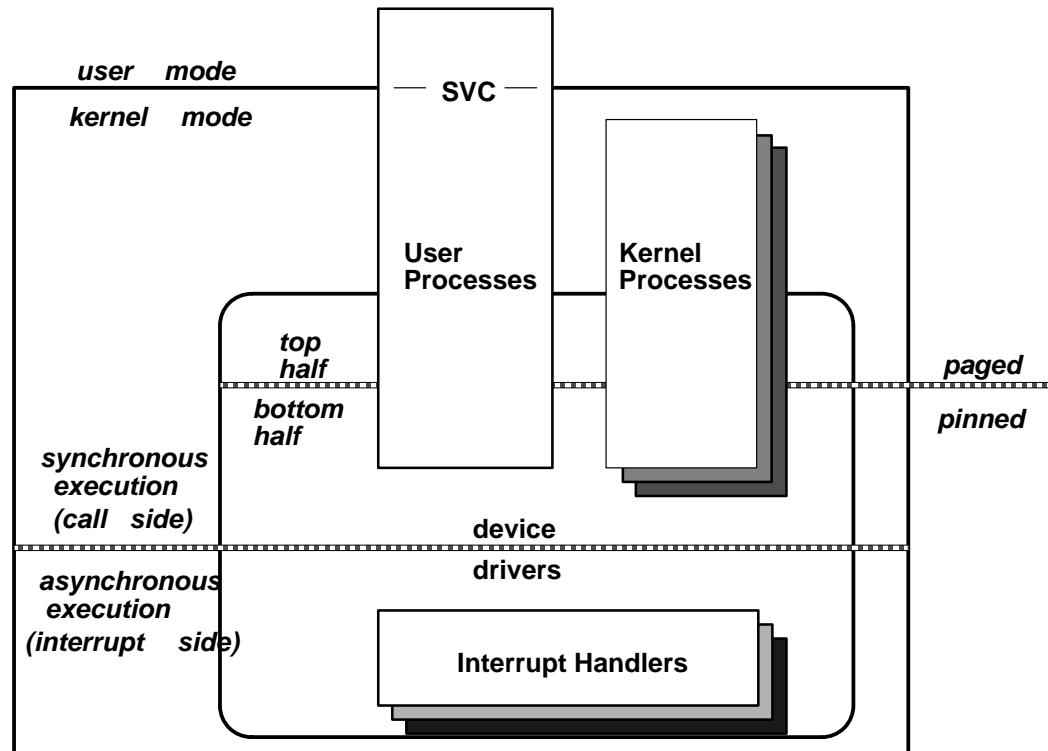
Aspects of the Kernel that Affect Device Drivers

There are a number of attributes of the AIX kernel that affect device drivers:

- Kernel routines can only call kernel services.
Kernel routines cannot call routines meant to execute in user mode. So, device driver routines are not linked with the C library `libc.a`, nor can they invoke system calls. There are some kernel DMA and timer routines in `libsys.a`, and there are some C library calls written to execute in kernel mode in the library `libcsys.a`. For more information on them, please refer to "Understanding Kernel Binding," in *AIX Version 4.1 Kernel Extensions and Device Support Programming Concepts*
- Kernel code and data that is not pinned (explicitly or implicitly) is paged into system RAM from a paging logical volume on disk.
So, one must distinguish between device driver routines that are to be collectively pinned, called *the bottom half* and those that are to be paged, called *the top half*. Because device driver routines on the interrupt side must be pinned, the phrases *in the bottom half* and *on the interrupt side* are sometimes used as if they are synonyms, but they really are not synonyms. Due to real-time concerns, sometimes routines on the call side are placed in the bottom half.
- Kernel routines are difficult to debug.
There is a kernel debugger, but it is not as easy to use as is `dbx`. References to improper addresses may cause data corruption (because the kernel is privileged) or a system crash. It is possible to trace the execution of a device driver with `printf`; but `printf` only prints to a native serial port, cannot be used in an interrupt handler, and may affect a device driver's timing.
- Execution of kernel routines, in the context of a process, can be preempted by the scheduler in favor of a process with greater priority.
This means that device driver routines cannot depend on disabling interrupts as being sufficient to avert concurrent access to shared resources. It also means that drivers (not just interrupt handlers) must be reentrant.
- The AIX kernel is dynamically extendible.
Object files acceptable to the kernel loader, are bound into the AIX kernel while the computer is still operating; there is no need to restart the system. A device driver's routines are typically linked to form two object files (one for the top half, the other for the bottom half), that can be loaded or unloaded by a user program invoking the kernel loader with the `sysconfig` system call. Loading (*configuring*) files into the kernel is

extending the kernel A kernel extension, such as a device driver, is configured into the kernel while starting the system, or while the system is operating.

The AIX Kernel figure summarizes some aspects of the AIX kernel that affect device drivers. Note that AIX enables *kernel processes* processes that execute entirely in kernel mode.



AIX Kernel

How Device Drivers Are Accessed

In many operating systems, a user who wants to transfer data to a device must execute a command specific to that device. In such an operating system, writing to tape requires a different command than writing to a terminal, or writing to a file.

A feature of any UNIX operating system, such as AIX, is that I/O to a device is made to look like I/O to a file in the system's directory tree. A device is made ready for I/O by having its corresponding file (usually placed in the directory */dev*) opened, and data is read from the device, or written to the device, by invoking **read** and **write** system calls on the corresponding file. The device (for example, a printer) is freed for access by other software by closing the file (called a *device special file*) associated with the device.

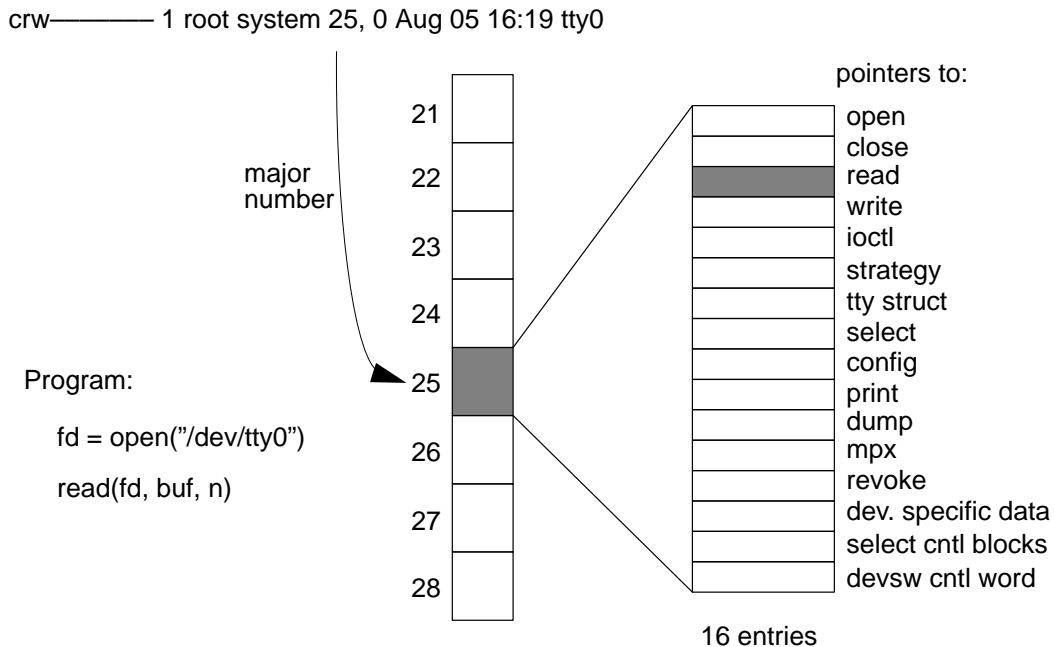
UNIX avoids the need to pass a device-specific parameter to system calls so that UNIX can present the same device interface to any user program accessing different devices. It does this by keeping device-specific data in the inode of the device special file. Such data includes:

- A flag marking the file as *special*
- A major number identifying which device driver is to be invoked (such as for a tape drive or printer)

- A minor number identifying a particular device among the several devices handled by the device driver associated with the major number (for example, selects tape1 or printer3).
- A flag marking the device type as *character* or *block*

The file, created by the **mknod** system call, is marked *special* so that system calls know to access a device, and not a file on disk.

The major number serves as an index into an array of structures. Each structure contains pointers to functions to be invoked when opening, closing, reading, writing, or performing whatever device operation the program requires. The array of these structures is called a *device switch table*. The Device Switch Table figure illustrates these structures.



Device Switch Table

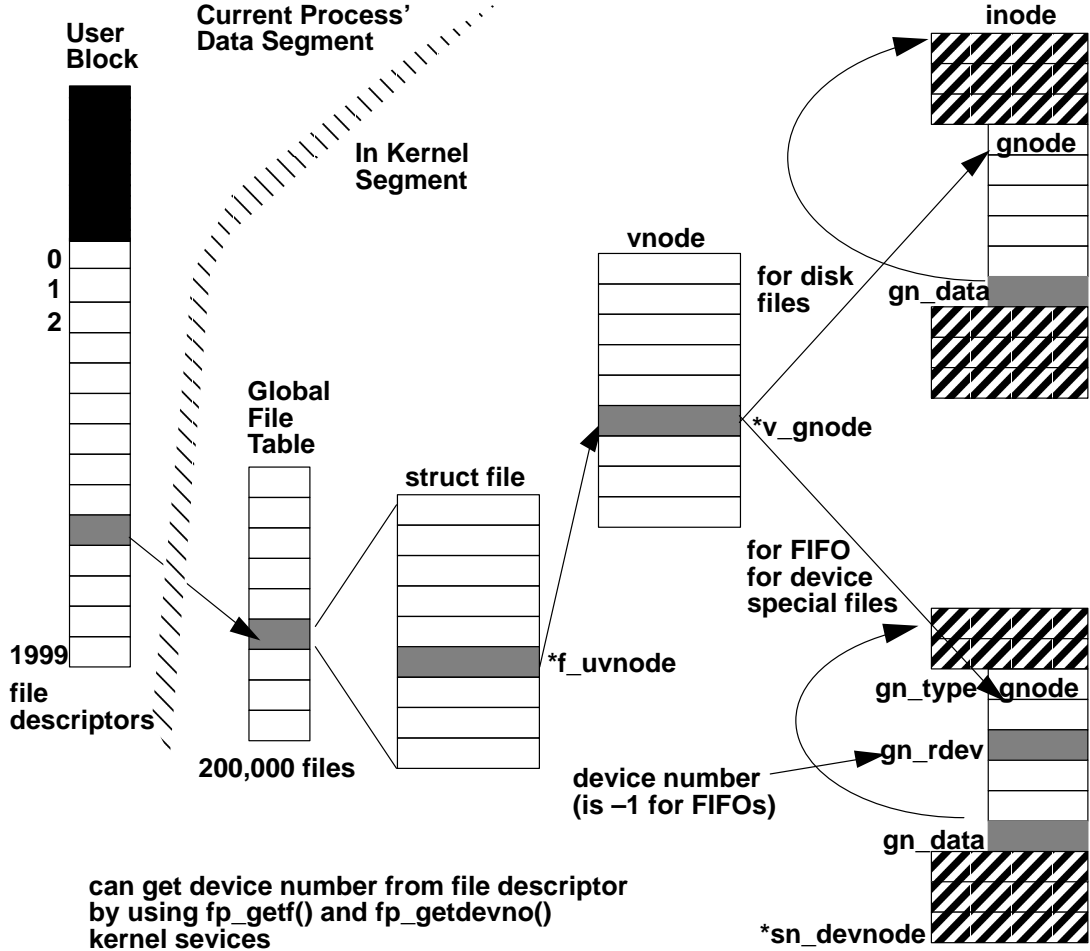
In AIX, the device switch table can have up to 256 such structures. When the driver is configured into the kernel, each function pointer in the structure is assigned the virtual address of the first instruction of an associated routine, that is, the virtual address of a routine's entry point. If a device driver does not define a particular routine (no drivers define all of them), use either a pointer to the function **nodev**, which returns ENODEV (in **sys/errno.h**), or a pointer to the function **nulldev**, which returns NULL.

In effect, the major number specifies which device driver to invoke. If there is a serial port with one driver, and, say, a multiport serial adapter which requires another driver, then each would have their own unique major number. On the other hand, if you have a high-density tape drive attached to one adapter, and a different, low-density tape drive attached to another adapter, but a driver supports both kinds of adapters, then there is only one major number needed to access either tape drive.

The minor number is used to distinguish between devices supported by the same driver. It typically serves as an index into an array, maintained by the driver, of structures containing device-specific information. For example, a terminal driver would need to keep track of the various baud rates or parity settings of each terminal.

Once a program has opened a file, it uses the file descriptor to determine the *device number* which combines the device's major and minor number into one integer. The program that

configures the device driver into the kernel allocates a device number that is unique for the system. The figure From File Descriptor to Device Number shows the data structures involved in determining a device number from a file descriptor associated with the device's special file.



From File Descriptor to Device Number

On some UNIX systems, there is a device switch table for character drivers, and one for block drivers, hence the need for a flag **c** or **b** to mark the distinction. Some UNIX systems access a stream head (a generic interface to STREAMS modules or drivers) through a third device switch table. However, in AIX, the entry points to character drivers, block drivers, and stream heads are all invoked through the same device switch table. The kinds of routines that are expected to be included as entry points for a device driver vary with the type of device driver needed.

Types of Device Drivers

There are three types of device drivers in AIX:

- A block device driver supports a device that reads and writes data buffers of a (large) fixed size

- A STREAMS device driver has some routines that are either invoked from a stream head, or from a STREAMS module, instead of from the device switch table.
- A character device driver is any driver that is not either of the other two types.

Block Device Drivers

Devices usually supported by a block device driver include: hard disk drives, diskette drives, CD-ROM readers, and tape drives. Block device drivers often provide two ways to access a block device:

raw access The buffer supplied by the user program is to be pinned in RAM as is.

block access The buffer supplied by the user program is to be copied to, or read from buffers in the kernel.

If the block device is accessed as *raw*, the driver can copy data from the pinned buffer to the device. In this case, the size of the buffer supplied by the user must be equal to, or some multiple of, the device's block size. The special file's name is usually prefixed by the letter *r* so that a user can tell which access type the block device has. For example, the name of a diskette drive's raw block special file is **rfd0**, and a special file name for a tape drive is **rmt0**.

Sometimes the term *character mode access* is used to mean raw access.

Otherwise, the block device is accessed as *block*. In this case, a **write** system call to such a device returns once the user buffer is copied to buffers in the kernel segment. The **write** call is asynchronous since the data in the kernel buffer is written out to the block device sometime after the **write** call returns. The size of the user buffer need not be a multiple of the device block size.

The term *buffer cache* refers to a collection of kernel buffers that are manipulated by some kernel services specifically associated with block devices. Although UNIX block device drivers have traditionally made use of the buffer cache, it is rarely used in AIX because buffering is more frequently done by memory mapping regions of the kernel segment with frames in RAM.

A block device driver with a block access type method is a driver that also provides a *strategy* routine to arrange accesses to device blocks so that overall access time is minimized. The strategy entry point is not invoked from a user program; rather, the entry point, which is in the device switch table, can be invoked by either of the following:

- Off-level interrupt handlers responsible for writing the buffer cache out to disk.
- The AIX Virtual Memory Manager to perform paging, that is, to page space for working segments, to disk files for memory-mapped files.

The block driver's strategy routine is intended for reading and writing buffers that are not necessarily contiguous on the device itself.

A tape device driver has a raw access method, but has no block access method. There is no use of kernel buffers, and there is no reason to provide a strategy entry point since tape does not lend itself to efficient random access.

For more information on implementing a block device driver, see Chapter 7, "Block Device Drivers."

Block devices are often intended to contain a file system. A block device driver that interfaces to the Logical Volume Manager (LVM) enables the block device to support a journaled file system (JFS). For more information on this, see "Logical Volume Programming," in *AIX Version 4.1 General Programming Concepts*, *Volume 1: Writing Programs*, and "Understanding Physical Volumes and the Logical Volume Device Driver" in *AIX Version 4.1 Kernel Extensions and Device Support Programming Concepts*.

For a block device to contain a file system that is not provided with the operating system, kernel routines that interface between the virtual file system (VFS) and the block device driver must be provided. For more information on this, see “Virtual File Systems,” in *AIX Version 4.1 Kernel Extensions and Device Support Programming Concepts*, Chapter 9, “Writing a Virtual File System.”

STREAMS Device Drivers

Devices that may be supported by a STREAMS driver include: any device connected to the serial port (such as a terminal), or any device attached to a LAN or WAN (such as an Ethernet adapter).

Such devices lend themselves to support from STREAMS drivers because the STREAMS facility is flexible and modular. These qualities that are well suited to implementing communication protocols.

Since the TTY subsystem in AIX Version 4.1 consists of STREAMS modules, if you want to support terminal processing from a serial adapter you must provide a STREAMS driver. For more information on this, see Chapter 10, “STREAMS-Based TTY Subsystem Interface.” Chapter 11, “Implementing Graphical Input and 2-D Graphics Device Drivers”, contains related information about the low-function terminal (LFT) subsystem that supports use of a console display.

AIX provides a STREAMS driver, the Data Link Protocol Interface (DLPI), which supports some LAN adapters. For more information on this, see Chapter 12, “Implementing a Network Device Driver.”

A stream is a linked list of kernel modules, and consists of a stream head at one end of the list and a STREAMS device driver at the other. To visualize how a stream works, see the STREAMS Driver Entry Points figure on page 1-13.

The stream head (supplied with the operating system as part of STREAMS, a device driver writer does not need to write a stream head) contains some routines that are invoked from the device switch table, so the stream head is associated with a device special file in the AIX file tree.

A STREAMS driver has some routines that are either invoked by the stream head, or by a STREAMS module that had been inserted into the stream between the stream head and the STREAMS driver. The driver may, or may not, have any routines that are invoked from the device switch table.

For more information on this, see “STREAMS,” in *AIX Version 4.1 Communications Programming Concepts* and *UNIX System V Release 4, Programmer's Guide: STREAMS*

Character Device Drivers

Devices that are supported by a character device driver include any device that reads or writes data a character at a time (such as printers, sound boards, or terminals). Also, any driver that has no associated hardware device (called a *pseudo-driver*) is treated as a character device driver. For example, **/dev/mem**, **/dev/kmem**, and **/dev/bus0** are *character pseudo-drivers*.

Graphics input devices and graphics capable displays are often supported by character device drivers. For more information on how to implement such drivers, see Chapter 11, “Implementing Graphical Input and 2D Graphics Device Drivers.”

A character special file that has the mode flag S_ISVTX (also called the *sticky bit* because it causes the text of an executable to remain in memory after use) set, is a multiplexed character file. Special files are created and deleted in **/dev** as needed to support multiple ports connected to that adapter. For example, a serial adapter that supports multiple

terminals would need to be a multiplexed character file. A multiplexed device driver contains an additional `xyzmpx` entry point. A pseudo-TTY (PTY) is an example of a multiplexed character device.

Device Driver Configuration

When a version of a device driver is written, to test the driver you need to create and load the driver's object files into the kernel. Sample code in this section, shown in several parts, shows the basic steps of compiling, linking, loading, and testing a pseudo-driver.

Assume that the device special file is `/dev/xyz`. The device driver's object files are usually kept in the directory `/usr/lib/drivers`, but for this simple example, the object file `xyz` is kept in the current directory. You can give a driver's object file any name permitted within an AIX file system.

To configure a device driver object file `./xyz` into the AIX kernel, a user program with root user authority, here written in C, extends the kernel:

```
struct cfg_load cfg;
cfg.path = "./xyz";
sysconfig(SYS_KLOAD, &cfg, sizeof(cfg));
```

Continuing the sample code, the configuring program needs to pick major and minor numbers that are not already in use:

```
majorno = 99;          /* To avoid ODM for now */
minorno = 0;
device_number = makedev(majorno, minorno); /* see sysmacros.h */
```

A program that configures a driver into the kernel does not really just guess which major and minor numbers to use. Associated with the Object Data Manager (ODM) are user routines for determining what numbers to use and for allocating the numbers. Use of ODM is described later, but major number 99 is picked in the example just as an illustration. Calling **makedev** combines the major number and the minor number into one integer, the device number.

A configuring program usually creates the device special file to be associated with the device:

```
mknod ("/dev/xyz", 0666 | _S_IFCHR, device_number);
```

Now, the configuring program invokes the device driver's configure entry point:

```
struct cfg_dd xyzcfg;
xyzcfg.kmid = cfg.kmid; /* kernel module ID from sysconfig */
xyzcfg.devno = device_number;
xyzcfg.cmd = CFG_INIT;
sysconfig(SYS_CFGDD, &xyzcfg, sizeof(xyzcfg));
```

Control now passes to the default entry point for the device driver module. When linking the device driver routines object file, the entry point specified should be the symbolic name of the configuration entry point, in this case `xyzconfig`. A more complete example is shown in the configuration program in "Sample Device Driver" on page 1-14.

STREAMS device drivers are configured into the kernel by invoking the **strload** and **str_install** commands after editing the `/etc/pse.cfg` file. The **strload** command extends the kernel by loading the Portable STREAMS Environment (pse) kernel extension. The **str_install** command extends the kernel by issuing a series of calls to **sysconfig** as indicated by entries in the file `/etc/pse.cfg`. For more information about configuring STREAMS modules and drivers into the AIX kernel, see "STREAMS" in *AIX Version 4.1 Communications Programming Concepts*.

Object Data Manager (ODM) Database

Many UNIX systems have ASCII stanza files that are edited as part of configuring a driver into the kernel. For example, in some UNIX systems you add stanzas to a file such as **/etc/system** or **/etc/master**. In AIX, such stanza files are kept in a directory, **/etc/objrepos**, called the ODM database. In AIX, you do not directly modify the files in the ODM database, but instead you call ODM routines, or execute ODM commands to modify the ODM database. The environment variable **ODMDIR** specifies which directory the ODM routines reference.

The files in an ODM database are indexed ASCII files called ODM object classes. Object classes can be thought of as tables where each row is an object and each column is a field within each object. For example, in the file **CuAt**, there is an object with name `tok0` that has attribute `dma_lvl` and value `0x5`. This object says the Token Ring card associated with `/dev/tok0` has DMA level 5.

The following object classes are significant for device drivers:

PdDv	Predefined (supported) devices. (Not actually installed on the system).
PdAt	Predefined attributes of the predefined devices
PdCn	Predefined connections/dependencies
CuDv	Customized (defined and/or available) devices
CuAt	Customized attributes of system and customized devices
CuDep	Customized dependencies, which devices/subsystems require which others
CuDvDr	Customized device driver resources. For example, ensures unique major numbers.
CuVPD	Customized vital product data (for Micro Channel adapters)
Config_Rules	List of configuration methods for cfgmgr command to execute product inventory LPP history.

A device method is an executable program, usually written in C, that modifies an ODM object class associated with a device. A device method is invoked by a user with root user authority, or when the command **cfgmgr** is called by `rc.boot` in RAM disk when the computer system is started.

The following types of device methods should be provided for a device driver. These methods are usually kept in the directory **/usr/lib/methods**.

- Define method (causes device to be defined)
 - A define method's main task is to retrieve device data from PdDv in ODM and create a CuDv object. Also, it ensures that a parent device exists in the CuDv object.
- Configure method (causes device to be available)
 - A configure method should perform the following steps:
 - a. Display LED value on system LED panel
 - b. Verify that a parent device is available (in ODM)
 - c. Verify that a device is present
 - d. Invoke the **busresolve** system call to get an interrupt level assigned to the device.
 - e. Extend the kernel by calling **sysconfig**.

- f. Generate major and minor numbers and create a special file entry in **/dev**.
- g. Build a device dependent structure (DDS).
- h. Invoke the device driver's config entry point by calling **sysconfig** and passing it the DDS.
- i. Downloading any microcode needed by the adapter.
- j. Updating CuDv with Vital Product Data and making the device available.

In addition to the define method and the configure method, other methods (for example, undefine, unconfigure, start, stop, and change methods) may also be needed.

For more information on configuration methods, see Chapter 6, "Device Driver Configuration" in this book or "Object Data Manager (ODM)" in *AIX Version 4.1 General Programming Concepts Volume 1: Writing Programs*

Device Driver Entry Points

We now consider how control passes from a user program to a device driver entry point associated with the system call that the user program invoked. As shown in the Device Switch Table figure, on page 1-4, there are up to eleven entry points listed in the device switch table for a particular driver: open, close, read, write, ioctl, strategy, select, config, dump, mpx, and revoke. The entry point, "print," is not used. Even though the strategy and dump driver routines have entry points in the device switch table, they cannot be invoked by a system call in a user program.

Those routines which can be invoked from a system call in a user program, whose entry points are listed in the device switch table, are collectively known as the *device head*, and are said to perform the *device head role* within the driver. These routines are expected to return control to the user program that invoked them once their task is complete. Even though these routines are placed in the same object module, they rarely interact with each other (for example, they don't call each other). These routines merely share resources such as kernel data structures and the device itself. Like any program, a device driver can define other routines as needed, routines that may be invoked by any other routine in the driver.

For information on how to write each routine associated with an entry point in the device switch table, please refer to a complete list of such routines "Device Driver Operations" in *AIX Version 4.1 Technical Reference, Volume 5: Kernel and Subsystems*. The following discussion focuses on routines common to most drivers.

Note that a routine's entry point is customarily labeled by prefixing a routine's function with some device-specific abbreviation based on the device special file name. For example, the entry point associated with a routine that opens a terminal device is labeled, **ttyopen**, and one that closes the device is, **ttyclose**. But, one is free to label entry points with any symbol that a compiler and linker will accept.

xyzconfig Entry Point

A device driver's configuration entry point is called when a program directs the kernel loader to configure the device driver's object file into the kernel. The configuration entry point is also called when the device driver is being removed from the kernel or when certain data is queried from the device.

Below are the commands one may pass to the **sysconfig** system call:

- CFG_INIT

In this case, the tasks that the configuration routine might perform are such things as, placing entry points to other routines into the switch table by invoking **devswadd**, or initializing and allocating kernel data structures associated with device driver, or initializing the adapter, or downloading microcode to the card attached to the device.

Also, the **sysconfig** system call usually passes a Device Dependent Structure (DDS) to the configuration routine when initializing the device. One defines the DDS in whatever fashion as is necessary for the driver. For example, a serial device driver might require the DDS to contain initial values for baud rate, bits per character, parity bit settings, and so on.

- **CFG_TERM**

In this case, the routine checks for any outstanding open file descriptors and releases any associated resources.

- **CFG_QVPD**

In this case, the routine returns Vital Product Data from the card attached to the device.

xyzopen and xyzclose Entry Points

These routines usually perform the following functions:

- Allocate or free resources for this device instance

This is often where data structures are allocated or freed from the kernel heap.

This is often where the bottom half is pinned or unpinned, and this is where interrupt handlers are often registered for use by the kernel, or removed from use by the kernel.

- Update use counts and semaphores if exclusive access required

The **xyzopen** routine is invoked by **open** or **creat** system calls issued from a user program, or is invoked from an **fp_open** or **fp_opendev** kernel service call issued from a kernel extension.

The **xyzclose** routine is invoked by the **close** system call issued from a user program, or is invoked from the **fp_close** kernel service call issued from a kernel extension.

xyzread Entry Point

This routine usually does the following:

- Returns a buffer of whatever data was collected from the device (via the device driver's interrupt handler)

A call to the **read** system call from a character device amounts to transferring data from a kernel buffer that had been populated by this device's interrupt handler to a buffer supplied by the user (which is usually outside that user's data segment)

- From a block device, initiates block I/O requests via the **uphysio** kernel service

This is referred to as *raw access* since no kernel buffers are being used for reading.

- From a streams device, causes the stream head or module to invoke the device driver's **xyzput** entry point

Depending on how the device special file was opened, the **xyzread** routine usually puts the calling process to sleep until the data requested is available.

The **xyzread** routine is invoked by the **read** system call issued from a user program, or is invoked by the **fp_read** kernel service call issued from a kernel extension.

xyzwrite Entry Point

This routine usually does the following:

- Outputs data to a block device
An explicit write to a block device, one not using the **xyzstrategy** routine, is raw access. This initiates block I/O requests via the **uphysio** kernel service
- Outputs data to a character device
A write to a character device transfers data from a buffer supplied by the user, which is usually out of that user's data segment, pointed to by the **uio** structure, to any buffer needed by the device adapter, usually one character at a time.
- Outputs data to a STREAMS device
A write to a STREAMS device causes the stream head (or some module in the Stream) to invoke the streams driver's **xyzwput** routine.

The **xyzwrite** routine is invoked by the **write** system call issued from a user program, or is invoked by the **fp_write** kernel service routine issued from a kernel extension.

xyzstrategy Entry Point

This routine usually schedules read or write requests of a block device. Such requests are added to a queue of pending I/O requests for the device. The queue can be sorted to optimize device access; for example, one may wish to have disk blocks organized so that any that are within the same cylinder (under the drive's read/write head) are input or output in one operation.

The buffer supplied to the **xyzstrategy** routine must be pinned in RAM, because the actual I/O with the buffer is asynchronous (it may happen after the routine exits).

This routine calls the **iodone** kernel service once it's finished.

The **xyzstrategy** routine is invoked indirectly by the following:

- The **uphysio** kernel service.
- The Logical Volume Manager (LVM).
- The Virtual Memory Manager (VMM). (For example, for handling page faults.)

More information about use of strategy routines for block devices, see Chapter 7, "Block Device Drivers."

xyzioctl Entry Point

This routine usually performs functions that are not done by any of the other routines mentioned up to now. Usually, the **xyzioctl** routine modifies, or inspects the state of the attached device, and reports any results back to the calling program. For example, a terminal driver would have its **ttzioctl** routine enable the calling program to modify baud rate, bits per character, and so on.

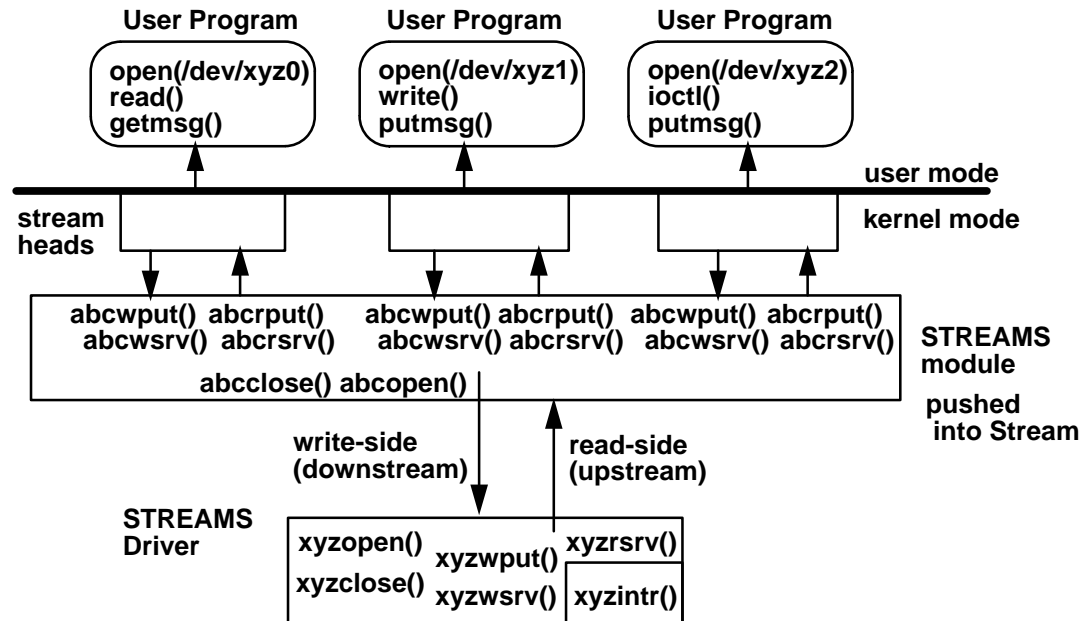
This routine need not be synchronous since a device may not permit immediate action. For example, a program that calls **ioctl(CIOSTART)** on a LAN adapter requires a subsequent **ioctl(CIOGETSTAT)** to determine whether the first **ioctl** call is complete.

Not every driver has this entry point, but if this routine is present, it must support the **IOCINFO** command option which tells **xyzioctl** to return a structure that describes the attached device.

The **xyzioctl** routine is invoked by the **ioctl** system call issued from a user program, or is invoked by the **fp_ioctl** kernel service call issued by a kernel extension.

STREAMS Entry Points

A STREAMS driver may include routines invoked directly from the device switch table. However, a STREAMS driver usually handles all processing by using a write-side (downstream) put routine, and one or two optional service routines. Like any device driver, a STREAMS driver also has an interrupt handler which performs functions similar to those performed by a read-side (upstream) put routine of a STREAMS module. These relationships are shown in the STREAMS Driver Entry Points figure.



STREAMS Driver Entry Points

For more information on how to write a STREAMS module or device driver, see *UNIX SYSTEM V Release 4, Programmer's Guide: STREAMS*

xyzwput Entry Point

The write-side put routine receives messages from the stream head or STREAMS module upstream. The stream head converts **write** and **ioctl** system calls, issued from the user program, into messages, and then sends the messages downstream by invoking the write-side put routine of whatever STREAMS module or driver is downstream. The stream head handles the **read** system call issued from the program; the STREAMS driver does not.

This routine is invoked by the stream head or STREAMS module upstream from the driver.

A STREAMS driver does not have a read-side put routine. The interrupt handler assumes the role of receiving data from the device's adapter.

xyzsvr, xyzrsrv Entry Points

A STREAMS driver can optionally support either a write-side or read-side service routine, neither, or both.

When a driver's write-side put routine determines that it must defer writing data to the device (flow control), it must place the data on a write-side queue. A kernel process, which is part of the Portable STREAMS Environment (PSE), eventually invokes the driver's write-side

service routine, which checks the write-side queue and attempts to write the data to the device, or enqueues the data on the write-side queue so it can make another attempt later.

Similarly, the device driver's interrupt handler may place a buffer of data coming in from the device on a read-side queue. The PSE kernel process eventually invokes the driver's read-side service routine, which checks the read-side queue and then processes the data. Once the processing is complete, the read-side service routine invokes the read-side put routine of the STREAMS module upstream.

The data to be handled by a service routine is placed on a queue, as follows:

- The write-side put routine must call **putq** to place messages on a write-side service queue, so that invoking the **xyzwsrv** routine later will process that message.
- The interrupt handler must call **putq** to place messages on a read-side service queue for deferred processing by the **xyzsrv** routine.
- A service routine gets messages off its own queue by calling **getq**.

Service routines are invoked by the STREAMS scheduler, implemented within a kernel process in AIX. The STREAMS scheduler is part of the **PSE** (Portable STREAMS Environment) kernel extension.

Sample Device Driver

For greatest simplicity, consider a pseudo-driver (one having no associated device), whose routines are minimal and yet demonstrate an outline of what a device head looks like. As with any sample code in this book, the following warning applies.

Warning: The source code examples provided by IBM are only intended to assist in the development of a working software program. The source code examples may not function as written: additional code is required. In addition, the source code examples may not compile and may not bind successfully as written. International Business Machines corporation provides the source code examples, both individually and as one or more groups, "AS IS" without warranty of any kind, either expressed or implied, including, but not limited to the implied warranties of merchantability and fitness for a particular purpose. The entire risk as to the quality and performance of the source code examples, both individually and as one or more groups is with you. Should any part of the source code examples prove defective, you (and not IBM or an authorized RISC System/6000 Workstation dealer) assume the entire cost of all necessary servicing, repair, or correction. IBM does not warrant that the contents of the source code examples, individually or as one or more groups, will meet your requirements or that the source code examples are error-free. IBM may make improvements and/or changes in the source code examples at any time. Changes may be made periodically to the information in the source code examples; these changes may be reported, for the sample device drivers included herein, in new editions of the examples. References in the source code examples to IBM or non-IBM products, programs, or services shall not be viewed as an endorsement of any kind and references do not imply that IBM intends to make these available in all countries in which IBM operates. Any reference to an IBM licensed program in the source code examples is not intended to state or imply that only IBM's licensed program may be used. Any functionally equivalent program may be used.

Files for Sample XYZ Device Driver

The sample device driver has the following files located in a user directory:

- | | |
|-------------------|---|
| aprogram.c | User program's source: opens, reads, writes to device |
| makefile | Command file that directs the make command |

xyz.c Source code for the device driver
xyz_cfg.c Source code for the configure program

makefile for Sample XYZ Device Driver

This file, makefile, contains commands for building the sample XYZ device driver:

```
# Once this is done, have root user run xyz_cfg -q to query the kernel
# to see that the driver is not loaded. Then run xyz_cfg -l to load it.
# Check again with xyz_cfg -q; if OK, then have non-root user run "aprogram."
# Then clean up by running xyz_cfg -u and verify absence of /dev/xyz.

#needed to get kernel services like devswadd()
KSYSLIST=/lib/kernex.exp

# needed for trace macro (containing a system call)
SYSLIST=/lib/syscalls.exp

all: aprogram xyz xyz_cfg

# import the kernel service calls and make xyzconfig() the entry point
xyz: xyz.o
    ld -e xyzconfig -o xyz -BI:$(KSYSLIST) -BI:$(SYSLIST) xyz.o

# It is necessary to use separate compile and link steps to avoid picking up
# routines like printf from libc.a. There is a kernel printf().
# _KERNEL needed to get trace macro, others determined by header files
xyz.o: xyz.c
    cc -c -D_ALL_SOURCE -D_POSIX_SOURCE -D_KERNEL xyz.c

# this is to create the driver's configure method that extends the kernel
xyz_cfg: xyz_cfg.c
    cc -o xyz_cfg xyz_cfg.c

aprogram: aprogram.c
    cc -o aprogram aprogram.c
```

Configuration Program for Sample XYZ Device Driver

This configuration program, xyz_cfg.c, is *not* recommended for configuring drivers in AIX. It avoids the use of ODM routines for simplicity.

```
/*
 * FUNCTION: Configure/Unconfigure program for bare bones driver, xyz
 * Normally, one has a configure program (run at boot time)
 * and then a separate unconfigure program (run interactively).
 * Run "xyz_cfg" to see parameters needed.
 */

#include <stdio.h>            /* for printf() */
#include <unistd.h>          /* for getopt() */
#include <stdlib.h>          /* for exit() */
#include <sys/types.h>        /* for dev_t and other declarations */
#include <sys/errno.h>       /* for perror() */
#include <sys/sysmacros.h>   /* for makedev() */
#include <sys/sysconfig.h>   /* for sysconfig() */
#include <sys/device.h>      /* for CFG_INIT & other flags */
#include <sys/mode.h>        /* for mknod() */

void main(int argc, char *argv[])
{
    struct cfg_load cfg;       /* to load kernel extension */
```

```

struct cfg_dd xyzcfg;          /* to invoke xyzconfig() */
int majorno, minorno;
dev_t device_number;
int ch;                       /* flag char returned by getopt */

extern int optind;            /* for getopt function */
extern char *optarg;         /* for getopt function */

/* normally call ODM routines to get values for these */
majorno = 99;
minorno = 0;
device_number = makedev(majorno, minorno);

/* parse command line */

if(argc <= 1)
{ printf("You must give an argument.\n");
  printf("arguments to xyz_cfg are:\n");
  printf("\t-l to load the driver and invoke xyzconfig() \n");
  printf("\t-u to invoke xyzconfig() and unload the driver \n");
  printf("\t-q to query the status of the kernel extension\n");
}

while ((ch = getopt(argc,argv,"luq")) != EOF)
{ switch (ch)
  { case 'l': /* load the driver--assume is first time (shouldn't!) */
    cfg.path = "./xyz"; /* path is local--usually /etc/drivers */
    if (sysconfig(SYS_KLOAD, &cfg, sizeof(cfg)) == -1)
    { perror("sysconfig SYS_KLOAD FAILED");
      exit(1);
    }
    xyzcfg.kmid = cfg.kmid; /* kernel module ID from SYS_KLOAD */
    xyzcfg.devno = device_number;
    xyzcfg.cmd = CFG_INIT;
    if (sysconfig(SYS_CFGDD, &xyzcfg, sizeof(xyzcfg)) == -1)
    { perror("sysconfig SYS_CFGDD FAILED");
      exit(1);
    }

    /* make /dev entry in honor of device switch table entry */
    if (mknod("/dev/xyz", 0666 | _S_IFCHR, device_number) == -1)
    { perror("mknod FAILED");
      exit(1);
    }

    break;

    case 'u': /* unload the driver */

    /* the kmid lost once this exits, so we requery the information */
    cfg.path = "./xyz";
    if (sysconfig(SYS_QUERYLOAD, &cfg, sizeof(cfg)) == -1)
    { perror("sysconfig SYS_QUERYLOAD FAILED");
      exit(1);
    }

    xyzcfg.kmid = cfg.kmid; /* kernel module ID from SYS_QUERYLOAD */
    xyzcfg.devno = device_number;
    xyzcfg.cmd = CFG_TERM;
    if (sysconfig(SYS_CFGDD, &xyzcfg, sizeof(xyzcfg)) == -1)
    { perror("sysconfig SYS_CFGDD FAILED");
      exit(1);
    }

    /* remove /dev entry...normally would use ODM's reldevno() */
    unlink("/dev/xyz");

```

```

    if (sysconfig(SYS_KULOAD, &cfg, sizeof(cfg)) == -1)
    { perror("sysconfig SYS_KULOAD FAILED");
      exit(1);
    }
    break;

    case 'q': /* query the status of the system call */
    cfg.path = "./xyz";
    if (sysconfig(SYS_QUERYLOAD, &cfg, sizeof(cfg)) == -1)
    { perror("sysconfig SYS_QUERYLOAD FAILED");
      exit(3);
    }
    printf("The kernel module ID is %d\n", cfg.kmid);

    xyzcfg.kmid = cfg.kmid; /* kernel module ID from SYS_QUERYLOAD */
    xyzcfg.devno = device_number;
    xyzcfg.cmd = CFG_QVPD;
    if (sysconfig(SYS_CFGDD, &xyzcfg, sizeof(xyzcfg)) == -1)
    { perror("sysconfig SYS_CFGDD FAILED");
      exit(1);
    }
    break;

    default:
    printf("arguments to xyz_cfg are:\n");
    printf("\t-l to load the driver and invoke xyzconfig() \n");
    printf("\t-u to invoke xyzconfig() and unload the driver \n");
    printf("\t-q to query the status of the kernel extension\n");
  } /* end switch on ch */
} /* end while getopt */
exit(0);
}

```

Source Code for Sample XYZ Device Driver

This file, `xyz.c`, contains source code for the sample device driver:

```

#include <sys/types.h> /* for dev_t and other types */
#include <sys/errno.h> /* for errno declarations */
#include <sys/sysconfig.h> /* for sysconfig() */
#include <sys/device.h> /* for devsw */
#include <sys/trchkid.h> /* for trace hook macros */

/*****
BARE BONES DRIVER

This shows the format of a minimal set of entry points for a
pseudo-driver.
*****/

/***** xyzopen *****/
int xyzopen(dev_t devno, ulong devflag, chan_t chan, int ext)
{
    TRCHKL5T(HKWD_USER1, 0x7, devno, devflag, chan, ext);
    return(0);
}

/***** xyzclose *****/
int xyzclose(dev_t devno, chan_t chan)
{
    TRCHKL3T(HKWD_USER1, 0x8, devno, chan);
    return(0);
}

```

```

/***** xyzread *****/
int xyzread(dev_t devno, struct uio *uiop, chan_t chan, int ext)
{
    TRCHKL5T(HKWD_USER1, 0x9, devno, uiop, chan, ext);
    return(0);
}

/***** xyzwrite *****/
int xyzwrite(dev_t devno, struct uio *uiop, chan_t chan, int ext)
{
    TRCHKL5T(HKWD_USER1, 0x0a, devno, uiop, chan, ext);
    return(0);
}

/***** xyzconfig *****/
int xyzconfig(dev_t devno, int cmd, struct uio *uiop)
{
    struct devsw dsw_struct;
    extern int nodev();
    int return_code;

    /* trace macro to print received parameters in hex */
    TRCHKL4T(HKWD_USER1, 0x1, devno, cmd, uiop); /* 0x01-tracept label */

    switch(cmd)
    {
    case CFG_INIT:
        dsw_struct.d_open      = xyzopen;
        dsw_struct.d_close    = xyzclose;
        dsw_struct.d_read     = xyzread;
        dsw_struct.d_write    = xyzwrite;
        dsw_struct.d_ioctl    = nodev;
        dsw_struct.d_strategy = nodev;
        dsw_struct.d_ttys     = NULL;
        dsw_struct.d_select   = nodev;
        dsw_struct.d_config   = xyzconfig;
        dsw_struct.d_print    = nodev;
        dsw_struct.d_dump     = nodev;
        dsw_struct.d_mpx      = nodev;
        dsw_struct.d_revoke   = nodev;
        dsw_struct.d_dsdptr   = NULL;
        dsw_struct.d_opts     = NULL;

        if((return_code = devswadd(devno, &dsw_struct)) != 0)
        {
            TRCHKL3T(HKWD_USER1, 0x2, devno, dsw_struct);
            return(return_code);
        }

        /* entry points now in device switch table */
        break;

    case CFG_TERM:
        TRCHKL1T(HKWD_USER1, 0x3);
        if((return_code = devswdel(devno)) != 0)
        {
            TRCHKL2T(HKWD_USER1, 0x4, devno);
            return(return_code);
        }

        break;

    case CFG_QVPD:
        TRCHKL1T(HKWD_USER1, 0x5); /* would normally handle this case too */
        break;

    default:
        TRCHKL1T(HKWD_USER1, 0x6);
        return(EINVAL);
    } /* end switch(cmd) */
}

```

```
    return(0);
}
```

User Program to Invoke Sample XYZ Device Driver

This file, `aprogram.c`, contains sample user code to invoke the sample device driver:

```
#include <stdio.h>          /* for printf() */
#include <fcntl.h>          /* for open(), close() */
#include <unistd.h>         /* for read(), write() */
#include <sys/errno.h>      /* for perror() */

int fd;                    /* file descriptor */
char buf[10];              /* read/write buffer */

void main()
{
    printf("buf pointer: 0x%x\n", buf);
    if((fd = open("/dev/xyz", O_RDWR)) == -1)
    { perror("open /dev/xyz FAILED");
      exit(1);
    }

    if(read(fd, buf, sizeof(buf)) == -1)
    { perror("read FAILED");
      exit(1);
    }

    if(write(fd, buf, sizeof(buf)) == -1)
    { perror("write FAILED");
      exit(1);
    }

    if(close(fd) == -1)
    { perror("close FAILED");
      exit(1);
    }
}
```

Running the Sample XYZ Device Driver

The lines prefixed by a # were commands executed as a user with root authority in a window. The lines prefixed by a > were commands executed as a staff user in another window. The order of the commands is as listed.

```
# trace -j'010' -l -s -a &
[1]      21971
# xyz_cfg -q
The kernel module ID is 0
sysconfig SYS_CFGDD FAILED: No such device
[1] + 21971      Done      trace -j'010' -l -s -a &
# xyz_cfg -l
# xyz_cfg -q
The kernel module ID is 21790464
> $ ls /dev/xyz
> /dev/xyz
# chmod 666 /dev/xyz
> $ aprogram
> buf pointer: 0x200516c0
# trcstop
# trcrpt -O'exec=y' -O'pid=n' -O'svc=y' -O'timestamp=1' > $HOME/frog
> $ ls /dev/xyz
> /dev/xyz
# xyz_cfg -u
# xyz_cfg -q
The kernel module ID is 0
sysconfig SYS_CFGDD FAILED: No such device
#
> $ ls /dev/xyz
> ls: 0653-341 The file /dev/xyz does not exist.
```

Trace Output for Sample XYZ Device Driver

Here is the output of **trcrpt** (abbreviated for space):

The hook data indicates calls to xyzconfig (load and query), open, read, write, and close.

ID	PROCESS NAME	I	SYSTEM CALL	ELAPSED	KERNEL	INTERRUPT
001	trace			0.000000	TRACE ON	channel 0
010	trace			19.362228	UNDEFINED	TRACE ID idx 0x21dc traceid 0010
			hookword 10E0000 type 0E			
			hookdata 0000 00000001 00630000 00000001 2FF97F1C 00000000			
010	trace			25.716153	UNDEFINED	TRACE ID idx 0x2230 traceid 0010
			hookword 10E0000 type 0E			
			hookdata 0000 00000001 00630000 00000003 2FF97F1C 00000000			
010	trace			25.716159	UNDEFINED	TRACE ID idx 0x224c traceid 0010
			hookword 10A0000 type 0A			
			hookdata 0000 00000005			
010	trace			99.695506	UNDEFINED	TRACE ID idx 0x24e8 traceid 0010
			hookword 10E0000 type 0E			
			hookdata 0000 00000007 00630000 00000003 00000000 00000000			
010	trace			99.706120	UNDEFINED	TRACE ID idx 0x2504 traceid 0010
			hookword 10E0000 type 0E			
			hookdata 0000 00000009 00630000 2FF97DC0 00000000 00000000			
010	trace			99.706318	UNDEFINED	TRACE ID idx 0x2520 traceid 0010
			hookword 10E0000 type 0E			
			hookdata 0000 0000000A 00630000 2FF97DC0 00000000 00000000			
010	trace			99.706446	UNDEFINED	TRACE ID idx 0x253c traceid 0010
			hookword 10E0000 type 0E			
			hookdata 0000 00000008 00630000 00000000 00000000 00000000			
002	trace			109.684978	TRACE OFF	channel 0

Routines on the Interrupt Side

When one of the system's processors receives an external interrupt, an AIX processor interrupt handler, for that particular processor interrupt, begins execution. This portion of the AIX kernel determines which device interrupt handler, or which collection of handlers, to invoke. For more information on interrupt processing, see Chapter 3, "Interrupts."

A card on the system bus that serves as an adapter between the system and the device may generate an interrupt on the bus. Its device driver will need to configure a routine to handle that interrupt, and may wish to add other routines to handle *off-level* interrupts, which are scheduled by the device interrupt handler as a way to defer interrupt processing.

xyzintr Entry Point

This routine is configured into the kernel by a call to the `i_init` kernel service. The `xyzopen` routine typically calls `i_init`, and pins the bottom half of the device driver. Because routines executing on the interrupt side cannot afford to be preempted by demand paging, they are placed within the driver's bottom half (they are pinned in RAM).

This routine usually does the following:

- Determines the slot of the card that generated the device interrupt.

This routine may also have to determine whether to process the interrupt. The interrupt level can be shared, so the interrupt may be intended for another device interrupt handler.

- Disables some other device interrupts to prevent this routine from being reentered if another interrupt for the same device driver is received.

This technique of avoiding concurrent execution can fail on computers with multiple processors.

Because interrupt handlers have to address concurrency by serializing interrupt processing in some way, the handler's **latency** (the time between when control of the system processor is assumed and when control can be relinquished) must be minimized in order to maximize the system's ability to respond to other device interrupts.

- Transfers data from the device into a buffer, or notifies a user program that something has happened

- Schedules an off-level interrupt handler to complete processing of input data

This can be done to minimize the handler's latency. As an alternative, the driver can depend on kernel processes to handle received data.

There are some special concerns that apply to routines on the interrupt side:

- Many kernel services can only be invoked from routines on the call side.
- These routines cannot go to sleep, though they can post events to (waken) routines on the call side.
- These routines cannot obtain or release locks.

Routines on the call side that are in the bottom half of a driver can adjust locks, provided the events or lockwords are also pinned in RAM.

- An interrupt handler is called at the priority registered when the handler is configured into the kernel via the `i_init` kernel service. The handler's execution can only be preempted by interrupts with a higher priority.

- Interrupt handlers have volatile data (local variables) in a pinned stack that is less than 4K bytes in size.
 - Note:** Because **xmalloc** cannot execute on the interrupt side, any non-volatile buffers that an interrupt handler needs must be previously allocated on the call side.
- Routines that handle off-level interrupts run at a less favored interrupt priority.

Pinning Device Driver Object Files

Sometimes a driver routine (such as an interrupt handler) and any associated data is required to be kept in RAM. This requirement might exist to avoid handling a page fault so the routine can execute within a fixed period of time, or to avoid receiving a page fault while interrupts are disabled.

Typically, a device driver writer compiles or links all routine and data definitions into one loadable file (the bottom half of the device driver), because the **pincode** kernel service marks each page of the loaded object file as being required to be kept in RAM. The method for identifying the object file is that, a pointer to a routine within the object file is an input value for the **pincode** kernel service.

Routines and data that can be subject to page replacement are typically collected into another loadable object (the top half of the device driver).

Routines that wish to allocate buffers from the kernel heap (by calling the **xmalloc** kernel service) must take care which heap the allocation is from. Routines in the bottom half allocate data from the kernel's pinned heap; and routines in the top half allocate data from either heap. For more information on this, see Chapter 4, "Memory Management."

Any routine in the bottom half that invokes a routine in the top half, or refers to data in the top half, negates any purpose for having been pinned in RAM. The routine may (intermittently) cause a page fault.

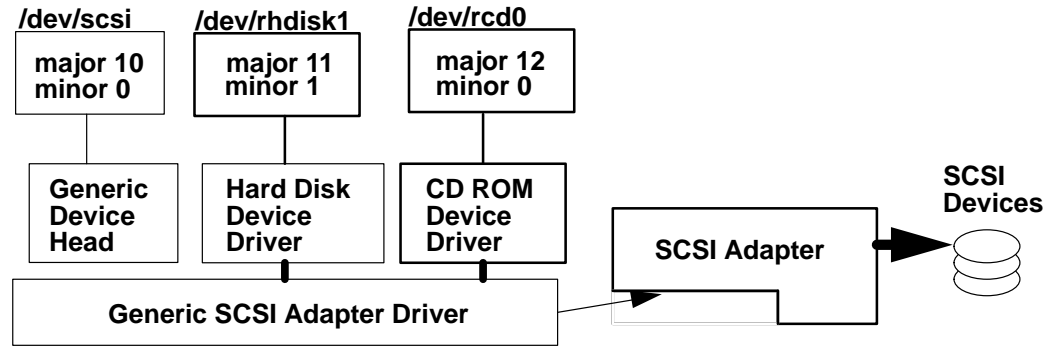
Typically, a device driver's open routine pins the device driver's bottom half when the device special file is first opened. The driver's close routine checks for any outstanding open calls on the file and then calls the **unpincode** kernel service.

Driving a SCSI Attached Device

A SCSI device driver converts I/O requests into SCSI commands, and passes these commands to the SCSI adapter driver provided in AIX. Therefore, a SCSI device driver only consists of routines that execute on the call side. The SCSI device driver passes commands to the adapter driver by direct invocation of its routines; therefore, a SCSI device driver must be integrated with the existing SCSI subsystem. The SCSI subsystem is illustrated in the SCSI Subsystem figure.

The generic SCSI device head associated with the special file `/dev/scsi#` (where # is 0, 1, 2, or some other number) permits a user program with root user authority to send SCSI commands through the associated SCSI adapter.

For more information on how to drive a SCSI attached device, see Chapter 8, "Integrating a SCSI Device Driver into the SCSI Subsystem."



SCSI Subsystem

Other Topics

It may be necessary to integrate a device driver into some existing subsystem in AIX. For example, you may need to integrate a driver into one of the following subsystems:

- TTY subsystem

See Chapter 10, "STREAMS-Based TTY Subsystem Interface," for more information on supporting terminal users through a serial device.
- Graphics device subsystem

See Chapter 11, "Implementing Graphical Input and 2D Graphics Device Drivers," for more information on supporting a pointing device, graphics monitor, keyboard, or other graphics related device.
- TCP/IP

See Chapter 12, "Implementing a Network Device Driver," for more information on supporting the TCP/IP socket interface or the Transport Layer Interface (TLI) through a network (Ethernet, Token Ring, or other) adapter.
- Socket Interface to your own protocol

See Chapter 13, "Network Interfaces and Protocols," for more information on implementing a protocol that is an alternative to TCP/IP, which uses a network adapter driver provided by AIX.

Once a device driver is written and configured into the kernel, it is necessary to test and debug the driver routines. The example driver given in this chapter demonstrates the use of trace macros, but you may wish to incorporate error logging capability into a driver so that a system administrator can be notified of any irregularities.

Furthermore, defects in a driver could cause the system to halt after having copied an image of the kernel (called a *kernel dump*) to a logical volume reserved for that purpose. The kernel dump can be inspected with the **crash** utility. It may also be necessary to debug a driver using the kernel debugger, which comes with AIX and is configured into the kernel with the **bosboot** command. For more information on the use of such tools, see Chapter 14, "Debugging Tools."

Once a driver is ready for delivery, you may wish to archive the resulting binary and character files into a format that the **installp** command can process. For information on

packaging files into **installp** format see “Software Product Packaging” in *AIX Version 4.1 General Programming Concepts Volume 1: Writing Programs*

In addition, you may wish to add an entry in the SMIT interface to enable users to install, configure, or remove your device in the same way that is done with devices supported by AIX. For more information on adding your own SMIT dialogs see “System Management Interface Tool (SMIT) for Programmers” in *AIX Version 4.1 General Programming Concepts Volume 1: Writing Programs*

Chapter 2. Device I/O

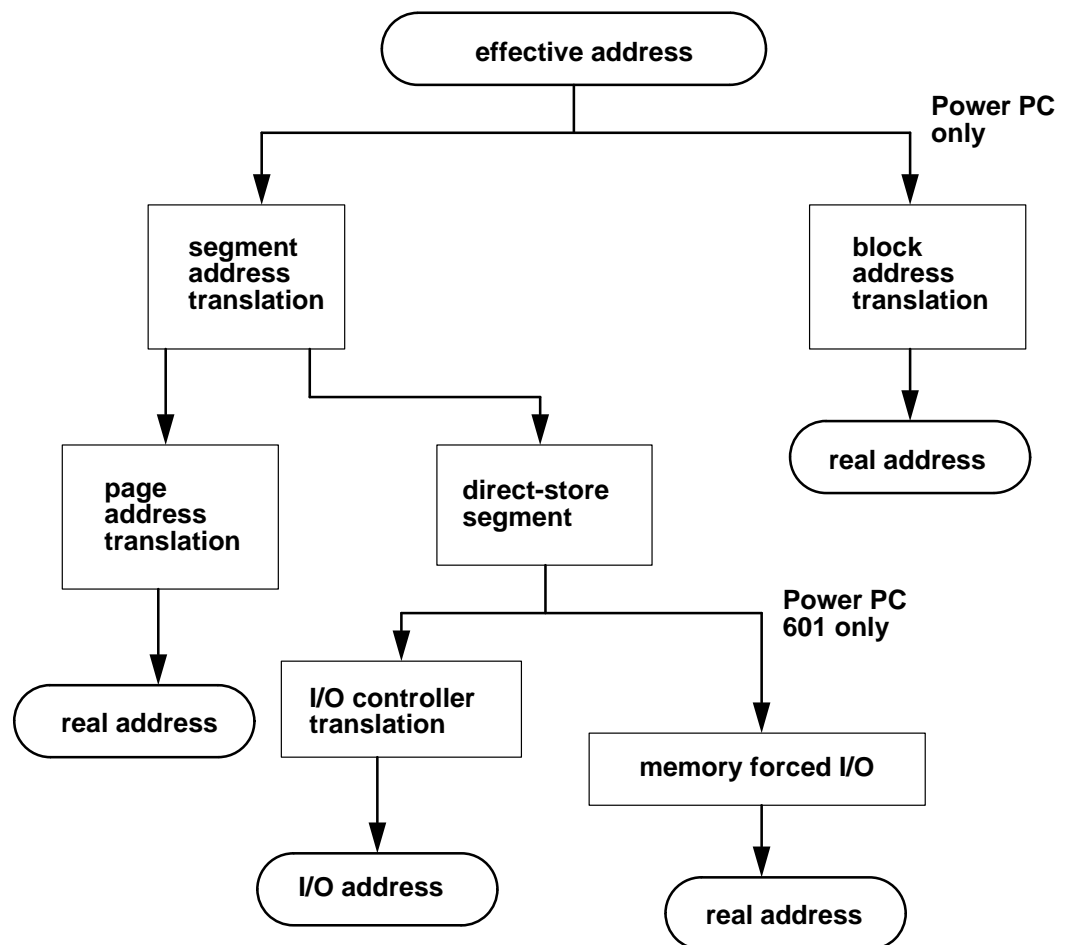
Even though a driver can perform many functions, one usually writes a driver to output data to a device or demand data from a device; in other words drivers usually perform device I/O. A driver may have to read from, or write to, registers on a card serving as an adapter between an I/O bus and a device connected to the card, or the driver may have to set up the means for data to be transferred in some other way.

Address Translation

There are two basic types of address translation implemented on system processors that AIX supports:

- Block address translation (on PowerPC only)
- Segment address translation

The various kinds of address translation are diagrammed in the Address Translation figure.



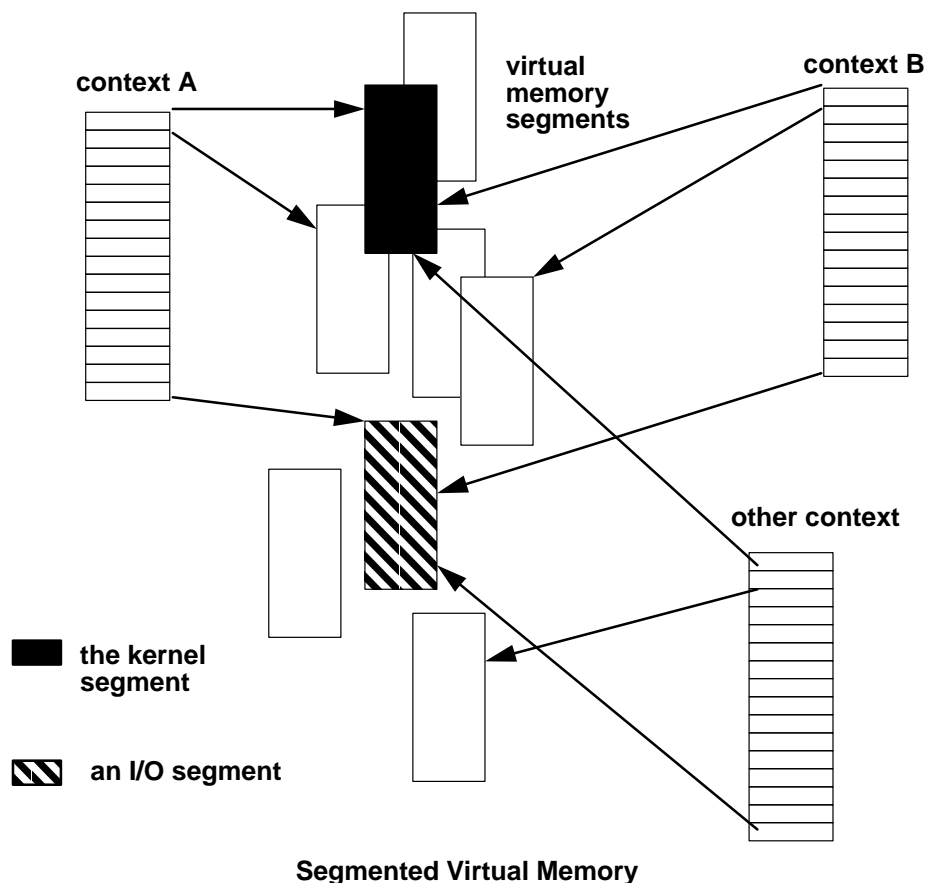
Address Translation

Block Address Translation

Block address translation (BAT) is a feature of the PowerPC architecture that is an alternative to page address translation. Pages have a fixed size, but blocks can be from 128 KB to 256 MB in size, although on the PowerPC 601 RISC Microprocessor, BAT areas can be no larger than 8MB. Whether this type of address translation is implemented in AIX depends on the system PowerPC processor. AIX does not use the BAT tables of the PowerPC 601 RISC Microprocessor, but instead uses memory-forced I/O, which is described later. For more information on block address translation, see *PowerPC Architecture* (order number SR28-5124).

Segment Address Translation

All three processor architectures that AIX supports, POWER, POWER2, and PowerPC, can translate byte addresses by using sixteen segment registers. The contents of these registers, together with that of some other processor registers (such as the general purpose registers, instruction register, and machine state register) make up the context in which an instruction executes. Changing the contents of some of these registers (having saved the former values somewhere) is a *context switch*. A device driver's call side routines typically execute in the context of a process, and its interrupt side routines execute in the context of an interrupt handler. The Segmented Virtual Memory figure shows how virtual memory consists of segments accessed by instructions executing in various contexts.



In AIX, instructions typically execute in a context whose segment registers contain the following:

Segment register 0	Kernel segment identifier (ID) (Shared by all)
Segment register 1	Text segment ID (Where instructions are fetched from)
Segment register 2	Data segment ID (Not shared by others)
Segment register 13	Shared text segment ID (For shared libraries)
Segment register 15	I/O segment ID (For I/O. Often used, but not required.)

For more details, see "Program Address Space Overview" in the chapter about shared libraries and shared memory in *AIX Version 4.1 General Programming Concepts, Volume 1: Writing Programs*

A segment register contains a segment identifier (segment ID) which determines, among other things, what kind of address translation is to be performed on addresses within that segment.

Consider the following kinds of segment address translation:

- Page address translation
- I/O controller interface translation
- Memory forced I/O (PowerPC 601 RISC Microprocessor)

Page address translation is usually performed on an address referenced by an instruction accessing system RAM. For the POWER and POWER2 architectures, this kind of address translation is detailed under "Memory Addressing" in the chapter about system processors in *POWERstation and POWERserver Hardware Technical Information-General Architectures*. For the PowerPC architecture, page address translation is detailed in *PowerPC Architecture* (SR28-5124).

A segment address is page translated if its corresponding segment register has the highest order bit clear (zero); otherwise, the address is given to the I/O controller for translation. The PowerPC 601 RISC Microprocessor has a feature, called memory-forced I/O, that enables the I/O controller to generate a real address. This address can be used when I/O space is memory mapped.

I/O Controller Types

Each processor architecture expects an I/O controller to interface between the system bus (the one the processors use to access RAM) and an I/O bus. A computer system may have more than one I/O bus, but each bus has its own I/O controller.

The following table shows processor types, I/O controllers, and I/O bus protocol combinations that are supported by AIX Version 4.1:

Processor	Controller	Bus Protocol	Address Space for I/O
POWER RS1	IOCC	Micro Channel	I/O address
POWER RSC	IOCC	Micro Channel	I/O address
POWER RS2	XIO	Micro Channel	I/O address
PowerPC	ASIC	Micro Channel	I/O address
PowerPC	PCIB/MC	PCI	real address (memory mapped I/O)

The following list explains some terms used in the preceding table:

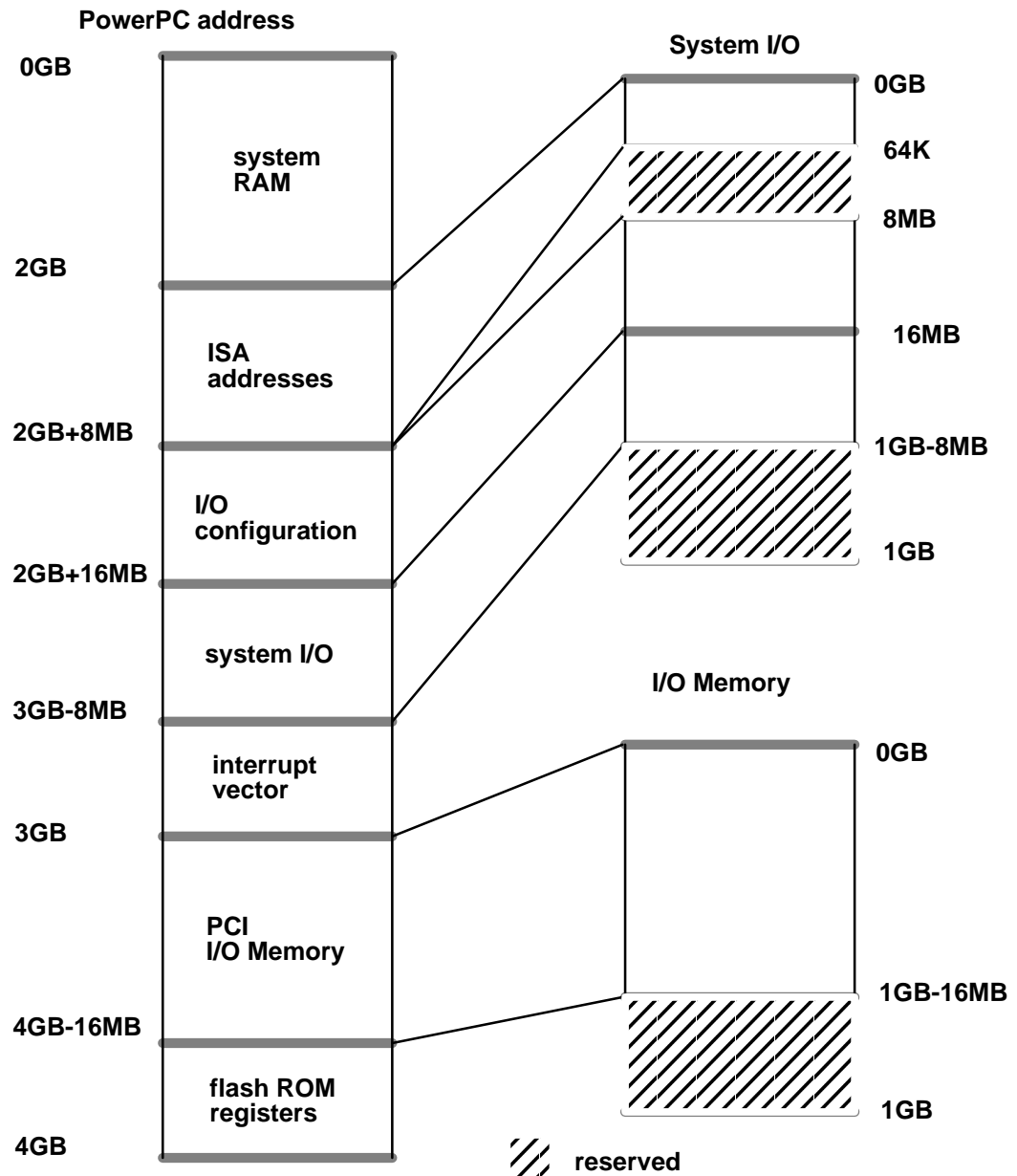
- Power RSC is POWER architecture on a single chip
- Power RS2 is POWER2 architected chip
- IOCC is I/O Channel Controller chip
- XIO is Extended IOCC
- ASIC is Application Specific Integrated Circuit chip
- PCIB/MC is Peripheral Component Interconnect (PCI) Bridge/Memory Controller complex.

The first four configurations, which interface to a Micro Channel bus, are quite similar. I/O space is accessed by passing an address to an I/O controller for translation; the address is then an *I/O address*. In this case, I/O space is mapped to I/O addresses by the I/O controller. Information about differing implementation details is contained in *POWERstation and POWERserver Hardware Technical Information-General Architectures*.

Configurations having a PCIB/MC, which interfaces to a Peripheral Component Interconnect (PCI) bus, are quite different. The I/O space is part of real address space, meaning that a memory controller translates real addresses so that certain address ranges access system RAM, and other address ranges access other system or bus attached devices. In this sense, I/O space is memory mapped. Different systems may have different address ranges mapped to different devices. For implementation details, see the hardware technical reference for the particular system.

I/O Space on PCI and ISA Systems

On PowerPC systems that do not have a Micro Channel bus, I/O space is memory mapped by the memory controller. So, I/O space is accessed by generating real addresses in one of two ways: either through block address translation, or memory-forced I/O (PowerPC 601 RISC Microprocessor only). The figure Controller's Memory Map (32 bit) shows an example of how the memory controller maps real addresses to bus addresses.



Controller's Memory Map (32 bit)

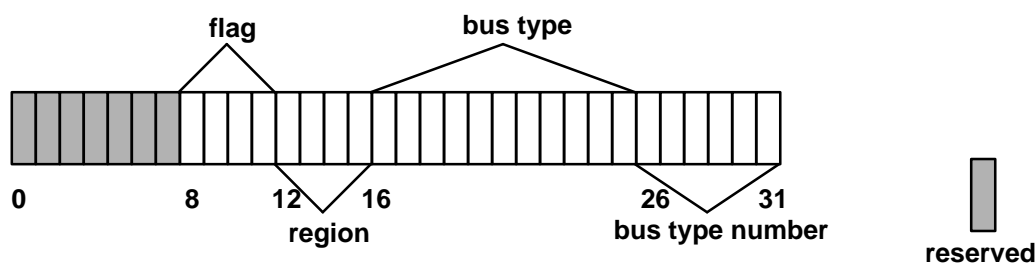
Programmed I/O to PCI and ISA Adapters

A routine is said to perform *programmed I/O* whenever it issues a load or store instruction with an address mapped to a bus or planar device. One distinguishes programmed I/O, where a system processor performs the data transfer, from direct memory access (DMA) where data is transferred by some other means.

For example, to output the value 0x12345678 to some register at offset 0x2f7:

```
volatile uchar *ioaddr;
struct io_map io_map;
...
ioaddr = iomem_att(&io_map);    /* do after io_map initialized */
*(ioaddr + 0x2f7) = 0x12345678;
eieio();                       /* ensures I/O instructions complete */
...
iomem_det(ioaddr);
```

The argument to **iomem_att** is a pointer to a structure **io_map** as defined in **sys/ioacc.h**. The calling routine provides the size of the address space needed, and a bus ID which specifies the bus type of region to be mapped. The figure Format of a Bus ID (real address, 32 bit) shows the format of a bus ID. The device driver's configuration entry point must get a valid bus ID from its configure method.



Format of a Bus ID (real address, 32 bit)

In this case, the bus type might be **IO_ISA**. The region might be **ISA_IOMEM** if I/O is to an ISA adapter's registers, or **ISA_BUSMEM** if I/O is to RAM on an ISA adapter.

A call to **iomem_att** returns a bus address in **io_map**. The bus address is to be passed to **iomem_det** after the I/O is complete.

There is no exception handling for programmed I/O on systems that do not have Micro Channel. An I/O exception causes the system to halt.

Direct Memory Access

A method of transferring data to or from a device without having a system processor issue load or store instructions is to use *direct memory access* (DMA), which relies on capabilities designed into the DMA controller and the adapter interfacing to the attached device.

There are two types of DMA:

- DMA Master

When an adapter arbitrates for the bus, and is able to transfer data directly by generating its own bus addresses and transfer lengths, then the transfer is *DMA master* and the card is a *DMA master adapter*.

- DMA Slave

When an adapter arbitrates for control for the bus, but lacks the ability to generate its own bus addresses to perform data transfer, so that a third party (a DMA controller) performs the data transfer, then the transfer is *DMA slave* and the card is a *DMA slave adapter*.

The steps for performing DMA differ between bus types, though the steps for performing DMA on ISA are similar to those for PCI. There is currently no DMA support for PCMCIA adapters.

DMA on POWER and POWER2 Architectures

The RIOS-1 and RIOS-2 architectures are not cache-consistent. Therefore, on these platforms, memory pages involved in a DMA transfer have to be made inaccessible to the processor (hidden), and the processor data cache must be flushed accordingly.

The IOCCs on these models are buffered (including the dual-buffered XIO) implementations, and thus require appropriate buffer flushing and invalidating.

On these platforms, TCE memory is disjoint from system memory and resides logically with the IOCC. Also, the TCE entries contain reference and change bits that must be maintained.

DMA slave operations are performed via the TAG Table implementation.

Authority checking is performed using a mask residing in the Channel Status Register for a DMA channel by checking that mask against the page protect key in the associated TCEs.

DMA on RSC (Single-Chip) Architectures

The RSC architecture is a cache-consistent architecture, therefore page hiding and data cache flushes aren't necessary.

The IOCC on this architecture is non-buffered, therefore does not require buffer flushes or buffer invalidation.

DMA on PowerPC Architectures

The PowerPC architecture is a cache-consistent architecture, therefore page hiding and data cache flushes aren't necessary. However, if a DMA operation modifies an instruction stream, instruction cache management is the responsibility of the application or instigator of the data transfer.

DMA Routines for PCI and ISA Adapters

The **dio** structure, which is defined in **sys/dma.h**, has many uses by a device driver. It is used both to pass a list of virtual addresses and lengths of buffers to the **d_map_list** and **d_map_slave** services, and to receive the resulting list of bus addresses (**d_map_list** only) for use by the device in the data transfer. Note that for calls to **d_map_slave** the driver does not need a **dio** bus list, since by nature the address generation for slaves is hidden.

Typically, a device driver will provide a **dio** structure containing only one virtual buffer and length in the list. Then, if the virtual buffer length spans many pages, the bus address list would contain multiple entries reflecting the physical locations making up the virtually contiguous buffer. The driver, however, can provide multiple virtual buffers in the virtual list which allows it to coalesce many buffer requests into one I/O operation. The device driver is responsible for allocating the storage for all **dio** lists that it will need. The driver will need at least two **dio** structures, one for passing in the virtual list, and another for accepting 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 function of the device and driver. The virtual list needs as many elements as are planned to be coalesced into one operation by the device. A formula for estimating how many elements the bus address list will need is the sum of each of the virtual buffers lengths divided by page size plus 2. One way to represent this formula is:

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

This is to handle a worst case situation, where, for a contiguous virtual buffer spanning multiple pages, each physical page is discontinuous, and neither the starting or ending addresses are page-aligned.

If the **d_map_list** service runs out of space while filling in the **dio** bus list, then an error, **DMA_DIOFULL**, is returned to the device driver and the **bytes_done** field of the **dio** virtual list is set to how many bytes were successfully mapped in the bus list. This byte count is a multiple of the **minxfer** field provided to the **d_map_list** or **d_map_slave** service. Also, the **resid_iov** field of the virtual list is set to the index of the first **d_iovec** entry representing the remainder of **iovecs** that could not be mapped. The device driver then can initiate a partial transfer on its device and leave the remainder on its device queue, or make another call to the **d_map_list** with new **dio** lists for the remainder, then setup its device for the full transfer that was originally intended. If the driver chooses not to initiate the partial transfer, it still must make a call to **d_unmap_list** to undo the partial mapping.

If **d_map_list** or **d_map_slave** encounter an access violation on a page within the virtual list, then an error, **DMA_NOACC**, 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**. Also 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 that in the **DMA_NOACC** case, the **bytes_done** count is *not* always a multiple of the **minxfer** field provided to the **d_map_list** or **d_map_slave** service, and no partial mapping is done. This means for slaves that no setup of the address generation hardware has been done, and 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.

Finally, if while mapping a transfer, either the **d_map_list** or **d_map_slave** services run out of resources to map the transfer, an error, **DMA_NORES**, 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 a multiple of the **minxfer** field provided to the **d_map_list** or **d_map_slave** service. 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 then can initiate a partial transfer on its device and leave the remainder on its device queue, or choose to leave the entire request on its device queue and wait for resources to free up (for example, after device interrupt from previous operation). 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.

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 modifies any of the other fields in the bus list. It is this list that the device driver uses to setup its device for the transfer, and it is this list that 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.

As far as the virtual list, the device driver sets up all of the fields except for the **bytes_done** and **resid_iov** fields which are set by the mapping service.

Notes:

1. More information describing the DMA services for PCI and ISA adapters is contained in the appendix.
2. These services are typically invoked by using macros (such as **D_MAP_INIT** and **D_MAP_LIST**) defined in **/sys/dma.h**.

Page Protection

Page protection checking is performed by the **d_map_page**, **d_map_list**, and **d_map_slave** services by calling the **xmемdma** kernel service for each page of a requested transfer. If the intended direction of a transfer is device to memory, then the page access permissions must allow writing to the page. If the intended direction of a transfer is from memory to device, then the page access permissions need only allow reading from the page.

In the case of a protection violation, a 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 allows a device driver to bypass the access checking functionality of these services; this flag should only be used for global system buffers such as mbufs or other command, control, and status buffers used by a device driver.

Peer-To-Peer DMA Support

Peer-to-Peer DMA is DMA transfer from one adapter's bus memory to another's, controlled by either a Bus Master or a DMA Slave with assistance of the I/O Controller. In either case, the transfer is independent of the CPU and does not involve system memory.

The flag **BUS_DMA** supports peer-to-peer DMA operations on calls to **d_map_page**, **d_map_list** and **d_map_slave**. This flag ensures that these services do not translate the provided address as a virtual address; instead it ensures that they treat the address as a bus address.

DMA Master I/O for an ISA Adapter

Because an ISA adapter can only generate 24 bit addresses, it can address up to 16MB. A device driver may need to ensure that DMA buffers are in low memory, unless there is no more than 16MB of RAM on the system. A driver allocates such buffers with the **rmalloc** kernel service, and later frees them with the **rmfree** kernel service. Then, the driver can “bounce” data from user supplied buffers (which may not be in low memory) into the buffer allocated by **rmalloc**. Then the driver could cause the data to be transferred to the ISA adapter with DMA.

For example, if the driver is to transfer data from a user buffer to a DMA master adapter, some steps for the transfer could be:

Steps for ISA DMA Master Transfer	
Module	What Device Driver Could Do
xyzopen	rmalloc driver buffer(s). Call DIO_INIT to allocate and initialize a dio structure. Call D_MAP_INIT which returns “d_handle” containing bus specific services. Call D_MAP_ENABLE if DMA not already enabled for this “d_handle”. Call D_MAP_PAGE if rmalloc 'ed buffer contained within a 4K page, or call D_MAP_LIST if rmalloc 'ed buffer(s) span multiple pages.
xyzread or xyzwrite	Perform PIO write to adapter register to start DMA transfer. Copy data to bounce buffer (writes only).
xyzintr	(Interrupt received from adapter when DMA transfer complete.) Copy data out of bounce buffer (reads only). Perform PIO read of adapter register to check DMA status.
xyzclose	Call D_UNMAP_PAGE or D_UNMAP_LIST. Call D_MAP_DISABLE if DMA not already disabled for this “d_handle”. Call D_MAP_CLEAR to free “d_handle” structure. Call DIO_FREE to free the dio structure, then rmfree the driver buffer(s).

Here is an example in pseudo-code of how a DMA master transfer could be done:

```
Initialization entry point: (xyzopen or xyzconfig)

    determine bus type for device from configuration information
    determine 64 vs.32 bit capabilities from configuration
        information
    call "handle = D_MAP_INIT(bid, DMA_MASTER|flags,
                            bus_flags, channel)"
    if handle == DMA_FAIL
        could not configure
    else
        call "D_MAP_ENABLE(handle)" (if necessary)

start_io entry point: (xyzread or xyzwrite)

    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, minxfer, vlist, 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

finish_io entry point: (xyzintr)

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

unconfigure code: (xyzclose or xyzconfig)

    call "D_MAP_DISABLE(handle)" (if necessary)
    call "D_MAP_CLEAR(handle)"
```

DMA Slave Transfers on an ISA Adapter

Here is an example in pseudo-code of how a DMA slave transfer could be done:

```
initialization entry point: (xyzopen or xyzconfig)

    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)

start_io entry point: (xyzread or xyzwrite)

    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

finish_io entry point: (xyzintr or xyzclose)

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

unconfigure entry point: (xyzclose or xyzconfig)

    call "D_MAP_DISABLE(handle)" (if necessary)
    call "D_MAP_CLEAR(handle)"
```

DMA Master Transfers on a PCI Adapter

The main difference between DMA on ISA and DMA on PCI is that there is no need for “bounce” buffers for PCI DMA; also, there are no DMA slave adapters on a PCI bus.

For example, if the driver is to transfer data from a user buffer to a DMA master adapter, some steps for the transfer could be:

Steps for PCI DMA Master Transfer (Short Term)	
Module	What Device Driver Could Do
xyzopen	Call DIO_INIT to allocate and initialize a dio structure. Call D_MAP_INIT which returns “d_handle” containing bus specific services. Call D_MAP_ENABLE if DMA not already enabled for this “d_handle”.
xyzread or xyzwrite	Call D_MAP_PAGE if buffer contained within a 4K page, or call D_MAP_LIST if buffer(s) span multiple pages. Perform PIO write to adapter register to start DMA transfer.
xyzintr	(Interrupt received from adapter when DMA transfer complete.) Perform PIO read of adapter register to check DMA status. Call D_UNMAP_PAGE or D_UNMAP_LIST.
xyzclose	Call D_MAP_DISABLE if DMA not already disabled for this “d_handle”. Call D_MAP_CLEAR to free “d_handle” structure. Call DIO_FREE to free the dio structure.

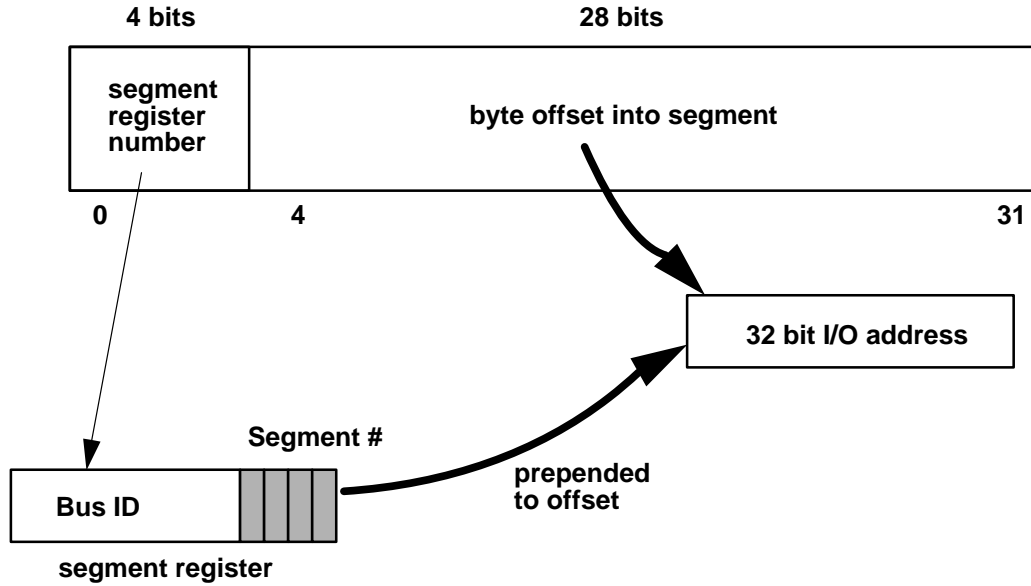
I/O Controller Interface Translation on Micro Channel Systems

A processor associates device registers or memory to particular addresses via *I/O controller interface translation*. When a processor instruction accesses an I/O bus, the address referenced by that instruction is said to *access an I/O controller interface* because all three processor architectures expect an I/O controller to interface between a processor and an I/O bus associated with that I/O controller.

I/O controller interface translation, like page address translation, is affected by the value of the segment ID contained in a segment register.

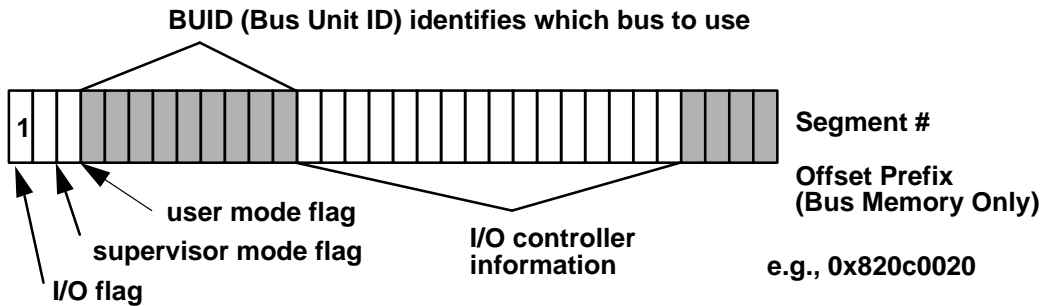
One performs I/O by executing load or store instructions referencing virtual addresses. Some other architectures (for example, the Intel x86 series) rely on special instructions to perform I/O to particular ports. The three processor architectures considered here translate addresses so that system device registers, adapter registers, adapter RAM, and so on, appear to be part of virtual memory.

The figure How an Address Specifies a Segment Register shows how the address specifies which segment register to use. For example, the address 0x50285740 says to look in segment register 5 for the segment ID, and it references byte 0x285740 within that segment. If the high order (leftmost) bit in segment register 5 is 0, then the address is translated as a page address, otherwise, the address is translated so that one can access an I/O bus through its I/O controller interface.



How an Address Specifies a Segment Register

A bus ID is a segment identifier associated with an I/O segment. The Format of a Bus ID figure shows the overall format of a bus ID contained in a segment register. Note that the high order bit is set to 1. This bit distinguishes which kind of address translation to perform: when set, it signifies that addresses in this segment require I/O controller interface translation.

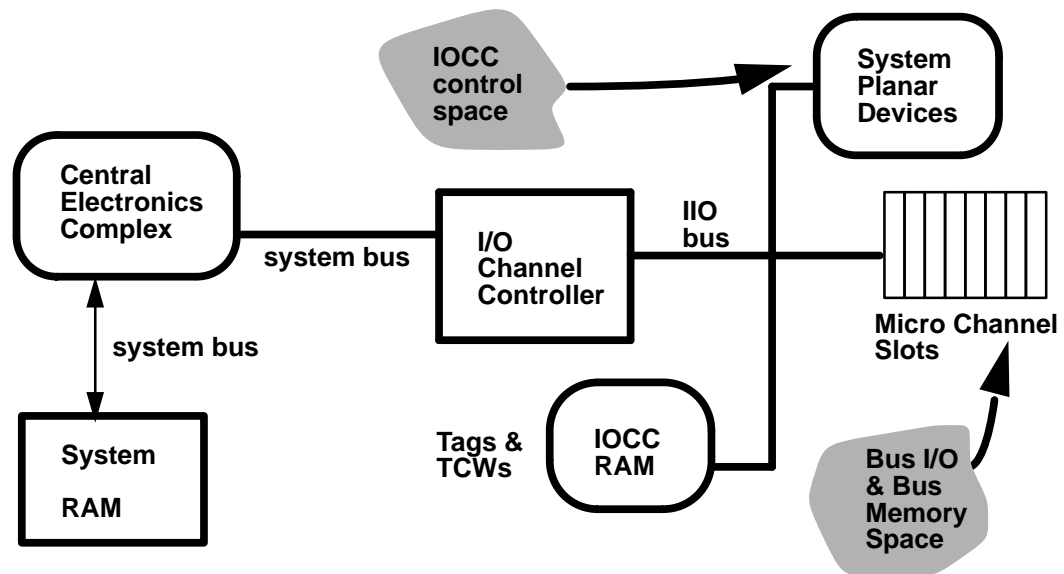


Format of a Bus ID

I/O controller interface translation differs from page address translation in that no page tables are searched. So there is no way to protect one program accessing an I/O bus from another program accessing the same bus.

I/O Address Spaces on Micro Channel Systems

The I/O Subsystem figure shows how portions of the I/O subsystem work with one another. Note that a system processor does not access an I/O bus directly, but relies on a special-purpose processor, mounted on the system planar with its own dedicated RAM, to serve as an interface between the system bus, which accesses system RAM, and the I/O bus.



The I/O Subsystem

Note that there are devices on the system planar itself, and there are devices plugged into slots accessing the Micro Channel bus. There are also resources on the IOCC itself that one may wish to modify or inspect. There is a distinction between I/O to a system planar device (or the IOCC itself) and I/O to a device attached to the Micro Channel bus. You specify which kind of I/O to do by selecting an I/O address segment, or “I/O space.”

There are several kinds of I/O address spaces defined by each processor architecture. The ones that affect a device driver are:

- IOCC control space (or system I/O space)
- Bus I/O space
- Bus memory space

Devices that are attached to the system planar board (also known as the *motherboard*) have addresses mapped to IOCC control space. Examples of such devices are: system timer, calendar, non-volatile RAM, the LED registers, programmable option select (POS) registers, and other planar device registers.

IOCC control space consists of one I/O segment called the IOCC control segment. An I/O address translates to this segment if the IOCC bit in the IOCC controller information part of a bus ID is set to 1. See the IOCC Control Space (Portion) figure for an example (for the POWER architecture) of what system planar devices are mapped to which addresses.

Both bus I/O space and bus memory space are collectively referred to as *standard I/O space*, and are accessed within one segment.

Any attempt to access IOCC control space without system privilege (as marked in the bus ID) causes an I/O exception.

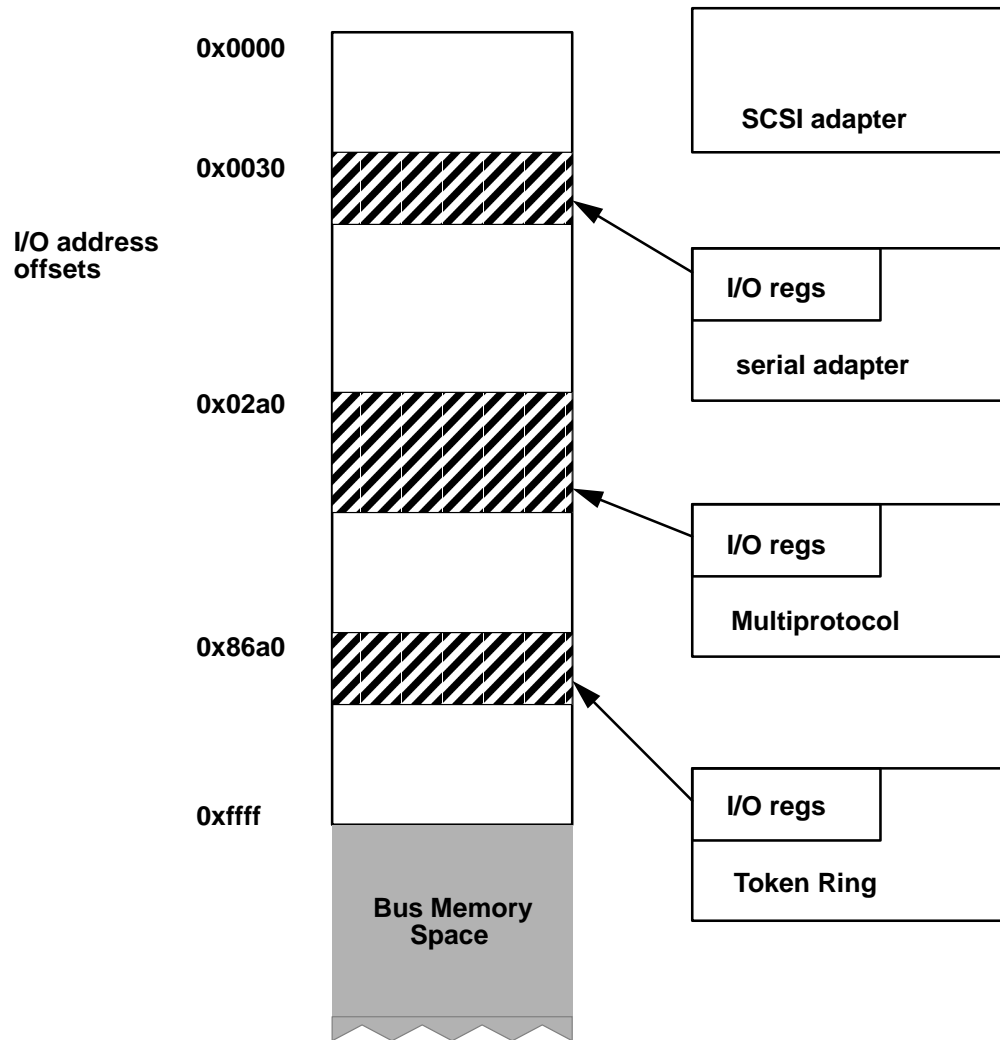
0x400020	bus status register	4 bytes
0x400024	TCW Itag anchor addr register	4 bytes
0x40002c	component reset register	4 bytes
0x40002f	standard I/O reset register	1 byte
0x400084	interrupt request register (IRQ)	4 bytes
0x4000c0	time of day clock	32 bytes
0x4000e0	system reset status register	1 byte
0x4000e4	power status register	4 bytes
0x4000ea	power reset register	2 bytes
0x4000ec	interrupt request register	1 byte
0x4000fc	I/O board EC level register	1 byte
0x400000	POS registers slot 0	
0x410000	POS registers slot 1	
0x4f0000	POS registers slot 15	
0xa00000	operator panel LEDs (2 bytes)	NVRAM, 2 MB
0xa00300		

IOCC Control Space (Portion)

Note: Some of the information in the preceding IOCC Control Space (Portion) figure applies only to the POWER architecture and not to the PowerPC architecture.

So, for example, the address 0xf0400020 refers to BUID in register 15, where the offset specifies the bus status register. Note that IOCC RAM is mapped to IOCC control space.

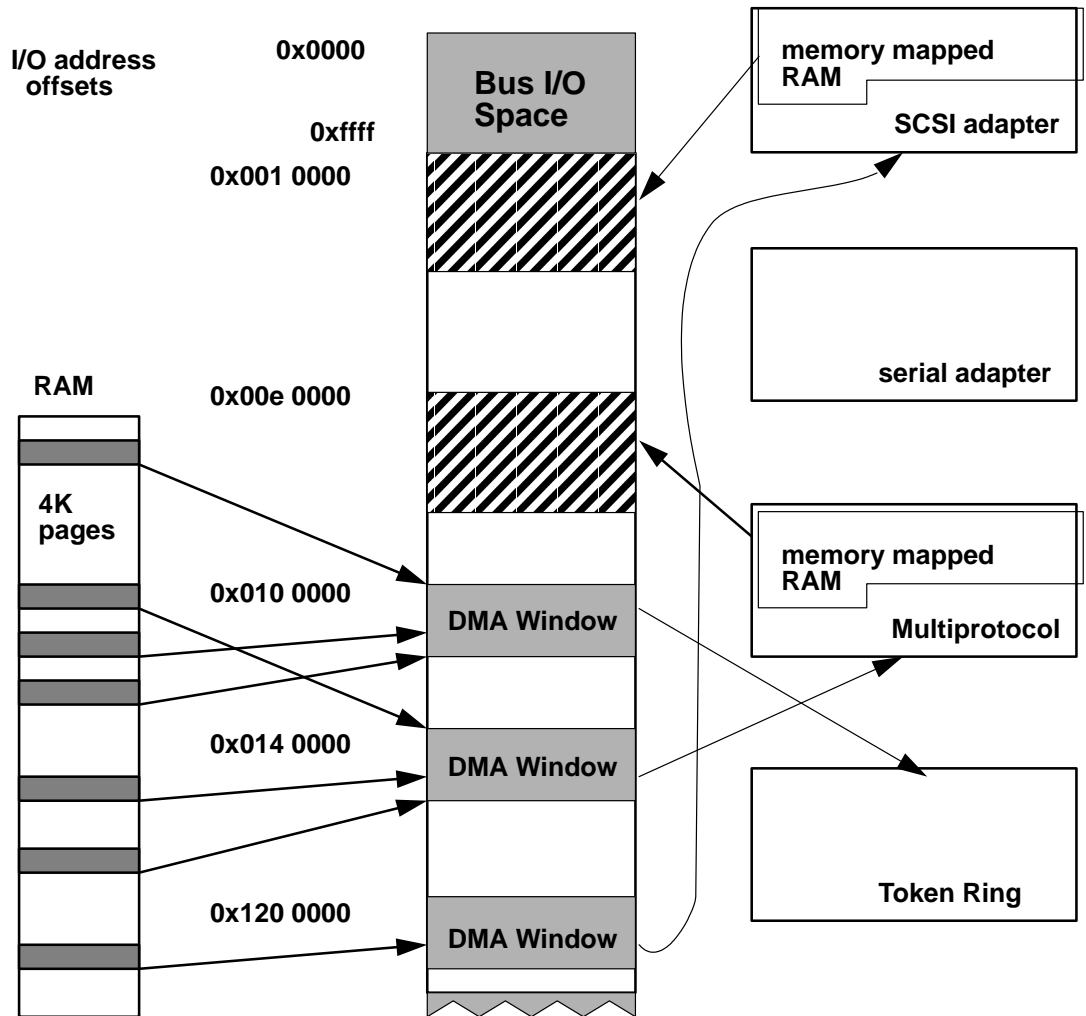
The register addresses start from the value of IO_IOCC as found in the header file `/usr/include/sys/iocc.h`. Here the value of IO_IOCC is 0x00400000 (4 MB).



Bus I/O Space (Example)

The Bus I/O Space (Example) figure shows an example of how adapters attached to the I/O bus can be mapped to bus I/O space. The IOCC checks each access to bus I/O space to see if the address is within range of some region mapped in bus I/O space. If not, it generates an I/O exception.

Note that the SCSI adapter has no registers mapped to addresses in bus I/O space, but the others do. For example, a load instruction from address 0x50000030 would (if segment register 5 has a bus ID for the standard I/O segment) get four bytes from the serial adapter's first register. Provided that segment register 15 has a standard I/O segment's bus ID, a store instruction to address 0xf00086a0 would place four bytes of data in the first register on the Token-Ring adapter.



Bus Memory Space (Example)

Each 4K page of bus memory space that is mapped for translation is associated with a Translation Control Word (TCW), which is a data structure kept in RAM dedicated for use by an IOCC. Any attempt to read or write to an address that is not within a page associated with a TCW causes an I/O exception. For more information on TCWs and protection, see "Translation, Protection, and the TCW Table," in "Programming Model: System I/O Structure" in *POWERstation and POWERserver Hardware Technical Information-General Architectures*.

The Bus Memory Space (Example) figure shows an example of how the same four adapters have bus memory addresses mapped to RAM either on the adapters or on the system itself. Note that in this case, the serial adapter has no addresses in bus memory space for its use.

In this example, a load instruction from address 0xf0010000 would read the first four bytes in RAM on the SCSI adapter. Similarly, a store instruction to address 0xf00e0000 would write four bytes to RAM on the multiprotocol adapter.

Once the mapping is in place, transfer from system memory to adapter RAM amounts to issuing a sequence of load and store operations from one region of bus memory space to another.

Programmed I/O to Micro Channel Adapters

Programmed I/O requires a properly formatted bus identifier in a segment register. The device driver's configuration method gets a bus identifier for its particular bus by reading its parent `bus_id` attribute out of the predefined attribute (PdAt) object class of the ODM. The configuration method passes the bus identifier, within the device dependent structure (DDS), to the device driver's configuration entry point where the bus identifier is further modified as needed. For example, a driver may wish to set the bus ID bits associated with *privilege key* and *IOCC space select* by the following statement:

```
bid |= (IOCCSR_KEY | IOCCSR_SELECT); /* see <sys/iocc.h> */
```

To perform I/O to IOCC space, one uses the following macros:

```
IOCC_ATT(bid, iocc_addr) /* to place Bus ID in a segment reg */
BUSIO_PUTLX(long_val, iocc_addr) /* to write four bytes */
BUSIO_GETLX(long_val) /* to read four bytes */
BUSIO_PUTSX(short_val, iocc_addr) /* to write two bytes */
BUSIO_PUTCX(char_val, iocc_addr) /* to write one byte */
IOCC_DET(iocc_addr)
```

To perform I/O to bus I/O space (adapter registers):

```
BUSIO_ATT(bid, io_addr) /* to place Bus ID in a segment reg */
BUSIO_PUTLX(long_val, io_addr) /* to write four bytes */
BUSIO_GETLX(long_val) /* to read four bytes */
BUSIO_PUTSX(short_val, io_addr) /* to write two bytes */
BUSIO_PUTCX(char_val, io_addr) /* to write one byte */
BUSIO_DET(io_addr)
```

To perform I/O to bus memory space (adapter RAM):

```
BUSMEM_ATT(bid, mem_addr)
BUS_PUTLX(long_val, io_addr) /* to write four bytes */
BUS_GETLX(long_val) /* to read four bytes */
BUS_PUTSX(short_val, io_addr) /* to write two bytes */
BUS_PUTCX(char_val, io_addr) /* to write one byte */
BUSMEM_DET(io_addr)
```

The macros suffixed with an 'X' have exception handlers built into them. The return value is nonzero if an I/O exception occurs during the data transfer. I/O to IOCC control space does not generate an I/O exception if there is an error. To verify that data has been written correctly, you need to read it after it is written.

It is possible to use macros that lack exception handlers, but if you don't provide your own exception handler, then an I/O exception would cause the system's default handler to halt the system. There are other macros defined in the header file `/usr/include/sys/ioacc.h`. For more information on this, see "I/O Bus Attach/Detach Macros for Programmed I/O (PIO) Operations," in the chapter about device drivers in *AIX Version 4.1 Kernel Extensions and Device Support Programming Concepts*

DMA Master I/O on a Micro Channel Adapter

As an example, if the driver is to transfer data from a user buffer to a DMA master adapter, some steps for the transfer could be:

Steps for DMA Master I/O for Micro Channel	
Module	What Device Driver Could Do
xyzopen	Call d_init to allocate and initialize DMA channel.
xyzread or xyzwrite	Call d_master to map pages of I/O space or RAM to TCWs. Perform PIO write to adapter register to start DMA transfer. Check for PIO exception and report any error.
xyzintr	(Interrupt received from adapter when DMA transfer complete.) Perform PIO read of adapter register to check DMA status. Check for PIO exception and report any error. Call d_complete to unhide pages and check for errors detected by IOCC.
xyzclose	Call d_clear to free DMA channel.

DMA Slave I/O on a Micro Channel Adapter

As an example, if the driver is to transfer data from a user buffer to a DMA slave adapter, some steps for the transfer could be:

Steps for DMA Slave I/O for Micro Channel	
Module	What Device Driver Could Do
xyzopen	Call d_init to allocate and initialize DMA channel.
xyzread or xyzwrite	Call d_slave to map pages of I/O space or RAM to Tags. Perform PIO write to adapter register to start DMA transfer. Check for PIO exception and report any error.
xyzintr	(Interrupt received from adapter when DMA transfer complete.) Call d_complete to unhide pages and check for errors detected by IOCC.
xyzclose	Call d_clear to free DMA channel.

Alignment Issues for DMA on Micro Channel

As described under “Hiding DMA Data” in “Understanding DMA Transfers” in the chapter on kernel services in *AIX Version 4.1 Kernel Extensions and Device Support Programming Concepts*, pages covering buffers for use by DMA services may be hidden unless `DMA_NOHIDE` is specified. If the buffer is supplied by a user, and if the buffer is not page aligned, then there is a chance that nearby data could be on the same page as the buffer being hidden. If a user routine attempts access to such data, its context will sleep until the page is unhidden. If the driver does not unhide the page, then such user routines would sleep indefinitely.

Any buffers that a driver allocates for DMA transfers may preferably be allocated from the kernel heap (via `xmalloc`) so that the buffer can be page aligned: this lessens the chance of nearby data being inadvertently hidden, and simplifies the coordination of data transfer with the adapter. Also, fewer pages may be used.

If a DMA buffer is to be shared between the system processor (the driver) and a DMA master adapter, then it would be registered with the flag `DMA_NOHIDE`. If a DMA buffer is to be so shared, or shared between two DMA master adapters (where page hiding is irrelevant), then there is a chance of data corruption due to race conditions (see Chapter 5, Synchronization). An I/O controller reads and writes data in chunks which are the size of its data cache line. If a DMA transfer starts at an address that is not cache-aligned, then nearby data is inadvertently accessed. A race condition can be avoided by synchronizing the access, or by ensuring that DMA accesses are cache-aligned and thereby non-overlapping. A buffer can be cache aligned with the kernel services `d_align` or `d_roundup`.

Here is an example of a call-side routine that sets up DMA master transfers on a Micro Channel adapter. Please note that the buffers are allocated by the driver itself and they are accessed by the system and the adapter, so they are not hidden (DMA_NOHIDE) from the system.

```
int dma_init()
{
    caddr_t busptr; /* points to bus address space */
    int rc; /* followed by other necessary declarations */

    scb.ccbp = (t_ccb *) xmalloc(PAGESIZE, PGSHIFT, pinned_heap);
    if (scb.ccbp == NULL)
    {
        DTRC(0x20, 0x43, (ulong) scb.ccbp); /* custom trace macro */
        return(ENOMEM);
    }

    /* clear the buffer for use */
    bzero(scb.ccbp, PAGESIZE);

    /* SCB - System Control Block /* you would define your own */
    * TSB - Termination Status Block
    * Separate regions must be set up for transferring data for DMA
    * reads and DMA writes. These separate regions should each consist
    * of an integral number of cache lines. The Command Control
    * Block (CCB) starts on the page boundary, which is also a cache
    * line boundary. The CCB and TSB must reside in separate 128-byte
    * regions for cache consistency.
    */
    scb.ccb_size = d_roundup(sizeof(t_ccb)); /* size of ccb struct */

    /*
    * The algorithm used by d_roundup() is :
    * outlength = (inlength + cache_line_size-1) & (cache_line_size-1);
    * where inlength is the length in bytes of the structure to be mapped.
    *
    * So the next available cache line boundary is at ccbp + ccb_size
    */
    scb.tsbp = (uint)scb.ccbp + (uint)scb.ccb_size;

    /* need the size of tsb also for flushing */
    scb.tsb_size = d_roundup(sizeof(t_tsb)); /* size of tsb struct*/

    /* the amount of data that is actually used is:(should be 256 bytes)*/
    scbctl_size = (uint)scb.tsb_size + (uint)scb.ccb_size;

    /* The ccb is mapped to the beginning of the bus memory window
    * configured for this device. The tsb follows.
    */
    scb.bus_ccbp = dds.dds_hdw.dma_bus_mem; /* device dep. struct */
    scb.bus_tsbp = (uint)scb.bus_ccbp + (uint)scb.ccb_size;

    /* initialize the ccb before mapping */
    .....
    /* set up cross memory descriptor for the ccb and tsb */
    scbctl_xmem.space_id = XMEM_INVAL;
    rc = xmattach((char *)scb.ccbp, scbctl_size, &scbctl_xmem, SYS_ADSPACE);
}
```

```

if (rc == XMEM_FAIL)
{
    DTRC(0x20, 0x43, (ulong) rc);
    xfree(scb.ccbp, pinned_heap);
    return(ENOMEM);
}

/* map the control area to the first page of bus dma memory space */
d_master(scb.chan_id, DMA_READ|DMA_NOHIDE, (caddr_t)scb.ccbp,
        scb.ctl_size, &scb.ctl_xmem, scb.bus_ccbp);

/* attach to the bus segment */
busptr = BUSIO_ATT(dds.dds_hdw.bus_id, (caddr_t)dds.dds_hdw.bus_io_addr);

/* write address of the ccb to the command port, it never changes */
port.command = (caddr_t) scb.bus_ccbp;
/* pio_retry() logs temporary failure & after 3rd fail gives up */
if(BUS_PUTLRX((busptr + HLZ1_COMMAND), port.command))
    rc = pio_retry(rc, PUTL, addr, port.command);

if (port.pio_fatal) /* check for a permanent PIO exception */
{
    DTRC(0x20, 0x45, (ulong) rc);
    xmdetach(&scb.ctl_xmem);
    xfree(scb.ccbp, pinned_heap);
    BUSIO_DET(busptr);
    return(EIO); /* must read this reg to continue */
}

/* set the dma busmem address for the input buffer mapping */
dma.indmap = dds.dds_hdw.dma_bus_mem + PAGESIZE;
dma.window_size = dds.dds_hdw.dma_bus_length - PAGESIZE;
dma.enddmap = dds.dds_hdw.dma_bus_mem + dds.dds_hdw.dma_bus_length;

/* initialize mapping control variables */
dma.last_inbuf = NULL;
dma.last_outbuf = NULL;

/* enable the device as a busmaster, enable interrupts, clear IV on read*/
port.sub_ctrl = .... /* device flags */;
if(rc = BUS_PUTCX((busptr + HLZ1_SUB_CTRL), port.sub_ctrl))
    rc = pio_retry(rc, PUTC, addr, port.sub_ctrl);

if (port.pio_fatal) /* check for a permanent PIO exception */
{
    DTRC(0x20, 0x46, (ulong) rc);
    xmdetach(&scb.ctl_xmem);
    xfree(scb.ccbp, pinned_heap);
    rc = EIO; /* must read this reg to continue */
}
BUSIO_DET(busptr);
return(rc); /* error already logged if PIO exception */
} /* end dma_init() */

```

Chapter 3. Interrupts

This chapter describes the issues associated with writing an interrupt handler for an adapter on AIX Version 4.1. It contains the following sections:

- Overview, on page 3-1
- Interrupt Hardware Support, on page 3-2
- Interrupt Levels, on page 3-2
- Interrupt Priorities, on page 3-4
- Interrupt Level Mapping, on page 3-6
- Interrupt Handling, on page 3-9
- Interrupt Management Kernel Services, on page 3-9
- Multiprocessor Interrupt Concerns, on page 3-10

Overview

Adapters on any bus can generate interrupts to the host processor. Each interrupt is associated with a particular level, sometimes called an Interrupt Request Level (IRQ). For example, one adapter can generate interrupts on IRQ2, although another can generate interrupts on IRQ3. Assignment of interrupt levels to adapters is done by one of the following methods:

- Setting POS registers on the adapter during device initialization. Some buses do not support this.
- Manually setting the interrupt level via hardware switches.

Interrupt levels on certain buses, like the Micro Channel or PCI, can be shared. This means that more than one adapter can generate interrupts at the same level. In this case, each adapter must provide a register that can be interrogated to determine if the current interrupt is due to this particular adapter. This enables the software to determine which adapter actually generated the interrupt. Having each adapter on a separate interrupt level usually gives better performance than interrupt-sharing on a particular level.

Each adapter also provides a procedure, which software must follow, that resets an interrupt indication. This must be done before any more interrupts are generated by the adapter.

Device drivers usually provide an interrupt handler to handle these situations. An interrupt handler is a function that is called by the AIX kernel whenever an interrupt occurs on a given IRQ level. The interrupt handler must first determine whether the interrupt was caused by the adapter this driver is managing. If not, the handler exits immediately indicating no action taken. If so, the interrupt handler performs whatever processing is needed to deal with the interrupt, resets the interrupt in the adapter, and returns to the AIX kernel.

Interrupt handlers, like most of the AIX kernel, can be preempted. They run with some interrupts enabled. When a device driver configures itself, it specifies the priority of interrupts from its associated adapter. When an interrupt occurs, only interrupts from devices at that priority level and below are disabled. Higher priority interrupts can still occur. Interrupt handlers for low priority devices (such as, a printer) can be preempted if an interrupt occurs on a high priority device (such as, an unbuffered high-speed communications link).

The AIX operating system provides routines so that interrupts at higher levels can be disabled and enabled during the operation of the interrupt handler. Use these routines with caution so that the higher priority interrupts can be serviced.

Interrupt Hardware Support

The processor recognizes the following types of hardware interrupts:

- SYSTEM RESET
- MACHINE CHECK
- DATA STORAGE
- INSTRUCTION STORAGE
- EXTERNAL
- ALIGNMENT
- PROGRAM
- FLOATING-POINT UNAVAILABLE
- DECREMETER (available on PowerPC only)
- SYSTEM CALL
- TRACE
- FLOATING-POINT ASSIST

The hardware interrupts in the preceding list are the only way that normal instruction execution can be interrupted. Thus all interrupt signals must be converted to one of these to interrupt the processor. An interrupt vector for each type of hardware interrupt is loaded in low memory and is executed upon receipt of the interrupt. The first two types of interrupts, SYSTEM RESET and MACHINE CHECK, may occur at any time regardless of the state of the interrupt mechanism and are processed immediately. The other interrupts are handled in the order listed from highest to lowest priority.

The hardware interrupt type EXTERNAL is the interrupt into which all bus level interrupts are multiplexed. This is the type of interrupt that a device driver will be required to handle.

For the PowerPC architecture, the Machine State Register (MSR) and the Save Registers (SRR0 and SRR1) are used to process external interrupts. The MSR register contains an enable bit which is used to enable or disable external interrupts. The SRR registers are used to save the processor state at the time of the interrupt. For more details on these registers and the other types of hardware interrupts, see *PowerPC Architecture* (SR28–5124).

The hardware interrupt mechanism between the processor and the device is bus architecture dependent. For Micro Channel architecture, the interrupt controller is embedded in the I/O Channel Controller (IOCC). For PowerPC based platforms with PCI and ISA buses, a System I/O Bridge chip (SIO) and a proprietary PCI Bridge/Controller are used to relay the interrupt signal to the processor as an external interrupt.

Interrupt Levels

The mechanism by which the processor knows where the interrupt originated depends on a software abstraction of the interrupt source referred to as *processor interrupt level* (*Processor interrupt level* is also called *software interrupt level* or just *interrupt level*). The interrupt level corresponds to a distinct hardware interrupt level present on a device. For example, each bus interrupt level that originates on the PCI bus is converted into a unique software interrupt level that the processor understands and uses to call the proper interrupt handler for that interrupt. This design isolates the software from the information the hardware presents to indicate the interrupt.

The number of interrupt levels that a processor recognizes is architecture dependent. The POWER machines (RS1, RSC, and RS2) support 64 interrupt levels, and the PowerPC machines can support up to 8000 interrupt levels. An interrupt must map into one of these

levels to be correctly serviced; that is, each hardware interrupt level must be assigned to one of the processor interrupt levels. The details of this mapping process is described in "Interrupt Level Mapping", on page 3-6.

Because individual adapters interrupt at various hardware levels of their own depending on what type of bus they are on, there is one more configuration issue to consider. Some bus implementations support assigning hardware interrupt levels at system configuration time. System configuration software determines which adapters are present and assigns an interrupt level to the device adapter using special bus commands. System configuration then sets the device configuration hardware levels and initialization data to reflect this assignment.

For example, on the Micro Channel bus, there are 16 bus interrupt levels. A typical Micro Channel adapter has several bus interrupt levels from which to choose. Before setting which hardware level to use, the adapter's configuration method calls the **busresolve** routine. The **busresolve** routine returns a bus interrupt level that does not conflict with any other adapters on the bus. The bus interrupt level returned by **busresolve** is written to the ODM database. The device driver open routine can refer to the ODM database to register the interrupt handler at the correct level.

However, some buses, like the ISA bus, do not support programmable assignment of hardware interrupt levels. The assignment of these hardware interrupt levels is usually hardwired or selected by a jumper on the adapter. In this case, the Predefined Attributes PdAt entries of the ODM database must be set properly by the user to reflect the manual settings (or the manual settings changed to match the ODM database), because there is no way for the values to be changed by software. Driver developers should remember this important difference about ISA bus configuration.

Warning: Do not change the Predefined Attributes (PdAt) object class by removing information that was shipped with the base operating system.

On some buses, such as the PCI and MCA buses, interrupt level sharing is supported. To ensure that **busresolve** handles these levels correctly, the ODM stanza for interrupt level must have the type field set to "I". Non-shareable interrupts (all those on the ISA bus, for example) are declared with type set to "N". Shared interrupt levels on the MCA buses *must* all request the same priority, but this is not required for the PCI bus.

The following sample ODM entry declares an interrupt level to be shareable by setting the attribute type to "I":

```
PdAt :
    uniquetype = "adapter/pci/sample"
    attribute = "intr_level"
    deflt = "15"
    values = "15"
    width = ""
    type = "I"
    generic = "D"
    rep = "nl"
    nls_index = 3
```

Here is an example of a non-shareable interrupt (type = "N"):

```

PdAt :
      uniquetype = "adapter/isa/sample"
      attribute = "intr_level"
      deflt = "9"
      values = "5, 7, 9"
      width = ""
      type = "N"
      generic = "DU"
      rep = "nl"
      nls_index = 4

```

Chapter 6, "Device Configuration Methods," has more information on the ODM database and adapter device attributes.

The following table summarizes the interrupt level properties for different types of buses.

Bus	Interrupt Sensitivity	Interrupt Levels Shareable?	Programmable Interrupt Levels?	Available Interrupt Levels
MCA	Level	Yes, but must have same priority.	Yes	1 – 16
PCI	Level	Yes	Yes	15
ISA	Edge	No	No	5,7,9,11,14,15 *

*IRQ 15 is used by PCI or ISA, but not both.

In some documentation, the distinction between hardware interrupt levels and processor interrupt levels is not made, because there is a one-to-one mapping between the two. Usually, just the term interrupt level is used.

For information on how interrupts are processed, see "Processing Interrupts" in *AIX Version 4.1 Kernel Extensions and Device Support Programming Concepts*

Interrupt Priorities

Each processor interrupt level (see "Interrupt Levels", on page 3-2) also has a interrupt priority associated with it. The following is a list of all the interrupt priorities (sorted so the interrupt priorities with higher priority appear first):

- INTMAX
- INTCLASS0
- INTCLASS1
- INTCLASS2
- INTCLASS3
- INTOFFL0
- INTOFFL1
- INTOFFL2
- INTOFFL3
- INTIODONE
- INTBASE

These interrupt priorities represent the order in which the kernel enables external interrupts to be processed. (The **sys/m_intr.h** file shows the numeric value associated with each priority. The higher priorities have lower numeric values.)

Although the physical mechanisms vary for the different architectures, a basic rule for all architectures is that the interrupt priority must be of a higher priority than the current processor priority for the interrupt to be serviced. For example, if the processor is currently

servicing an interrupt of priority INTCLASS3, any off-level interrupts will not be serviced until the INTCLASS3 interrupt handling is finished, and the processor priority is lowered.

The choice of what priority to run your device driver interrupts is based on two criteria:

- The maximum interrupt latency requirements.
- The interrupt execution time of the device driver

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 degraded. The interrupt execution time is the number of machine cycles required by the device driver to service the interrupt.

A device with a short interrupt latency time must have a short interrupt service time. In other words, a device that *loses* data if not serviced quickly must have a higher priority interrupt level. This in turn requires that it spends less time in the interrupt handler. The following list contains general guidelines for interrupt service times:

Interrupt Priority	Service Time
INTMAX	All interrupts disabled
INTCLASS0	Less than 200 cycles
INTCLASS1	Greater than 200 but less than 400 cycles
INTCLASS2	Greater than 400 but less than 600 cycles
INTCLASS3	Greater than 600 but less than 800 cycles
INTOFFL0	Less than 1500 cycles (off-level priority)
INTOFFL1	Greater than 1500 but less than 2500 cycles (off-level priority)
INTOFFL2	Greater than 2500 but less than 5000 cycles (off-level priority)
INTOFFL3	5000 cycles or greater (off-level priority)
INTIODONE	I/O completion processing (lowest off-level priority)
INTBASE	All interrupts enabled

Note: To find out your interrupt service time, you can put a trace hook with a time stamp into both the entry and the exit points of your interrupt handler. The resulting trace tells you the cumulative time that was spent in the handler. Divide this time by the cycle speed of your system to get the number of cycles used.

Once an interrupt priority has been determined, the information is written to the system via an ODM entry. The following sample entry assigns an interrupt priority INTCLASS2, which has the value 3, to a device:

```
uniquetype = "adapter/mca/hscsi"  
attribute = "intr_priority"  
deflt = "3"  
values = "3"  
width = ""  
type = "P"  
generic = "D"  
rep = "nl"  
nls_index = 5
```

Chapter 6, "Device Configuration Methods," has more information on the ODM database.

Interrupt Level Mapping

The mapping of interrupt source to a specific processor interrupt level is dependent upon the hardware architecture. This section briefly describes the mappings for the POWER, POWER2, and PowerPC architectures and, also, discusses software interrupt priority.

The POWER architecture has a simple static mapping that directs the interrupt sources to one of the 64 available processor levels. Since there are at most 2 IOCCs, each bus has its 16 hardware levels directed to 32 of the existing processor levels and the rest are used for other interrupt sources. The POWER Interrupt Level Mapping figure shows the details of the processor level mapping scheme.

IOCC 0				E P O W	S G A	S L A	IOCC 1				E X T	D E C		O F F
0	15	24	25	28	32	47	48	62	63					

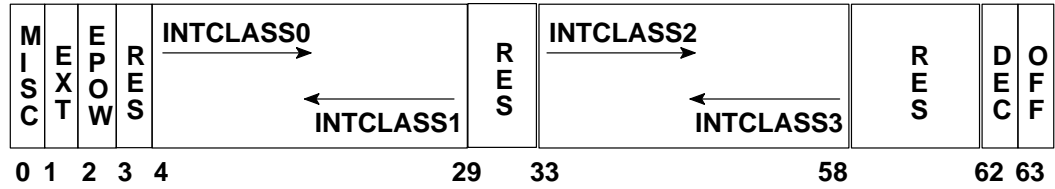
<u>LEVEL</u>	<u>ASSIGNED</u>
0–15	IOCC 0
24	Early Power-Off Warning (EPOW)
25	SGA – integrated graphics adapter
28	SLA0
29	SLA1
30	SLA2
31	SLA3
32–47	IOCC 1
48	External Check
62	Decrementer
63	Off-level hardware assist

POWER Interrupt Level Mapping

The interrupts are grouped into priorities by using a mask built and maintained by the interrupt subsystem called the External Interrupt Mask (EIM). This mask is changed whenever a new priority is set, thereby enabling the correct processor levels.

Although the POWER2 architecture also has 64 processor levels, it does not have the priority flexibility of the POWER architecture. There is a direct relationship of processor level and priority. Of the 64 processor levels, the hardware presents interrupts from the most favored (level 0) to the least favored (level 63). Thus, interrupt levels must be dynamically assigned based on priority.

The interrupt level allocation algorithm groups the processor levels into ranges corresponding to software priorities, and assigns them as needed. This ensures that higher priority interrupts are serviced before those of lesser priority. In addition to these dynamically allocated levels, there are other interrupt levels that are pre-allocated to specific sources. These can be seen in the POWER2 Interrupt Level Mapping figure.

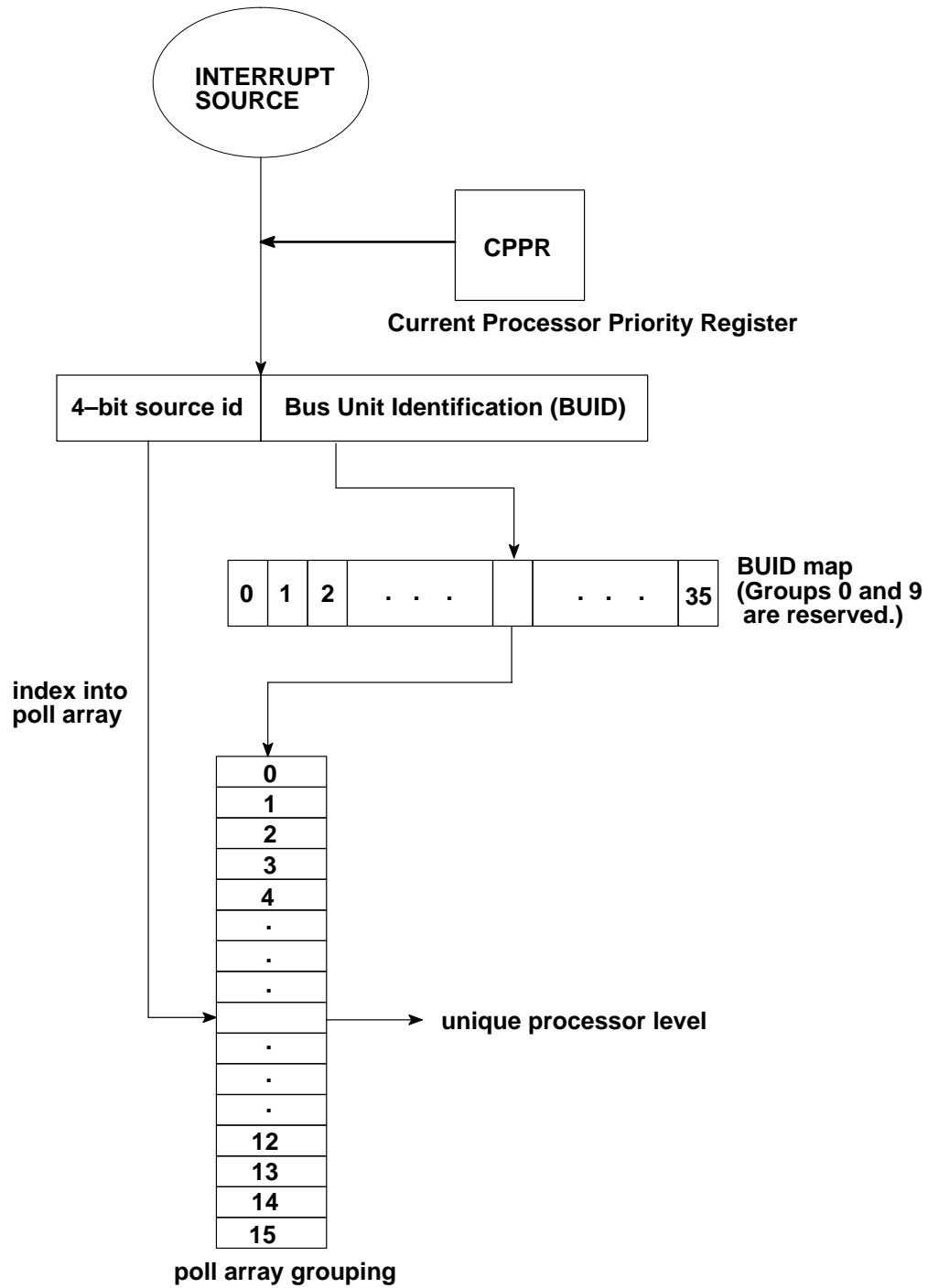


<u>LEVEL</u>	<u>ASSIGNED</u>
0	Miscellaneous IOCC interrupts
1	External Check
2	Early Power-Off Warning (EPOW)
3	Reserved
4–28	INTCLASS0, INTCLASS1 dynamically assigned interrupts
29–32	Reserved
33–57	INTCLASS2, INTCLASS3 dynamically assigned interrupts
58–61	Reserved
62	Decrementer
63	Off-level hardware assist

POWER2 Interrupt Level Mapping

The PowerPC interrupt logic is very different from the previous architectures. Each hardware level interrupt has a priority associated with it and will only interrupt on processors if this priority is more favored than the current processor priority. The processor priority is set by software for each processor.

Attached to each interrupt source is the Bus Unit Identifier (BUID) and a 4-bit code that indicates the source on the specific Bus Unit Controller (BUC). These two values are used by the interrupt subsystem to determine the interrupt level. On the PowerPC architecture, there are a possible 8000 levels, of which 160 are currently being used. These levels are assigned to groups that contain 16 interrupt levels. The poll groups are dynamically assigned to a BUID through the `buid_map` when the interrupt handler is registered. This mapping is shown in the PowerPC Interrupt Level Mapping figure.



PowerPC Level Mapping

Interrupt Handling

When an external interrupt is first detected, the system immediately calls the external interrupt first-level interrupt handler (FLIH), which queries the hardware registers. At this point on POWER and POWER2 machines, the interrupt is directly serviced. On PowerPC machines, however, the FLIH will enqueue the interrupt based on level and priority. This raises a flag indicating that there is pending work to be done and it will be serviced later, thus the PowerPC essentially enqueues all interrupts.

The kernel detects queued interrupts at various key times, such as when enabling to a less-favored priority from a more favored one. Once a queued interrupt is detected and the processor is executing at (or about to be enabled to) a priority that lets the pending interrupt be serviced, the current machine state is saved, and interrupt processing is started.

Then the kernel begins calling the interrupt handlers that are registered at the specified level. Because interrupt levels can be shared on certain buses, the adapter which caused the interrupt is not necessarily known at this stage. The order in which the kernel calls interrupt handlers at a certain level is the order in which they were initially registered. This ordering does not change as long as the interrupts are registered. Thus, there is no way to “steal” interrupts from a previously loaded handler at the same level and priority.

Once your second-level interrupt handler (SLIH) is called, it must determine whether the associated adapter caused the interrupt. If the interrupt was caused by the adapter, the interrupt handler does its necessary work, possibly schedules more off-level work, and returns `INTR_SUCC`. If the interrupt was not caused by the adapter, `INTR_FAIL` is returned, and the kernel calls the next handler in the list.

Interrupt Management Kernel Services

The following list contains interrupt management kernel services. For more detailed information on the syntax and return codes, see *AIX Version 4.1 Technical Reference, Volume 5: Kernel and Subsystems*

i_init	Defines an interrupt handler to the system, connects it to an interrupt level, and assigns an interrupt priority to the level.
i_clear	Removes an interrupt handler from the system. Note: A system assert will occur if this service is called for an undefined interrupt handler.
i_disable	Raises the interrupt priority to the specified level, thus disabling all interrupt levels at a less-favored interrupt priority.
i_enable	Restores the interrupt priority to a less-favored interrupt priority, thus enabling all interrupt levels of a higher priority.
i_mask	Disables the specified bus interrupt level.
i_unmask	Enables the specified bus interrupt level.
i_sched	Schedules an off-level interrupt handler to be executed.

The **i_reset** kernel service, which was used to reset a bus interrupt level in releases before AIX Version 3.2.5, no longer needs to be explicitly called. Its function has been moved to the first-level interrupt handler (FLIH), and thus interrupt levels are automatically reset.

Multiprocessor Interrupt Concerns

Ensuring proper synchronized access to the interrupt handlers is of major concern when writing handlers for a multiprocessor (MP) environment. AIX Version 4.1 provides two kernel services that enable you to make your interrupt handling MP-safe:

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.

Essentially, these kernel services should be used wherever **i_enable** and **i_disable** are normally used in a uniprocessor interrupt handler. The simple lock kernel services should *not* be called directly, use **disable_lock** and **unlock_enable** to ensure that a thread will never be interrupted while it holds a simple lock. Allowing a thread which holds a simple lock to be interrupted can deadlock the system.

In an MP environment it is not necessary for all the interrupt handlers to be MP-safe; however, this is a desired goal for increased performance. If an interrupt handler that is not MP-safe is registered, the allocated interrupt level is marked for funnelled operation, which is serviced by the master processor only. Additionally, all the other interrupt handlers at this level will be routed to the master processor. In other words, all interrupt handlers sharing the same level are either considered all MP-safe or all funnelled. Once all interrupt handlers that are funnelled are removed from a level (by calling **i_clear**), any remaining MP-safe handlers on that level will be allowed to run non-funnelled.

For more information and definitions of multiprocessor terms, see “Understanding Multiprocessor-Safe Device Drivers” in *AIX Version 4.1 Kernel Extensions and Device Support Programming Concepts*

Chapter 4. Memory Management

This chapter discusses the various system calls typically utilized by kernel extensions and device drivers in the manipulation of kernel memory. These services include:

- Memory Allocation Services, on page 4-1 (allocate and free kernel memory)
- Memory Pinning Services, on page 4-3 (pin and unpin kernel memory)
- Memory Access Services, on page 4-5 (transfer data between user memory and kernel memory)
- Virtual Memory Management Services, on page 4-6 (manage virtual memory)
- Cross-Memory Services, on page 4-12 (perform cross-memory transfers)

This chapter discusses only the most commonly used kernel memory services available to device driver developers. For more information about these kernel services and additional kernel services, see the Hypertext Information Base Library 1.1 for AIX, which is viewed with InfoExplorer, or *AIX Version 4.1 Technical Reference, Volume 5: Kernel and Subsystems* and *AIX Version 4.1 Technical Reference, Volume 6: Kernel and Subsystems*

Memory Allocation Services

The following are common memory allocation services:

- **xmalloc**
- **xmfree**
- **init_heap**

xmalloc

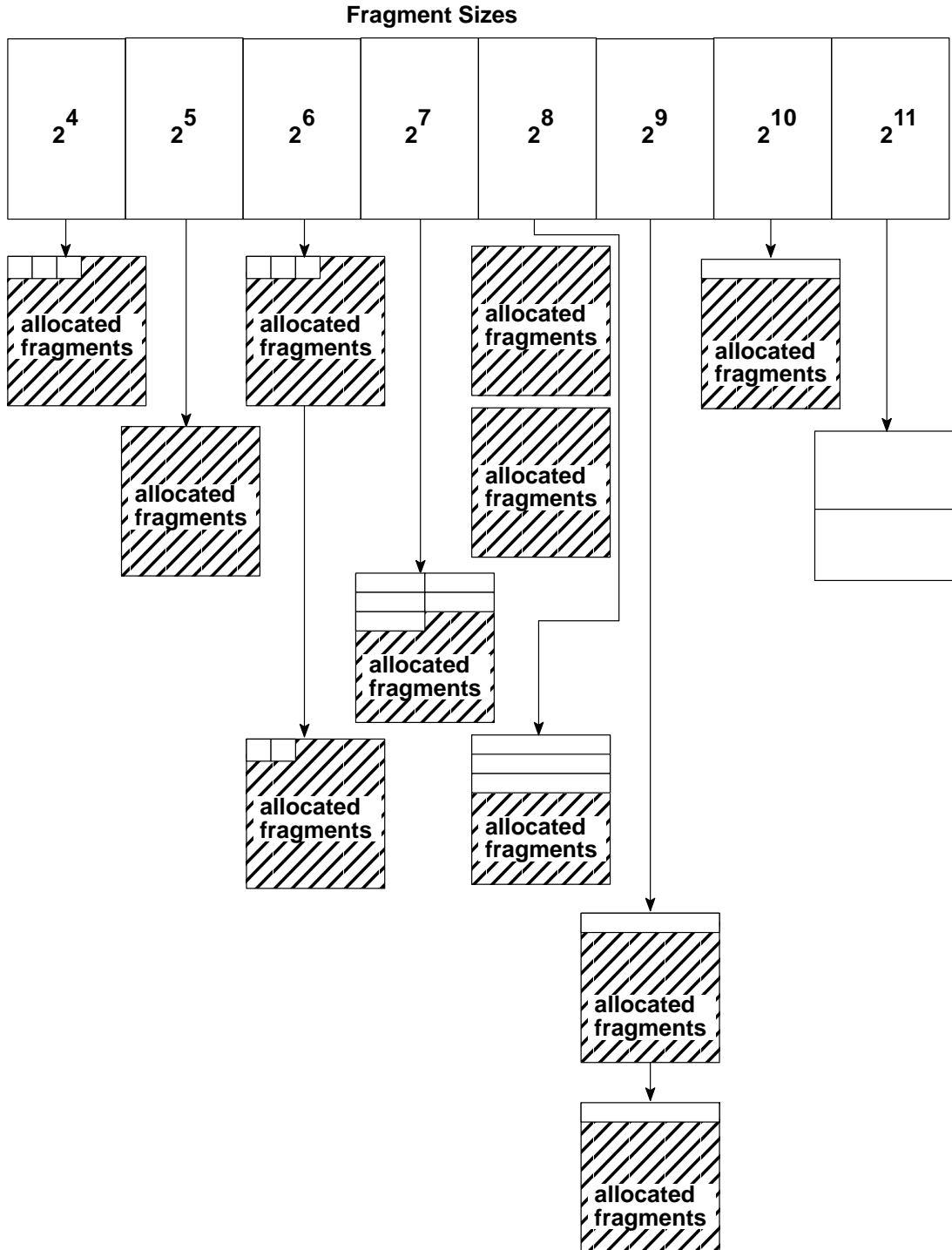
The **xmalloc** kernel service allocates an area of memory from either the kernel heap or the pinned kernel heap. Memory should be allocated from the pinned heap if it is intended to always remain pinned or remain pinned for a long period of time. If the allocated memory can be paged out, it should be requested from the kernel heap. Any unpinned memory can be pinned with the **pin** and **unpin** system calls at a later time if necessary.

The memory area returned by this service can be allocated on a boundary that is a power of 2 bytes from 16 bytes up to a page boundary of 4096 bytes (16-byte boundary, 32-byte boundary, and so on up to a 4096-byte boundary.)

For requests less than one page, **xmalloc** rounds up the request to the next higher power of 2. This implies that a request of just more than half a page is rounded up to one page. **xmalloc** also allocates a minimum of 16 bytes at a time due to the allocation algorithm it utilizes.

The **xmalloc** kernel service allocates requests differently based on the size of the request. For requests of half a page or less (requests greater than half a page are rounded up to one page), a vectored array is kept of elements that represent powers of 2 up to half a page. Each element is an anchor to a list of allocated pages that are individually divided into fragments. These fragments are equal to the size represented by the anchor array element. Once a page is used up (all fragments are allocated), a new page is then grabbed and divided and placed onto the appropriate list. This algorithm allows for fast allocation of areas less than half a page because lists of fragments of the correct size have already been allocated and are immediately ready for use by the caller. Although this method may seem wasteful since it preallocates a whole page of fragments for each size, the worst case only requires 8 pages to be pre-allocated with all free fragments.

The following figure shows the layout of the array used to keep track of the pages that are divided into the various fragment sizes. Note that as long as a page contains a free fragment, it is kept on the linked list for its size. Once all fragments on a page are allocated, the page is no longer linked to the array. When one fragment comes free on a fully allocated page, that page is reinserted onto the list corresponding to its fragment size.



Xmalloc Array of Fragments Less Than 1 Page

For allocations of greater than half a page, **xmalloc** allocates enough pages to contain the requested allocation size.

xmalloc cannot be called with interrupts disabled.

xmfree

Essentially, the **xmfree** kernel service frees up memory areas allocated through the **xmalloc** kernel service. In the case of allocated regions one page or greater, **xmfree** merely frees up the allocated pages.

In the case of memory fragments less than one page, **xmfree** must first find the page to which the fragment belongs as it resides in the vectored array as described in the discussion of **xmalloc**. Once the fragment is freed, the page to which it belongs is also freed if the freed fragment completes an entire page of freed fragments.

xmfree cannot be called with interrupts disabled.

init_heap

The **init_heap** kernel service allows a device driver or kernel extension to set aside an area of memory as a heap for private use. This reserved area must be page aligned and may be a subset of another heap. The **xmalloc** and **xmfree** kernel services are used to allocate and deallocate memory from this private heap.

init_heap cannot be called with interrupts disabled.

Memory Pinning Services

The following are common memory pinning services:

- **ltpin**
- **pin**
- **pincode**
- **pinu**
- **ltunpin**
- **unpin**
- **unpincode**
- **unpinu**

ltpin

The **ltpin** kernel service enables device drivers and kernel extensions to perform long-term pinning of pages of kernel memory. This prevents them from being paged out. All pages touched by the specified address range are pinned, not just the range itself. If the defined range overlaps just slightly into one page, the entire page is pinned.

The **ltpin** service increments a pin count for each page. A page is not a candidate to be paged out until its pin count reaches zero, typically by later calls to the **ltunpin** service.

The major difference between this service and the short-term **pin** kernel service is that long-term pinned pages have their corresponding paging space allocation freed but short-term pinned pages retain their paging space allocation.

There is also a limit to the number of long-term pinned pages allowed in the system. The maximum number of long-term pinned pages is just under 32767 pages.

ltpin cannot be called with interrupts disabled.

pin

The **pin** kernel service enables device drivers and kernel extensions to perform short-term pinning of pages of kernel memory. This prevents them from being paged out. All pages touched by the specified address range are pinned, not just the range itself. Therefore, if the defined range overlaps just slightly into one page, that entire page is pinned.

The **pin** service increments a pin count for each page. A page is not a candidate to be paged out until its pin count reaches zero, typically by later calls to the **unpin** service.

There is also a limit to the number of short-term pinned pages allowed in the system. The maximum number of short term pinned pages is just under 32767 pages.

pin cannot be called with interrupts disabled.

pincode

The **pincode** kernel service supports long-term pinning of an entire object module's code and data. This service calculates the address and length of a module's code and data and calls the **ltpin** kernel service with these values to perform the actual pinning.

The **pincode** service increments a pin count for each page. A page is not a candidate to be paged out until its pin count reaches zero, typically by later calls to the **unpincode** service.

pincode cannot be called with interrupts disabled.

pinu

The **pinu** kernel service provides short-term pinning of memory in either user or kernel space. This routine calls the **pin** kernel service after performing an attach to the user's address space.

The **pinu** service increments a pin count for each page. A page is not a candidate to be paged out until its pin count reaches zero, typically by later calls to the **unpinu** service.

pinu cannot be called with interrupts disabled.

ltunpin

The **ltunpin** kernel service unpins the memory pages long-term pinned by the **ltpin** kernel service.

This service only decrements the pin count of one or more pages. A page is only a candidate to be paged out when its pin count reaches zero.

The **ltunpin** kernel service cannot be called with interrupts disabled because it may need to allocate paging space for the newly unpinned memory area.

unpin

The **unpin** kernel service unpins the memory pages short-term pinned by the **pin** kernel service.

This service only decrements the pin count of one or more pages. A page is only a candidate to be paged out when its pin count reaches zero.

The **unpin** kernel service can be called with interrupts disabled.

unpincode

The **unpincode** kernel service unpins a module's code and data that was pinned by the **pincode** kernel service. Once it has calculated the addresses of a module's code and data segments, **unpincode** calls the **ltunpin** kernel service to perform the actual unpinning.

This service only decrements the pin count of one or more pages. A page is only a candidate to be paged out when its pin count reaches zero.

unpinu

The **unpinu** kernel service unpins memory pinned by the **pinu** kernel service. As in the **pinu** kernel service, memory from either user or kernel space can be freed. This routine calls the **unpin** kernel service to perform the actual unpin.

This service only decrements the pin count of one or more pages. A page is only a candidate to be paged out when its pin count reaches zero.

This routine must be called with all interrupts disabled when unpinning memory in user space. This routine can only unpin memory in kernel space if interrupts are disabled.

Memory Access Services

The following are common memory access services:

- **copyin**
- **copyinstr**
- **copyout**
- **uiomove**

copyin

The **copyin** kernel service copies data from user memory to kernel memory with appropriate exception handling.

copyin cannot be called with interrupts disabled.

copyinstr

The **copyinstr** kernel service copies a string from user memory to kernel memory with appropriate exception handling. This service copies up to the number of bytes specified or until a NULL byte is reached.

copyinstr cannot be called with interrupts disabled.

copyout

The **copyout** kernel service copies data from kernel memory to user memory with appropriate exception handling.

copyout cannot be called with interrupts disabled.

uiomove

The **uiomove** kernel service copies data between kernel memory and either user or kernel memory as defined by a **uio** structure. When using cross-memory descriptors, **uiomove** calls either **xmemin** or **xmemout** for the actual transfer. Cross-memory transfers also require a valid `uio_xmem` pointer to an array of **xmem** structures.

uiomove cannot be called with interrupts disabled.

Virtual Memory Management Services

Virtual memory describes the hierarchy of both main system memory and secondary storage such as hard disks. This two-tier structure allows actual physical memory to be divided and allocated to several processes at the same time. It also allows programs to address memory far larger than the actual size of physical system memory because secondary storage can be used as an extension of the system memory.

Virtual memory objects are often used by device drivers to represent data that can be in either system memory or on a secondary storage device. The virtual memory address space is divided into *segments*, which are actually 256MB contiguous areas of this address space. Process addressability is managed at the segment level. These segments can then be kept private to a specific process or can be shared among several processes. Segments themselves are further divided into *pages*, which are currently 4096 bytes.

There are several types of segments. A *working* segment is a segment that does not have corresponding permanent storage space. For example, the stack and data region for a process are mapped to a working segment because they exist only within the context of the existence of the process. Once there is not enough physical memory to hold a working segment's pages, some of these pages must be transferred, or *paged out* to paging space, which is a temporary area of secondary storage (for example, hard disk). Once these pages are again needed by the process, they can then be *paged in* or transferred back into system memory.

Persistent segments are mapped segments that have a corresponding permanent storage area in secondary storage. Files on disk that are mapped are examples of *mapped* segments. When a page of a persistent segment must be paged out, the actual page must be written to secondary storage if it has changed while being mapped. If the page has not changed, it is merely discarded and its place in physical memory is now open for another page to be paged in.

Client segments are persistent segments that map files remotely, such as over an NFS mount. When pages for these types of segments are paged out, they are written out over the network.

Virtual memory objects are often used by device drivers to do things such as mapping a file that exists on a hard disk in secondary storage.

The following is an example of the usual steps to create and manipulate a virtual memory object that maps in an NFS mounted file:

```

#include <sys/vmuser.h>
#include <sys/types.h>
#include <sys/m_types.h>

vmid_t      vmid;
vmhandle_t  demo_vm_handle;
int         type;
struct gnode *gn;
int         size;
int         num_vm_bufs;
int         rc;
extern int  demo_nfs_strategy();
caddr_t     buffer;

num_vm_bufs = 50;      /* 50 is arbitrary for the number of buf
                        structures to allocate */
                        /* demo_nfs_strategy() is the strategy
                        routine that handles the pageouts */
rc = vm_mount(D_REMOTE, demo_nfs_strategy, num_vm_bufs);
if (rc)
    return (rc);
    . . .

type = V_CLIENT; /* going to work with a "client" segment */
/* gn is a pointer to the gnode of the file to be mapped */
/* size is the size in bytes of the file to be mapped */
if (rc = vms_create(&vmid, type, gn, size, 0, 0))
    return (rc);
/* vms_create() actually returns a segment ID in vmid */

demo_vm_handle = vm_handle(vmid, 0);
/* vm_handle() actually returns a segment register value */
    . . .

buffer = vm_att(demo_vm_handle, 0);
                        /* we're just picking 0 as the offset */
/* Note that the kernel will panic if it runs out of spare
   segment registers but there should be enough to go around.
   As a general rule, a device driver should limit itself to 2
   attaches at a time. */

/* The buffer can now be written to and read from */
    . . .

vm_det(buffer);
    . . .
if (rc = vms_delete(vmid))
    return (rc);
    . . .
vm_umount(D_REMOTE, demo_nfs_strategy);

```

The kernel services used to manage virtual memory generally have names beginning “vms_” (for example, **vms_create**, **vms_delete**, and **vms_iowait**) or names beginning “vm_” (for example, **vm_handle**, **vm_att**, and **vm_cflush**.)

vms_create

The **vms_create** kernel service creates a virtual memory object of a specific size and type. Several types and options are available when creating virtual memory objects.

Segments can have the following types:

V_WORKING	Working storage segment
V_CLIENT	Client segment
V_PERSISTENT	Persistent storage segment
V_MAPPING	Mapping segment

Working storage segments can have the following options:

V_UREAD	Only reads are allowed in user region with access key 1
V_PTASEG	Segment is a page table area segment
V_SYSTEM	Segment is a system segment.
V_PRIVSEG	Segment is a process private segment
V_SHRLIB	Shared library segment
V_SPARSE	Segment has pages sparsely populated

Persistent segments can have the following options:

V_JOURNAL	Segment is also journalled
V_DEFER	Segment is deferred update
V_LOGSEG	Segment is used for a log
V_SYSTEM	Segment is a system segment

Client segments can have the V_INTRSEG option, indicating the segment is an interruptable segment.

There are no options for mapping segments. All parameters other than segment ID (*sid*) and type should be 0.

vms_create cannot be called with interrupts disabled.

vms_delete

The **vms_delete** kernel service deletes a virtual memory object created through the **vms_create** kernel service. Even though this service completes asynchronously, notification of completion is given synchronously. This is because the allocated segment is not truly freed until all paging I/O has completed. Until the segment has been marked as freed, it remains accessible even though a successful return code may have been received.

vms_delete cannot be called with interrupts disabled.

vm_handle

The **vm_handle** kernel service creates a virtual memory handle for a virtual memory object, as created by **vms_create**, for use by the **vm_att** kernel service. The handle is created from a concatenation of the segment ID and a protection key.

vm_handle cannot be called with interrupts disabled.

vm_att

The **vm_att** kernel service maps a virtual memory object, created by the **vms_create** kernel service, by allocating a free segment register and returning a 32-bit address made up of the

segment register value and the offset within that segment. This 32-bit address has the segment register number in the upper 4 bits and the offset in the lower 28 bits.

If the system has no segment registers to allocate, it will panic.

vm_att can be called with interrupts disabled.

vm_cflush

The **vm_cflush** kernel service flushes out to memory all processor instruction and data cache lines that contain the memory boundaries specified by the given address and range.

vm_cflush can be called with interrupts disabled.

vm_det

The **vm_det** kernel service essentially unmaps the virtual memory object originally mapped by the **vm_att** kernel service. It does so by releasing the segment register associated with the given virtual address of the memory object.

vm_det can be called with interrupts disabled.

vm_mount

The **vm_mount** kernel service allocates an entry in the paging device table (PDT) for a file system and also allocates the required **buf** structures for processing by the strategy routine.

vm_mount cannot be called with interrupts disabled.

vm_umount

The **vm_umount** kernel service removes an entry from the paging device table (PDT) for a file system and also deallocates the required **buf** structures that were previously allocated by the **vm_mount** kernel service. If paging activity has not yet ceased when **vm_umount** is called, the service waits until all I/O has completed before completely freeing all resources.

vm_umount cannot be called with interrupts disabled.

vm_move

The **vm_move** kernel service will transfer data from a virtual memory object, created by the **vms_create** kernel service, and a buffer pointed to by a **uio** structure.

vm_move cannot be called with interrupts disabled.

The following is an example of using **vm_move** on a permanent hard disk file:

```
#include <sys/vmuser.h>
#include <sys/types.h>
#include <sys/m_types.h>

vmid_t    vmid;
int       type;
int       size;
dev_t     dev_num;
int       rc;
struct uiop *uiop;
. . .

type = V_PERSISTENT; /* we want a "persistent" segment */
/* dev_num is the devno of the block device on which
   the inode resides */
/* size is the size in bytes of the file */
if (rc = vms_create(&vmid, type, dev_num, size, 0, 0))
    return (rc);
if (rc = vm_move(vmid, uiop->uio_offset, uiop->uio_resid,
    UIO_READ, uiop));
    return (rc);
if (rc = vms_delete(vmid));
    return (rc);
. . .
```

vm_write

The **vm_write** kernel service initiates a page-out for all pages touched by a specified range in a virtual memory object. The memory range is defined by a beginning address and length. All pages touched by this range are marked for page-out.

A *force* parameter, applicable only to journaled segments, is provided to force a page-out on a page even though it may have been recently modified. Typically, if a page has been recently modified, it will most likely be modified again. By delaying page-outs on recently modified pages, time is saved by writing the page only after several modifications rather than after each modification. The *force* parameter overrides this process.

vm_write cannot be called with interrupts disabled.

vm_writep

The **vm_writep** kernel service is similar to the **vm_write** kernel service except that a page range, rather than a memory range, is specified for page-outs.

No *force* parameter is provided.

vm_writep cannot be called with interrupts disabled.

vms_iowait

The **vms_iowait** kernel service waits until all page-out I/O for a virtual memory object is complete.

vms_iowait cannot be called with interrupts disabled.

vm_release

The **vm_release** kernel service releases all pages touched by a specified range in a virtual memory object. In essence, an address range is defined by a beginning address and length. All pages touched by this range are freed.

vm_release cannot be called with interrupts disabled.

vm_releasep

The **vm_releasep** kernel service will release the specified pages in a virtual memory object. It is similar to the **vm_release** kernel service except that pages are specified rather than an address range.

vm_releasep cannot be called with interrupts disabled.

Example Using Virtual Memory Management Services

The following example shows typical use of the **vm_wriep**, **vms_iowait**, and **vm_releasep** kernel services:

```
#include <sys/vmuser.h>
#include <sys/types.h>
#include <sys/m_types.h>

vmid_t    vmid;
int       type;
int       size;
dev_t     dev_num;
int       pfirst;
int       npages;
int       rc;
. . .

type = V_PERSISTENT; /* we want a "persistent" segment */
/* dev_num is the devno of the block device on which
   the inode resides */
/* size is the size in bytes of the file */
if (rc = vms_create(&vmid, type, dev_num, size, 0, 0))
    return (rc);
. . .

/* pfirst is the page number of the first page to pageout */
/* npages is the number of pages we want to pageout */
if (rc = vm_wriep(vmid, pfirst, npages));
    return (rc);
rc = vms_iowait(vmid);
vms_releasep(vmid, pfirst, npages);
if (rc)
    return (rc);
. . .
```

Cross-Memory Services

The cross-memory kernel services enable device drivers to access user-mode data. This is often required by the interrupt handler section of device drivers or when performing DMA transfers to and from a user buffer.

The following are common cross-memory services:

- **xmattach**
- **xmdetach**
- **xmemin**
- **xmemout**
- **xmemdma**

The **xmattach** kernel service provides a descriptor to access the user-mode memory region while the **xmemin** and **xmemout** kernel services transfer data to and from this region. The **xmdetach** kernel service then deallocates resources used to access the user-mode region and prevents further accesses to that region.

The following example shows a typical transfer of user data to a kernel buffer using the **xmemin** kernel service:

```
#include <sys/xmem.h>

caddr_t  user_buffer;
caddr_t  local_buffer;
uint     buffer_len;
struct xmem  dp;
int      rc;

bzero(&dp, sizeof(struct xmem));
/* aspace_id field must be initialized to XMEM_INVAL */
dp.aspace_id = XMEM_INVAL;

/* user_buffer should have the address of the user buffer */
/* user_buffer_len should have the length of the user buffer */
if (rc = xmattach(user_buffer, user_buffer_len, &dp, USER_ADSPACE))
    return (rc);

/* allocate word-aligned kernel bufr to hold copy of user data */
local_buffer = xmalloc(user_buffer_len, 2, pinned_heap);
if (local_buffer == NULL)
    return (ENOMEM);
if (rc = xmemin(user_buffer, local_buffer, buffer_len, &dp))
    return(rc);
/* The data in local_buffer can now be sent to the device
   if desired (eg. through PIO) */
. . .

xmfree(local_buffer, pinned_heap);
xmdetach(&dp);
```

xmattach

The **xmattach** kernel service gives a device driver access to a user buffer without having to execute under the process that initiated the I/O. It returns a cross-memory descriptor if the attach is successful. The **xmemin** and **xmemout** kernel services can then be used to transfer data to and from the attached user buffer.

xmattach cannot be called with interrupts disabled.

xmdetach

The **xmdetach** kernel service detaches a user buffer, previously attached by the **xmattach** kernel service, from a device driver. This prevents a device driver from further accesses to a user buffer.

xmdetach can be called with interrupts disabled.

xmemin

The **xmemin** kernel service transfers data from an attached user buffer to a kernel buffer. The device driver performing the transfer should have previously attached to the user buffer using the **xmattach** kernel service.

xmemin can be called with interrupts disabled.

xmemout

The **xmemout** kernel service transfers data from a kernel buffer to an attached user buffer. The device driver performing the transfer should have previously attached to the user buffer using the **xmattach** kernel service.

xmemout cannot be called with interrupts disabled.

xmemdma

The **xmemdma** kernel service prepares a page for DMA I/O or processes a page after DMA I/O is complete. Even though **xmemdma** can be called with interrupts disabled, the page being processed must be in memory and either pinned or in the pager I/O state.

The following flags are used when preparing the page for DMA I/O:

XMEM_HIDE Flushes the cache and invalidates the page if this is the first hide for the page. (On a PowerPC machine, this flag instructs the kernel service to return the calculated real address for the page and if the page is not read-only, set the modified bit.)

XMEM_ACC_CHK Checks the page protection bits for this page before the flush and hide.

XMEM_WRITE_ONLY When used with the **XMEM_ACC_CHK** flag, indicates that page access is read-only and DMA transfers from this page will only be outward.

The following flag processes a page after DMA I/O:

XMEM_UNHIDE Decrements the hide count on this page. If this is the last unhide and the page is not in the pager I/O state, processes waiting on the page are taken off the wait queue and the page's modified bit is set unless it was a read-only page. (On a PowerPC machine, this flag causes the kernel service to return 0.)

xmemdma can be called with interrupts disabled.

The following is an example of the use of the **xmemdma** kernel service to unhide a page after a DMA transfer using the **d_master** kernel service:

```
int          channel_id;
caddr_t     buf_addr;
int         buf_size;
struct xmem xmem_buf;
caddr_t     dma_addr;
. . .
/* We're assuming that the channel_id, buf_addr, buf_size,
   xmem_buf, and dma_addr have already been setup */
d_master(channel_id, DMA_READ, buf_addr, buf_size, &xmem_buf,
         dma_addr);
xmemdma(&xmem_buf, buf_addr, XMEM_UNHIDE);
```

Chapter 5. Synchronization and Serialization

It is often necessary for a device driver to synchronize its access to a resource that is shared with a device or with kernel routines executing in other contexts. Kernel routines in all contexts share the same kernel segment and all data structures within it. Drivers for many devices share the same I/O space and system bus. A device driver's routines all access the same device, so they must synchronize their access to that device.

To synchronize access with another device, either the driver may *poll* the device by repeatedly checking to see if the resource is available, or the device may notify the driver that the resource is available by sending an interrupt whose handler posts an *event* that the driver is waiting for.

If a driver routine accesses a resource that it shares with a kernel routine concurrently executing in another context, then both routines must address various concurrency issues. For example, you want to avoid the following conditions:

- race conditions** Multiple contexts access the same shared resource with different results depending on the order in which the accesses are made. (Accesses to a shared resource are serialized if they are coordinated so there is no race condition.)
- deadlock** Each context awaits a notification the other never sends.
- livelock** Each context polls a condition the other never sets. (This is sometimes called a type of deadlock.)
- starvation** Other contexts monopolize access to the shared resource.

A portion of a routine that accesses a shared resource, a routine's *critical section* makes the access *atomic* relative to routines executing in other contexts whenever it synchronizes its own access to that resource with respect to any access attempts by other routines, so that all accesses are serialized (made one after the other). A routine synchronizes access with other routines either by using semaphores (such as locks) or by invoking *monitors*, which are routines that ensure atomic access.

A routine is *reentrant* if its critical sections make resource accesses that are atomic relative to the same routine executing in some other context. In other words, a routine is reentrant if it can safely execute in many contexts concurrently.

If a routine maintains state information somewhere (for example, in static or global data) so that a subsequent call to the same routine would access information modified by a previous call, then the routine is not reentrant: the state information is the shared resource, and that portion of the routine that accesses the information is that routine's critical section. If a routine enables a caller to access any resources that can be shared, such as by returning a pointer to static or global data, then the routine itself is not reentrant; it cannot safely execute in more than one context.

Recall that an interrupt handler routine must be reentrant because a device may generate an interrupt while another interrupt (possibly from the same device through another system processor) is being handled. In AIX, device driver routines on the call side also must be reentrant because a call-side driver routine executes in the context of a kernel thread whose execution can be preempted in favor of another thread of greater priority, which might be executing the same call-side driver routine.

For information on making interrupt-side routines reentrant, see Chapter 3, "Interrupts."

Timer Services

A driver routine can poll a device actively by checking the shared resource's availability and then doing something else (keeping a system processor), or it can poll a device by checking and then relinquishing its system processor for a certain period of time (going to sleep). Because the time between resource availability checks with active polling depends on the speed of a system processor, and because that time period may not be compatible with device being polled, timer-based polling is more robust.

Each processor that AIX supports has a decremator that periodically generates *timer interrupts* whose handler, among other things, can invoke routines specified to be called once a certain amount of time has passed. The kernel services and data structures associated with timer interrupt handler are *timers*. There are basically two kinds of timers: watchdog timers and real-time timers. A driver routine sets up a timer by specifying a *callback routine* which is an interrupt-side routine that is to be invoked by the timer interrupt handler. The driver routine also specifies when to invoke the callback routine as either an absolute time (from January 1, 1970) or a time relative to when the timer is started.

Watchdog Timers

Driver routines typically use timers to check if device I/O requests had successfully completed. A *watchdog* timer may be satisfactory to poll a device, as it is fairly simple to use and has low overhead. However, it involves writing an interrupt-side callback routine that is little more than a timed "go to" routine: a watchdog timer handler receives no user defined parameters. Also, a watchdog timer can only be set to the nearest second.

Here is an example of how to set up a watchdog timer. The following code sections are written to augment the sample device driver shown in Chapter 1, "Device Driver Overview."

The callback routine, which is invoked any time the watchdog timer expires, executes on the interrupt-side; so it will have to be in the device driver's bottom half. In other words, the callback routine will have to be explicitly pinned in RAM.

Here is the bottom half, `xyz_bot.c`:

```
#include <sys/types.h> /* for dev_t and other types */
#include <sys/trchkid.h> /* for trace hook macros */
/*****
WATCHDOG TIMER CALLBACK ROUTINE
This callback routine is invoked every time a timer expires.
This executes in the context of decremator interrupt handler.
*****/
void watchdog_callback()
{
    TRCHKL1T(HKWD_USER1, 0xb);
}
```

Since the callback routine has no parameters, it must access some shared resources to do anything useful. It may read an adapter's registers to poll a device, or it may see if some other resource is available for use. Portions of the callback routine that access such shared resources are critical sections shared among threads and interrupt handlers and require appropriate synchronization.

The driver's top half, called `xyz_top.c`, is the same as the sample driver in Chapter 1 with several additions.

Add the following to the declaration section:

```

#include <sys/malloc.h> /* for xmalloc() */

/* for watchdog timer support */
#include <sys/watchdog.h> /* for watchdog structure */
extern void watchdog_callback();
struct watchdog *watchp;

```

Note that the pointer to the watchdog timer structure is global and is therefore a shared resource requiring synchronization if the top half were to be reentrant. For simplicity, access to the pointer is not serialized in this example.

Add the following to the `xyzopen` entry point:

```

watchp = (struct watchdog *) xmalloc(sizeof(struct watchdog),
                                     0, pinned_heap);
if(watchp == (struct watchdog *) NULL)
{ setuerror(ENOMEM); /* failed, no space on pinned heap */
  return(-1); /* setuerror() used: there is no ublock access */
}

watchp->func = (void (*)(void *))watchdog_callback;
w_init(watchp);

```

Add the following to the `xyzclose` entry point:

```

w_stop(watchp);
w_clear(watchp);
xmfree(watchp, pinned_heap);

```

Add the following to the `xyzread` entry point:

```

watchp->restart = 3; /* 3 seconds */
w_start(watchp);

```

Add the following to the `CFG_INIT` portion of the `xyzconfig` entry point:

```

if((return_code = pincode(watchdog_callback)) != 0)
{ setuerror(return_code);
  return(-1);
}

```

And, add the following to the `CFG_TERM` portion of the `xyzconfig` entry point:

```

if((return_code = unpincode(watchdog_callback)) != 0)
{ setuerror(return_code);
  return(-1);
}

```

The corresponding stanzas of the makefile need to be changed:

```

xyz: xyz_top.o xyz_bot.o
      ld -e xyzconfig -o xyz -bI:$(KSYSLIST) -bI:$(SYSLIST)
xyz_top.o xyz_bot.o                # have on one line

xyz_top.o: xyz_top.c
      cc -c -D_ALL_SOURCE -D_POSIX_SOURCE -D_KERNEL xyz_top.c

xyz_bot.o: xyz_bot.c
      cc -c -D_ALL_SOURCE -D_POSIX_SOURCE -D_KERNEL xyz_bot.c

```

The user program, `aprogram`, only needs to have `“sleep(5);”` placed after the call to **read** so that the timer can go off before the program terminates.

The following sample trace report shows that the callback routine (hookdata 0xB) executed a little less than three seconds after the call to `xyzread` (hookdata 0x9):

```
trace -j010 -l -l -s -a
```

ID	PROCESS NAME	I	SYSTEM CALL	ELAPSED	APPL	SYSCALL	KERNEL	INTERRUPT
001	trace	0.000000	trace	ON	channel	0		
010	trace	20.121544	UNDEFINED	trace	ID	idx 0x2100	trace	id 0010
	hookword	10E0000	type	0E				
	hookdata	0000 00000001	00630000	00000001	2FF97F1C	00000000		
010	trace	26.447790	UNDEFINED	trace	ID	idx 0x2154	trace	id 0010
	hookword	10E0000	type	0E				
	hookdata	0000 00000001	00630000	00000003	2FF97F1C	00000000		
010	trace	26.447795	UNDEFINED	trace	ID	idx 0x2170	trace	id 0010
	hookword	10A0000	type	0A				
	hookdata	0000 00000005						
010	trace	63.647987	UNDEFINED	trace	ID	idx 0x22c4	trace	id 0010
	hookword	10E0000	type	0E				
	hookdata	0000 00000007	00630000	00000003	00000000	00000000		
010	trace	63.648657	UNDEFINED	trace	ID	idx 0x22e0	trace	id 0010
	hookword	10E0000	type	0E				
	hookdata	0000 00000009	00630000	2FF97DC0	00000000	00000000		
010	trace	66.057338	UNDEFINED	trace	ID	idx 0x2314	trace	id 0010
	hookword	10A0000	type	0A				
	hookdata	0000 0000000B						
010	trace	70.649700	UNDEFINED	trace	ID	idx 0x2348	trace	id 0010
	hookword	10E0000	type	0E				
	hookdata	0000 0000000A	00630000	2FF97DC0	00000000	00000000		
010	trace	70.649833	UNDEFINED	trace	ID	idx 0x2364	trace	id 0010
	hookword	10E0000	type	0E				
	hookdata	0000 00000008	00630000	00000000	00000000	00000000		
002	trace	77.524587	trace	OFF	channel	0		

There are many similarities between routines using watchdog timers and those using real-time timers. This is shown in the following discussion of real-time timers.

Real-Time Timers

If the device requires setting a timer to within something finer than a second, or if the callback routine needs user defined parameters, then the driver will use real-time timers, which make use of Timer Request Blocks (TRBs) as defined in the file **/usr/include/sys/timer.h**.

The real-time kernel services are:

talloc	Allocates a TRB
tstart	Submits the timer request (to place the TRB on a timer queue)
tstop	Removes the timer request
tfree	Frees up resources allocated for the TRB

Here is an example of a way to set up a real-time timer. The following code sections are written to augment the sample device driver shown in Chapter 1, "Device Driver Overview."

As with watchdog timers, the callback routine, specified to be called once the timer expires, executes on the interrupt-side and so will have to be in the device driver's bottom half.

Here is the bottom half, `xyz_bot.c`:

```

#include <sys/types.h> /* for dev_t and other types */
#include <sys/errno.h> /* for errno declarations */
#include <sys/trchkid.h> /* for trace hook macros */

/* for timer support */
#include <sys/timers.h> /* for itimerspec, includes time.h */
#include <sys/timer.h> /* for TRBs */

/*****
TIMER CALLBACK ROUTINE
This callback routine is invoked every time a timer expires.
This executes in the context of decremter interrupt handler.
*****/
void timer_callback(struct trb *timer_req_blk_ptr)
{
    TRCHKL2T(HKWD_USER1, 0xb, timer_req_blk_ptr->t_func_sdata);
}

```

The callback routine, `timer_callback`, has a TRB, which can contain some user supplied data, passed to it; in this case `t_func_sdata` is a signed integer containing, say, the device number. Any data that can't be passed as a parameter will have to be globally accessible; portions of a routine that access such shared resources are critical sections shared among threads and interrupt handlers and require appropriate synchronization.

The driver's top half, called `xyz_top.c`, differs from the sample driver in Chapter 1 in several ways described in the following discussion. One difference is the include files and declarations would include:

```

/* for timer support */
#include <sys/timers.h> /* for itimerspec, includes time.h */
#include <sys/timer.h> /* for TRBs */
#include <sys/m_intr.h> /* for INTTIMER priority */

extern void timer_callback();
struct trb *trbp; /* single pointer to TRB structure */

```

Note that the pointer to TRB is global and is therefore a shared resource requiring synchronization if the top half is to be reentrant. For simplicity, access to the pointer is not serialized in this example.

Here is what is added to the entry point `xyzopen`:

```

if((trbp = talloc()) == (struct trb *) NULL)
{
    setuerror(ENOMEM); /* talloc failed, no space on pin heap */
    return(-1); /* setuerror() used: there is no ublock access */
}

trbp->id = thread_self(); /* or getpid() */
trbp->timerid = TIMERID_REAL;
trbp->eventlist = -1; /* no events being waited on */
trbp->func = (void (*)(void *))timer_callback;
trbp->t_func_sdata = (int) devno; /* or whatever data */
trbp->ipri = INTTIMER;

```

Note that `talloc` allocates space from the pinned heap. The fields of the TRB listed are the only ones a device driver would set; even so, most fields only need to be filled out if the callback routine requires them. `id`, `timerid`, `eventlist`, `t_func_sdata` are all for use by the callback routine. The callback routine, `timer_callback`, listed previously, currently makes no use of any of these fields.

Here is what is added to the entry point `xyzclose`:

```
tstop(trbp);
tfree(trbp);
```

Here is what is added to the entry point `xyzread`:

```
struct timespec threeSec;
struct itimerspec timeStruct;
...
threeSec.tv_sec = 3; /* 3 seconds */
threeSec.tv_nsec = 0; /* and no milliseconds */

timeStruct.it_interval = threeSec;
timeStruct.it_value = threeSec;

trbp->timeout = timeStruct;

tstart(trbp);
```

So, if a program issues a `read` on this device, it starts a timer that causes the callback routine to execute three seconds later.

Here is what is added to the `CFG_INIT` case section of `xyzconfig`:

```
if((return_code = pincode(timer_callback)) != 0)
{ setuerror(return_code);
return(-1);
}
```

This pins the bottom-half, containing the callback routine, in RAM.

Here is what is added to the `CFG_TERM` case section of `xyzconfig`:

```
if((return_code = unpincode(timer_callback)) != 0)
{ setuerror(return_code);
return(-1);
}
```

This enables the pages containing the bottom half to be subject to page replacement.

The `makefile` and `aprogram.c` are the same as in the watchdog timer example.

Now, when the kernel is extended, and the user program, `aprogram`, is run, a system trace shows that the callback routine executes three seconds after **read** is called:

```
010 trace 22.459560 UNDEFINED TRACE ID idx 0x21d4 traceid 0010
hookword 10E0000
type 0E
hookdata 0000 00000009 00630000 2FF97DC0 00000000 00000000

010 trace 25.459661 UNDEFINED TRACE ID idx 0x2210 traceid 0010
hookword 10E0000
type 0E
hookdata 0000 0000000B 00630000 00000000 00000000 00000000
```

Event Notification

An alternative to polling a device for when some data is available, or when some status is achieved, is to make use of any interrupts that the device generates and notify the driver that an interrupt was processed. In this case, the driver routine registers itself for an *event* and causes its thread to relinquish the system processor. Whenever the event is posted, the kernel awakens the driver's thread and the driver routine resumes execution.

The kernel service **et_wait** causes the calling thread to relinquish the system processor until that thread receives an event. The kernel service **et_post** enables a kernel process to post an event to some thread, so that the thread can resume execution. The *event* is a bit in a bit mask that is passed as a parameter to each call.

The following is an example of adding event notification to the real-time timer sample code given previously.

In `xyz_bot.c`, add to the declarations:

```
/* for event handling */
#define THE_EVENT (1 << 31) /* this is the high order bit */
```

There is an event mask that can be used to ensure that an event is not reserved for the base operating system kernel. Be aware that this event is specific to a 32-bit system.

Add the following line to the callback routine:

```
et_post(THE_EVENT, timer_req_blk_ptr->id);
```

The field `id`, in the TRB, contains the ID of the thread that had set the timer.

In `xyz_top.c` add to the declarations:

```
/* for event handling */
#include <sys/sleep.h> /* for EVENT_SIGRET flag */
#define THE_EVENT (1 << 31) /* this is a high order bit */
```

And add the following to the `xyzwrite` entry point:

```
et_wait(THE_EVENT, THE_EVENT, EVENT_SIGRET); /* clear event
upon receipt */
TRCHKL1T(HKWD_USER1, 0x0a);
```

Keep `aprogram.c` the same as in Chapter 1. A portion of the trace looks like:

```
010 trace 18.261389 UNDEFINED TRACE ID idx 0x21c8 traceid 0010
hookword 10E0000
type 0E
hookdata 0000 00000000A 00630000 2FF97DC0 00000000 00000000

010 trace 21.261233 UNDEFINED TRACE ID idx 0x21fc traceid 0010
hookword 10E0000
type 0E
hookdata 0000 00000000B 00630000 00000000 00000000 00000000

010 trace 21.261919 UNDEFINED TRACE ID idx 0x2218 traceid 0010
hookword 10A0000
type 0A
hookdata 0000 00000000A
```

The trace shows entry into `xyzwrite`, the callback, and then the completion of `xyzwrite`.

Instead of posting events to a driver routine in a thread, a device driver may need to wait on an event and meanwhile make a shared resource available for another routine to lock. To do so, the driver can invoke the kernel service **e_sleep_thread**, which takes a pointer to an event and a pointer to a lock as parameters. If the event is to be posted from the interrupt side, be careful that the pointers refer to data in pinned memory .

An interrupt handler (or some other routine) may call the kernel service **e_wakeup**, to cause any threads sleeping on this event to resume execution and reacquire a lock on a shared resource.

If multiple locks are to be released and then reacquired, a driver will use the kernel services, **e_assert_wait**, **e_block_thread**, and **e_clear_wait** and **e_wakeup** instead.

For more information on event handling services, see *AIX Version 4.1 Technical Reference, Volume 5: Kernel and Subsystems* and *AIX Version 4.1 Technical Reference, Volume 6: Kernel and Subsystems*

Serialization Services

Sharing data among multiple concurrent processes causes an inherent problem with maintaining data consistency. Consider, for example, a doubly linked list structure. In order to remove an element from this list, it is necessary to update pointers in both the preceding and succeeding the elements. During this update sequence the list is in an inconsistent state. If additional process activity to add or remove elements is allowed to proceed while in this state, the results would be unpredictable. Access could be serialized through the use of locks.

Device drivers have the following types of critical code sections:

- Critical code sections shared among process threads. (Interrupt disabling is not required.)
- Critical code sections shared among process threads and interrupt handlers. (Interrupt disabling is required.)

Uniprocessor (UP) Serialization

On a uniprocessor, concurrent access to shared data is achieved through preemption. A process can be preempted by another higher priority process, or by an interrupt routine.

Protection for access to data shared between process threads is provided by using locks. Protection for access to data shared between the base level and the interrupt level is provided by disabling interrupts.

Multiprocessing (MP) Serialization

Multiprocessing involves the use of more than a single processor. AIX implements shared memory symmetric multiprocessing, that is, all processors are functionally equivalent and can perform I/O and computations. AIX manages a pool of identical processors, any of which may be used to control I/O devices or reference any memory unit. Conflicts between processors attempting to access the same data at the same time are resolved by hardware instruction synchronization. Conflicts in access to shared memory structures are resolved by software synchronization techniques; most notably locking.

The multiprocessing discussion introduces the following new terminology:

Master Processor

The default processor for of funneled code execution. This is usually the IPL boot processor.

Funneling

Funneled code runs only on the master processor. This enables a UP driver to run unchanged by funneling its execution through one specific processor (the Master processor).

MP-safe

Device driver code that can run on any processor. The code for the device driver provides for locks to serialize the device drivers execution. This provides for a coarse level of locking to enable parallel execution.

MP-efficient

Device driver code that can run on any processor. The code for the device driver provides for locks to serialize access to devices and data structures. This provides for a finer level of locking to further enable parallel execution.

Masking interrupts or disabling the processor by calling the **i_disable** kernel service to serialize with an interrupt handler is no longer sufficient with symmetric multiprocessing (except for funneled code). The I/O is symmetric and interrupts can be routed to any of the processors in the complex; masking interrupts will not prevent the interrupt routine from executing simultaneously on another processor. Serialization in a symmetric multiprocessor system, therefore, requires the use of locks in addition to using **i_disable**.

Lock Overview

A number of serialization techniques are available defined by the intent to serialize access to critical code sections, or data items. Locking critical sections of code is a coarse locking level that supports MP-safe program execution. Locking at the data level is a finer locking level that supports either MP-safe or MP-efficient program execution. Additionally, atomic operations such as **fetch_and_add** and **compare_and_swap** can be used as an alternative to locks to provide reliable access to a single shared variable for reading or writing.

Lock contention wait can result in a process spinning on a busy lock or sleeping until the lock is granted.

AIX provides for the following types of locks:

Simple locks Spin locks that provide exclusive ownership and are not recursive. These locks are used for serialization among threads, serialization among threads and interrupt handlers, and serialization among threads and interrupt handlers.

Complex locks Sleep locks that provide read or write access and are recursive on request. These locks are used only for serialization among threads.

lockl locks Sleeping, mutual exclusive locks provided for compatibility with AIX Version 3. These should not be used for newly written code.

Locking the global `kernel_lock` is not recommended (discouraged) because this lock has the potential to block the system scheduler. Also, this lock may be removed at some future time.

To support collecting and monitoring statistics about lock activity, AIX implements lock instrumentation. This function is activated or deactivated on an IPL boot basis via the **bosboot** command. Lock instrumentation criteria requires that locks be allocated before being used and deallocated when no longer needed. Instrumentation is not provided for **lockl** and **unlockl** lock services. Allocation calls and deallocation calls result in no operation if instrumentation is disabled.

Device Driver Lock Models

The device driver lock model in an MP environment can take on one of the forms previously mentioned (funneled, MP-safe, or MP-efficient). The implementation model used by the device driver needs to be communicated to the kernel for the MP-safe or MP-efficient models. This is done by setting necessary flags in the appropriate data structures supporting the **devswadd**, **i_init**, and **iodone** kernel services. The funneled implementation is assumed by default.

The funneled model is essentially a UP implementation. This is appropriate for low throughput devices or migrating device drivers from the UP environment. This is the implementation model that is assumed if no changes are made to an existing UP device driver migrated to an MP environment. Serialization of critical code sections is accomplished by using the **lockl** and **unlockl** kernel services for sections shared among threads, and the

i_disable and **i_enable** kernel services for sections shared among threads and interrupt handlers.

Although **lockl** locks are provided for compatibility with earlier releases, new code written for AIX Version 4.1 should use simple locks instead. This enables gathering lock instrumentation statistics for these locks.

The MP-safe model is intended for medium throughput devices and provides a coarse level of MP serialization without the complexity and overhead of MP-efficient implementation. Simple lock kernel services are provided for both serialization of critical sections shared among threads and serialization of critical sections shared among threads and interrupt handlers. This implementation enables critical sections of code to run on any processor in the complex, but not at the same time. The locking granularity is usually a single global lock at the device driver level.

The MP-efficient model is intended for high throughput devices by optimizing the parallel execution capabilities of an MP processor. A finer level of locking capability is available by using a hierarchy of locks based on data access serialization. Complex lock kernel services are provided for shared read and exclusive write access to data structures shared among multiple-thread code sections. Simple lock kernel services are provided for serialization among thread and interrupt handler code sections.

The implementation model used depends on such factors as performance needs and whether the device driver is being migrated from a UP environment. If designing, or moving to an MP-safe or MP-efficient model, the recommendation is to start by using simple locks (a single global lock initially) and gradually refine the serialization as needed after identifying all the data structures shared between the top and bottom half of the driver. Also, analysis of the lock statistics provided by lock instrumentation can be used to identify contention bottlenecks that may point out the need for additional locking or the need for the read and write locking capability provided by complex locks.

For a complete discussion of the MP locking considerations for a device driver see "Understanding Multiprocessor-Safe Device Drivers" in *AIX Version 4.1 Kernel Extensions and Device Support Programming Concepts*

MP-Safe Coding Sample

```
/* Skeleton code sample using the following lock related      */
/* kernel services:                                          */
/*                                                            */
/* lockl()           - conventional-lock lock request        */
/* lock_alloc()      - allocate simple lock                 */
/* simple_lock_init() - initialize simple lock              */
/* lock_free()       - release simple lock                  */
/* unlockl()         - conventional-lock unlock request     */
/* devswadd()        - devsw table lock model options      */
/*                                                            */
/* (driver configuration lock definition and use)           */
...
lock_t config_lock = {LOCK_AVAIL}; /* define lockl lock      */
Simple_lock dd_lock; /* define Simple MP lock */
...
int
dd_config(int cmd,
          struct uio *uiop)
{ /* start dd_config() */
  struct devsw dd_devsw;
  int rc;
  ...
}
```

```

rc = lockl(&config_lock, LOCK_SHORT); /* dd_config lock */
if (rc != LOCK_SUCC)
return(EINVAL);
lock_alloc(&dd_lock,          /* allocate Simple MP lock */
           LOCK_ALLOC_PIN,
           DD_LOCK,
           -1);
simple_lock_init(&dd_lock); /* initialize Simple MP lock */
...
switch(cmd)
{
case CFG_INIT:
{
dd_devsw.d_open = dd_open;
dd_devsw.d_close = dd_close;
dd_devsw.d_read = dd_read;
dd_devsw.d_write = dd_write;
dd_devsw.d_ioctl = dd_ioctl;
dd_devsw.d_strategy = dd_strategy;
dd_devsw.d_ttys = 0;
dd_devsw.d_select = nodev;
dd_devsw.d_config = dd_config;
dd_devsw.d_print = nodev;
dd_devsw.d_dump = nodev;
dd_devsw.d_mpx = nodev;
dd_devsw.d_revoke = nodev;
dd_devsw.d_dsdptr = NULL;
dd_devsw.d_selptr = NULL;
dd_devsw.d_opts = DEV_MPSAFE; /* register as MP SAFE */
rc = devswadd(devno,&dd_devsw); /* devsw table entry */
...
} /* end case CFG_INIT */
break;
case CFG_TERM:
{
...
lock_free(&dd_lock); /* release MP lock */
...
} /* end case CFG_TERM */
break;

} /* end switch (cmd) */
...
unlockl(&config_lock); /* unlock dd_config lock */
return(SUCCESS);
} /* end dd_config() */
/* Skeleton code sample using the following lock related */
/* kernel services: */
/* */
/* simple_lock() - set simple lock */
/* simple_unlock() - unlock simple lock */
/* */
/* (thread -- thread process serialization) */
/* */
int
dd_read (dev_t devno, struct uio *uiop)
{ /* start dd_read() */
simple_lock(&dd_lock); /* set simple lock */
... /* (critical section) */

```

```

    simple_unlock(&dd_lock); /* unlock simple lock          */
    return(SUCCESS);
} /* end dd_read */
/* Skeleton code sample using the following lock related  */
/* kernel services:                                     */
/*                                                       */
/* disable_lock() - disable interrupts and set simple lock */
/* unlock_enable() - unlock simple lock and enable interrupts*/
/*                                                       */
/* (thread -- interrupt process serialization)          */
/*                                                       */
int
dd_intr (void)
{ /* start dd_intr() */
    int old_level; /* interrupt priority level save word */
    ...
    old_level = disable_lock(CURR_LEVEL, /* disable interrupts & */
                             &dd_lock); /* set simple spin lock */
    ... /* (critical section) */
    unlock_enable(old_level, /* unlock simple lock & */
                  &dd_lock); /* enable interrupts */
    return(SUCCESS);
} /* end dd_intr */

```

MP-Efficient Coding Sample

```
/* Skeleton code sample using the following lock related      */
/* kernel services:                                          */
/*                                                          */
/* lockl()           - conventional-lock lock request        */
/* lock_alloc()      - allocate Complex lock                 */
/* lock_init()       - initialize Complex lock               */
/* lock_free()       - release Complex lock                  */
/* unlockl()         - conventional-lock unlock request      */
/*                                                          */
/* (driver configuration lock definition and use)            */
...
lock_t config_lock = {LOCK_AVAIL}; /* dd_config lock        */
struct { /* Adapter Control Block structure                */
    Complex_lock adapter_lock; /* R/W Adapter lock         */
    struct Device_Ctl *next_device /* device structure chain*/
} Adapter_Ctl;
struct { /* Device Control Block structure                */
    Complex_lock device_lock; /* W/Exclusive Device lock*/
    struct statistics device_stats; /* performance statistics */
    short dev_minor_no; /* Device minor number    */
} Device_Ctl;
int
dd_config(int cmd,
          struct uio *uiop)
{ /* start dd_config() */
    /*
     * Use lockl operation to serialize the execution of the
     * config commands.
     */
    if ((rc = lockl(&config_lock, LOCK_SHORT)) != LOCK_SUCC) {
        return(EBUSY);
    }
    switch(cmd) {
        case CFG_INIT:
            {
                ...
                /*
                 * Define the locks in the adapter/device blocks
                 */
                lock_alloc(&Adapter_Ctl.adapter_lock,
                           LOCK_ALLOC_PIN,
                           ADAPTER_CTL_LOCK,
                           -1);
                lock_init(&Adapter_Ctl.adapter_lock,
                          TRUE);
                lock_alloc(&Device_Ctl.device_lock,
                           LOCK_ALLOC_PIN,
                           DEVICE_CTL_LOCK,
                           Device_Ctl.dev_minor_no);
                lock_init(&Adapter_Ctl.adapter_lock,
                          TRUE);
                ...
            } /* end case CFG_INIT */
            break;
        case CFG_TERM:
            {
                /*
                 * Free the locks in the adapter/device control blocks

```



```

        */
        ...
        lock_free(&Adapter_Ctl.adapter_lock);
        lock_free(&Device_Ctl.device_lock);
        ...
    } /* end case CFG_TERM */
    break;
    ...
}
unlockl(&config_lock); /* unlock dd_config lock */
return (SUCCESS);
} /* end dd_config() */
/* Skeleton code sample using the following lock related */
/* kernel services: */
/* */
/* lock_read() - lock Complex lock for shared read */
/* - or - write exclusive access */
/* lock_write() - lock Complex lock for w/exclusive access */
/* lock_done() - release Complex lock */
/* */
/* (thread - thread process serialization only) */
/* */
int
dd_ioctl(dev_t dev, int cmd, caddr_t arg, uint flag,
        chan_t chan, caddr_t ext)
{
    int rc; /* return code */
    ...
    /*
     * Lock adapter lock in shared read access (this is done to
     * protect against removal of device structure)
     */
    lock_read(&Adapter_Ctl.adapter_lock);
    ...
    switch(cmd) {
    /*
     * Return current device performance statistics
     */
        case DD_GET_STATS:
        {
            ...
            /*
             * Lock device control block for read access to gather
             * statistics for caller
             */
            lock_read(&Device_Ctl.device_lock);
            ... /* (critical section) */
            lock_done(&Device_Ctl.device_lock);
            rc = (SUCCESS);
            break;
        }
    /*
     * Update device performance statistics
     */
        case DD_UPDATE_STATS:
        {
            ...
            /*
             * Lock device control block for write access to update

```

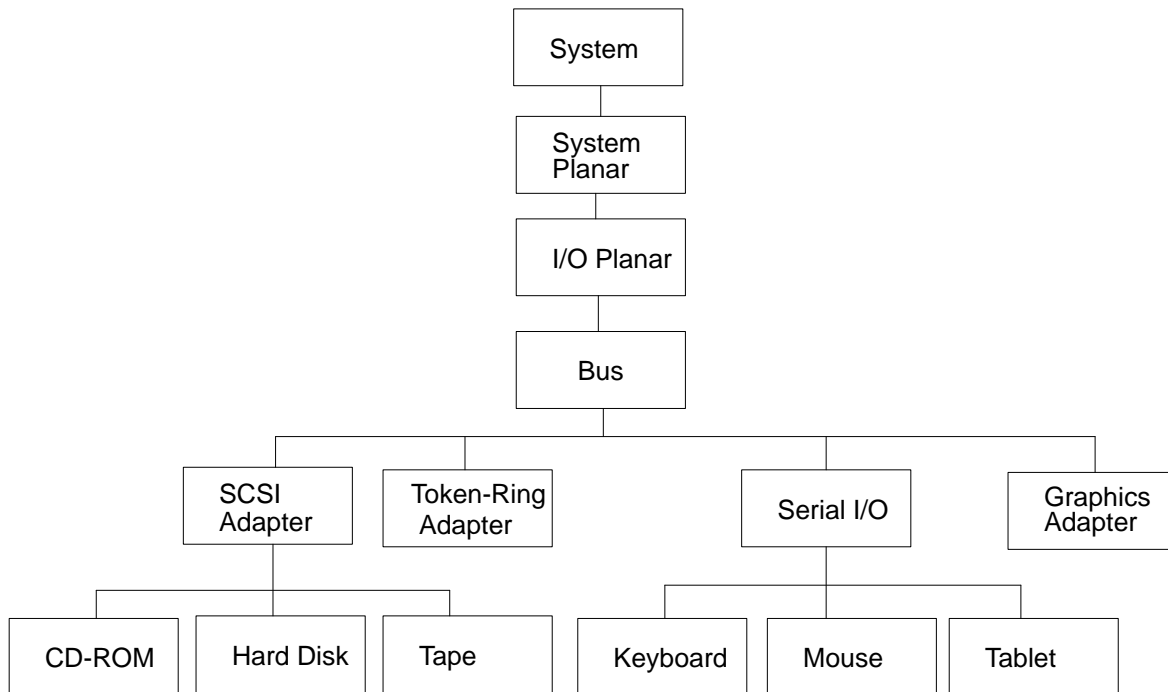
```
        * statistics for caller
        */
        lock_write(&Device_Ctl.device_lock);
        ... /* (critical section) */
        lock_done(&Device_Ctl.device_lock);
        rc = (SUCCESS);
        break;
    }
    ...
}
lock_done(&Adapter_Ctl.adapter_lock);
return(rc);
}
```

Chapter 6. Device Configuration Methods

The dynamically loadable and unloadable aspect of the AIX Version 4.1 kernel requires that all device drivers have configuration methods to support the ability to load and unload them from the kernel. *Configuration methods* are sets of executables including a Define method, a Configure method, a Change method, an Unconfigure method, and an Undefine method. A *Configure method* is part of a set of *configuration methods*. Be careful to not confuse these terms. AIX Version 4.1 also has provisions for a Stop method and a Start method but these are seldom needed in the case of device drivers.

The System Device Hierarchy figure shows how relationships of devices to other devices on the system can be seen as a tree of parent-child relationships. Parents detect their children and then execute the appropriate configuration methods to introduce them to the system. These methods will then load the appropriate device drivers and make the device available for use.

The system's Configuration Manager (started by the **cfgmgr** command) actually oversees the entire configuration process by starting configuration methods, interpreting errors, and managing the configuration of child devices.



System Device Hierarchy

Device States

The system considers each device to be in one of several different *states* indicating the device's availability for use. The state of each device is stored in a database by the Object Data Manager (ODM). There are several databases used to maintain the configuration data for each device. The following is a list of states and their meanings to the system.

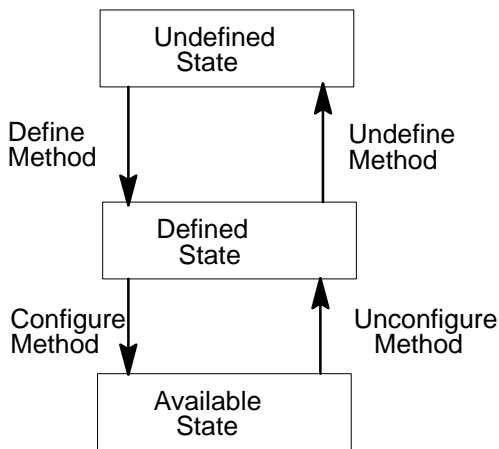
Undefined The device is not known to the system. There is no information about the device in the ODM database

Defined The device is known to the system but currently does not have its device driver loaded and is therefore not available for use

Available The device has its device driver loaded and is available for use

The various configuration methods will take a device from one state to another by detecting any attached devices, manipulating the databases, and loading or unloading the appropriate device driver. The names of a device's configuration methods are stored in the device's entry in the Predefined Device (PdDv) database.

The Device States and Methods figure shows how various methods change the state of a device.



Device States and Methods

ODM Configuration Databases

The databases used by the ODM (Object Data Manager) include the Predefined Devices (PdDv), Predefined Attributes (PdAt), Predefined Connections (PdCn), Customized Devices (CuDv), Customized Attributes (CuAt), Customized Dependencies (CuDep), Customized Device Drivers (CuDvDr), Customized Vital Product Data (CuVPD), and Config Rules (Config_Rules) object classes. The actual database files exist in the directories **/etc/objrepos** and **/usr/lib/objrepos**. The databases keep track of the devices defined to the system, the relationships among devices, and the various attribute values used to configure each device. The purpose of the various object classes is explained in the following list:

PdDv Contains definitions for the entire set of devices that could be supported by the system

PdAt Contains attributes for the devices listed in PdDv

PdCn Contains the supported connection points for parent-child devices

CuDv Contains definitions for devices that have been introduced to the system and may contain devices that are in either the DEFINED or AVAILABLE state

CuAt Contains the attributes required by the devices in CuDv

CuDep	Contains devices that are dependent on other devices
CuDvDr	Contains the major numbers of drivers loaded into the kernel as well as the device number (major, minor) of each device
CuVPD	Contains the Vital Product Data of a device if it contains such data
Config_Rules	Contains a list and order of config methods to be run

Several commands and system calls exist to support manipulation of the various configuration databases. The commands include **odmget**, **odmadd**, **odmdelete**, **odmchange**, and **odmshow**. For more information on these commands, see the *AIX Version 4.1 Commands Reference*. For a detailed description of the individual fields for the object classes, see reference information for the object classes in *AIX Version 4.1 Technical Reference, Volume 6: Kernel and Subsystems*.

The **odmshow** command is useful for displaying the format of the various databases.

The **odmget** command is useful for extracting current entries in each of the databases. These examples can then be used as templates for new entries that can be added with the **odmadd** command. When writing a new device driver package, new entries only need to be written for the PdDv, PdAt, PdCn, and possibly Config_Rules object classes.

The **odmadd** and **odmdelete** commands are useful for adding or deleting any new entries.

The system calls available for manipulating the configuration databases, usually called from within configuration methods, include **odm_initialize**, **odm_add_obj**, **odm_rm_obj**, **odm_change_obj**, **odm_get_first**, **odm_get_obj**, and **odm_terminate**. For a detailed explanation of these and other system calls designed for manipulating the ODM databases, see *AIX Version 4.1 Technical Reference, Volume 2: Base Operating System and Extensions*. The various system calls will typically be used by the various configuration methods to view, add, change, and remove objects from the configuration databases.

Define Methods

Devices are first introduced to the system via a define method. A generic define method (**/usr/lib/methods/define**) on the system can be used for many devices. This method is designed to cover a wide variety of devices, but if a specific device requires special processing, a new specific define method will have to be written for it.

In general, a define method will take as parameters, the name of the new device, its *class*, *subclass*, and *type*, the name of its parent, and the new device's connection point. The define method will then verify these parameters and create a new entry in the CuDv object class for the device using the **odm_add_obj** subroutine. At this point, the device is only marked as Defined and the driver has not yet been loaded. A device's define method is typically invoked from the command line via the **mkdev** command. This command will actually look up the name of the define method to run from the PdDv database for the device being defined. The following is a description of the various flags for the generic define method:

-c class	String that designates the class of the device
-s subclass	String that designates the subclass of the device
-t type	String that designates the type of the device
-p parent	String that designates the logical name of the device's parent
-w connection	String that designates the connection point
-l name	String that designates the new logical name of the device

- u** Indicates that the **-l** flag is not allowed (the define method will generate a new logical name based on the device's prefix stored in PdDv)
- n** Indicates that the device has no parent or parent connection
- o** Indicates that only one of these types of devices can exist and that the **-l** flag is not allowed
- k** Indicates that the device can only be defined if one does not already exist at the specified connection point

For a more detailed discussion on writing a define method, see the "Device Configuration Subsystem" chapter in *AIX Version 4.1 Kernel Extensions and Device Support Programming Concepts*. For more information on the **odm_add_obj** subroutine, see *AIX Version 4.1 Technical Reference, Volume 2: Base Operating System and Extensions*

Configure Methods

Config methods are designed to resolve a device's attributes, build device-dependent structures, find any child devices, and load a device's driver. A device may have attributes that could conflict with other devices that already exist on the system. For example, an adapter's bus-memory address cannot conflict or overlap the address of another adapter. It is the duty of the configure method to ensure that this does not occur. The **busresolve** routine (discussed in "Adapter Device Attributes and busresolve," on page 6-7), is a system call provided to ensure that bus attributes do not conflict.

A device dependent structure (DDS) is used by the configure method to pass vital information down to the device driver. This structure may contain information such as the bus slot location, DMA level, and interrupt level of an adapter, or various attributes needed to customize operation of the device. The DDS is built by reading in information from the various databases. For example, the slot number can be read from the `connwhere` value stored in CuDv. Various attributes required by the device driver can also be read from either PdAt or CuAt using the **getattr** subroutine. Because CuAt contains non-default attribute values, the **getattr** subroutine queries it first before retrieving values from PdAt.

The DDS structure is typically passed in to the driver through the driver's `dd_config` routine. After using the **loadext** routine to load the driver, the configure method should build the DDS structure and pass it to the driver using the **sysconfig** subroutine call and a command of `CFG_INIT`. The following is an example:

```
sysconfig (SYS_CFGDD, &dds, sizeof(dds));
```

The **sysconfig** routine will call the driver's **dd_config** routine with the `CFG_INIT` command and pass in the DDS structure.

For more information on the **sysconfig** routine, see *AIX Version 4.1 Technical Reference, Volume 2: Base Operating System and Extensions*. For more information on the **loadext** subroutine, see *AIX Version 4.1 Technical Reference, Volume 6: Kernel and Subsystems*

If the device can have child devices connected to it, the configure method must detect them and run the child device's define method in order for the child to be added to the CuDv database. Once the define method has been successfully run, using the **odm_run_method** subroutine, the logical name of the newly created child should be printed to **stdout** so that the Configuration Manager can catch it and run the child's configure method.

Once the configure method has successfully loaded and configured the device driver, it should change the state of the device in the CuDv database to available using the **odm_change_obj** subroutine.

For a detailed discussion of configure methods, see the “Device Configuration Subsystem” chapter in *AIX Version 4.1 Kernel Extensions and Device Support Programming Concepts*. For reference information on the **odm_change_obj** subroutine, see *AIX Version 4.1 Technical Reference, Volume 2: Base Operating System and Extensions*.

Change Methods

Change methods are designed to modify a device’s attributes to values other than the default. They can also alternately be designed to relocate a device to a different location or parent device. A generic change method **/usr/lib/methods/chgdevice** is provided on the system. A device’s default attribute values are stored in the Predefined Attribute (PdAt) object class. The Customized Attribute (CuAt) object class is used to store values that are different from the default.

A change method should first verify that the attributes being modified are being changed to valid values. If the device is currently configured (available), it should be unconfigured by using the **odm_run_method** subroutine to execute the unconfigure method. This is required so that the device driver can later be initialized again with the new attribute values.

If an attribute is being changed to a value other than the default as listed in PdAt, a new attribute value should be added in CuAt. If a CuAt value already exists, it can be changed to the new value with the **odm_change_obj** subroutine.

If an attribute is being changed back to the default PdAt value, the CuAt entry should be deleted.

Once all the attributes have been changed, the device driver’s configure method should be executed using the **odm_run_method** subroutine so that the device driver can be reloaded with the new attribute values. The following is a list of the various parameters for the generic change method:

- l name** String that designates the logical name of the device
- p parent** String that designates the logical name of the device’s new parent
- w connection** String that designates the device’s new connection point
- a attr=val, attr=val**
Strings that designate the attribute name and new value

For a more detailed discussion on writing a change method, see the “Device Configuration Subsystem” chapter in *AIX Version 4.1 Kernel Extensions and Device Support Programming Concepts*. For reference information on the **odm_run_method** and **odm_change_obj** subroutines, see *AIX Version 4.1 Technical Reference, Volume 2: Base Operating System and Extensions*.

Unconfigure Methods

Unconfigure methods merely undo what a configure method has done. As in the case of the define and change methods, the generic unconfigure method (**/usr/lib/methods/ucfgdevice**) that has been provided should be satisfactory for many devices. However, if a specific device requires additional functionality, a new unconfigure method will have to be written.

In general, an unconfigure method should first verify that the device is in the correct state (AVAILABLE). It should then verify that any children of the device are also in the correct state (DEFINED). The unconfigure method should then call **sysconfig** with a command of CFG_TERM to instruct the device driver to terminate (the **sysconfig** routine will call the

driver's **dd_config** routine with the CFG_TERM command). Once the device driver has successfully terminated, the driver can be removed with the **loadext** system call. As a final step, the unconfigure method should change the state of the device in CuDv from AVAILABLE to DEFINED. The following is a list of the various parameters for the generic unconfigure method:

-l <name> String that designates the logical name of the device

For more information on the **sysconfig** system call, please refer to the book *AIX Version 4.1 Technical Reference, Volume 2: Base Operating System and Extensions*. For reference information on the **loadext** subroutine, see *AIX Version 4.1 Technical Reference, Volume 6: Kernel and Subsystem*. For a more detailed discussion on writing an unconfigure method, see the "Device Configuration Subsystem" chapter in *AIX Version 4.1 Kernel Extensions and Device Support Programming Concepts*.

Undefine Methods

Undefine methods will take a device from the defined state to the undefined state by removing its entry from the CuDv database. In this state, all previous information about where the device was connected and what logical name was attached to it is removed. Any customized attributes for the device are also lost. As in the case of the change, define, and unconfigure methods, a generic define method (**/usr/lib/methods/undefine**), which should be satisfactory for many devices, has been provided.

In general, an undefine method should first verify that the device is in the correct state (DEFINED). It should then verify that child devices, if any, have been undefined (they don't exist in the CuDv database). The undefine method should then make sure that this device is not dependent on any other device, by checking the CuDep database, and removing all of its customized attributes from CuAt. The final step the undefine method should perform is to release the device number (devno) and remove any CuDvDr and CuDv entries. The following is a list of the various parameters for the generic undefine method:

-l <name> String that designates the logical name of the device

For a more detailed discussion on writing an undefine method, see the "Device Configuration Subsystem" chapter in *AIX Version 4.1 Kernel Extensions and Device Support Programming Concepts*.

Configuring Devices with No Parent

The configuration hierarchy is designed so that many devices are actually *children* of other *parent* devices. These child devices are automatically detected and defined by their parent configure methods. There are cases, however, of devices that are not connected to any parent devices. One example is the system node, **sys0**, which is a device that has no parent but is instead the parent and grandparent to most of the devices on the system. Many third-party vendor devices fall into this category of parent-less devices. In this case, these devices cannot rely on a parent because they are new to the system and the current OS either cannot detect them or will not recognize them.

When devices cannot rely on a parent to run their define method, they must use the Config_Rules database to initiate execution of their define method. Very often this define method will also have to detect the device being defined, which requires writing and using a new custom define method in place of the generic method. This is necessary because no parent exists or no parent has been written to perform the normal detection. The following is an example of a Config_Rules entry:

```

phase= 2
seq   = 50
boot_mask = 0
rule = "/usr/lib/methods/defdevice"

```

In the previous example, the define method `/usr/lib/methods/defdevice`, would be run during phase 2 with a sequence number of 50. Phase 1 is used to configure basic essential devices needed to boot the system. Phase 2 is used to boot most base system devices. Note that a phase of 3 indicates methods to be run in phase 2 service mode. A high sequence number was chosen in case the device depended on other devices to be available. A high value causes the method to be run late in the configuration process. A nonzero boot mask, as defined in `/usr/include/sys/cfgdb.h`, indicates the type of boot (for example, disk, diskette, tape, or network) to which the method applies.

For more information on the `Config_Rules` object class, see *AIX Version 4.1 Technical Reference, Volume 6: Kernel and Subsystems*

Adapter Device Attributes and busresolve

When configuring a device, care must be taken not to configure a device in a manner that conflicts with an existing device on the system. For example, a SCSI device cannot be configured with the same ID and LUN (logical unit number) as another device. In the case of adapters, not only must their connection locations (slot numbers) be unique, but several of their attributes must also be unique. For example, one adapter's allocated bus memory range must not overlap another's.

The **busresolve** routine is available to detect and resolve any possible conflicts with bus resource attributes. It is usually called from within an adapter's configure method. This routine needs to be called only if the configure method is being executed at run time. The bus configure method calls **busresolve** at boot time to properly resolve the attributes for all adapter devices.

busresolve will scan the CuAt and PdAt databases for all attributes with types that designate them as bus resource attributes. The following is a list of different bus attribute types:

Type	Description
O	Indicates an address and width for bus I/O
M	Indicates an address and range for DMA transfers
B	Indicates an address and range for non-DMA transfers
A	Indicates a DMA arbitration level
I	Indicates a sharable interrupt level
N	Indicates a non-sharable interrupt level
P	Indicates an interrupt priority class
W	Used to specify bus I/O and memory widths if not already specified by the address attribute
G	Indicates a bus resource that must be assigned as part of a group
S	Indicates a bus resource that must be shared with another adapter

busresolve will adjust the attribute values within each attribute's specified allowable ranges. **busresolve** never adjusts the attributes of an AVAILABLE device. It only adjusts the attributes of DEFINED devices. This implies that devices configured first will have a greater

probability of obtaining their default attribute values. When a configure method calls **busresolve** during run time, it should pass in the device's logical name as the *logname* parameter. This ensures that **busresolve** will adjust attributes only for the specified DEFINED device. Passing in NULL for *logname* causes **busresolve** to attempt resolution of attribute values for all Defined devices. This is typically done by the bus configuration manager at boot time and, therefore, does not need to be done by each individual adapter configure method. The CuAt and PdAt databases are automatically updated by **busresolve** once all attributes have been resolved.

For reference information on **busresolve**, see *AIX Version 4.1 Technical Reference, Volume 6: Kernel and Subsystems*

Configuration of Devices on PCI and ISA Bus Systems

The configuration process for adapters attached to systems that contain PCI and ISA buses remains similar to the process for systems with a Micro Channel bus. The same ODM and configuration methods scheme is used with one additional step required from the system administrator when configuring adapters attached to the ISA bus. The PCI and ISA bus configuration also differs from the Micro Channel bus in that the ISA is actually treated as a child of the PCI bus even though the ISA bus has adapter children of its own. During configuration of the PCI bus, the ISA bus is detected and then a separate configure method is called for the ISA bus.

Because ISA adapters do not implement the notion of POS registers as in the case of Micro Channel adapters, an additional step must be performed by the system administrator. In the case of PCI and Micro Channel adapters, through the machine device driver, the bus configure method can detect the type of adapter attached and then dynamically set certain parameters such as the interrupt level or DMA arbitration level. In the case of ISA adapters, these types of parameters are set via jumpers or DIP switches that physically reside on the adapter. The system administrator installing the card must first set these parameters to avoid conflicts with other adapters before physically installing the adapter into the system. Since these settings cannot be detected on the ISA bus, the system administrator must manually add the CuDv (Customized Device) and CuAt (Customized Attribute) entries for the adapter. The CuDv entries should reflect that the adapter is a child of the ISA bus and its location code must correspond to the correct slot in which it resides. The CuAt entries must also reflect the actual settings on the adapter itself and must also not conflict with other adapters on the ISA bus. The **generic** field of the CuAt and PdAt entries should also contain a second character **U** to indicate that the particular attribute value is user modified. Newer ISA adapters that provide dynamic detection of modification and card settings (plug-and-play ISA adapters) are currently not supported.

Once the ISA specific requirements have been met, the configuration process of the ISA and PCI busses is identical to that of the Micro Channel bus as previously discussed. The **busresolve** routine is called by the corresponding bus configure method during boot so that adapter configure methods are not required to call it. The only time an adapter configure method should call **busresolve** is during a runtime configure when the current attributes of the adapter being configured must be validated against values already resolved during boot.

Chapter 7. Block Device Drivers

A block device driver interacts with a special facility in the kernel called the *buffer cache*. Special entry points in the driver are provided because of this interaction. A block device driver may also support character type interaction through read and write operations referred to as *raw I/O*. The principal characteristic of block devices is to perform I/O operations using system facilities such as buffer cache management and paging.

Introduction

Data read from character devices is *not* stored in a cache for subsequent reading from system buffers. For block device drivers, data *is* stored in a cache. Block devices interact with the system to keep the cache containing information that a process (or multiple processes) can read from at any time. If the information is not in the cache, the system (not the user) requests the data from the block device driver.

Like all devices, the interaction with block devices is through shared memory. In addition, there are routines to indicate when data in the shared memory (called **buf** structures) has completed I/O processing.

The following sections cover the entry points responsible for the movement of data to and from block devices. This includes control information, the shared memory facilities, the mechanisms for programs to share the information, and the use of the kernel cache.

Finally, it may be necessary to talk to the device directly without interacting with the system buffer cache. This topic is presented in “Character Access to Block Device Drivers” on page 7-6.

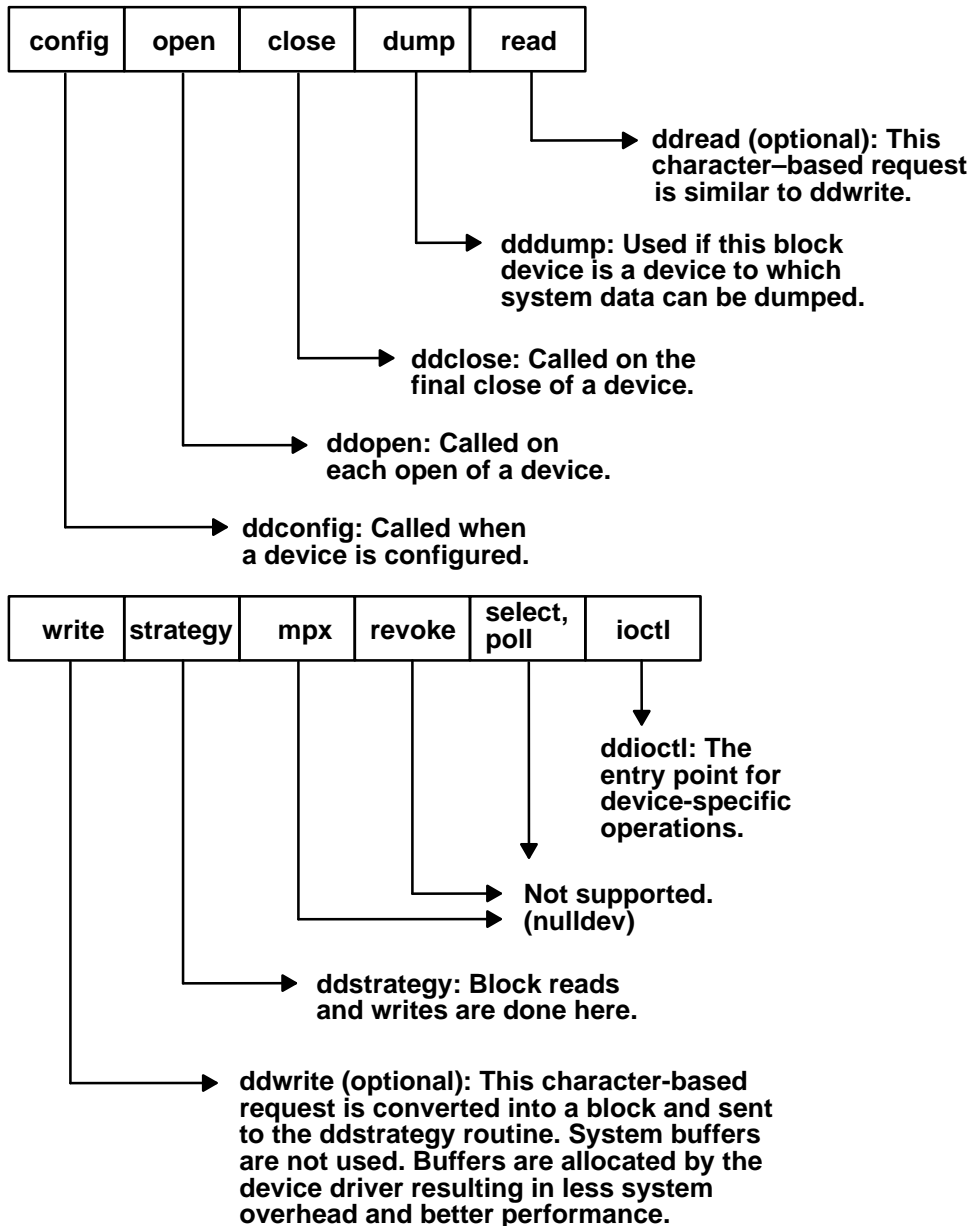
Block I/O Device Driver Entry Points

The device switch table contains the entry point addresses of the interface routines for each device driver in the system, just as it does for the character device drivers. The following Entry Points for the Block Device Driver figure shows the entry points for a block device driver.

Like the character device driver, the block device driver must supply a **config** routine for configuration support as well as an **open** and a **close** routine. The **open** routine is called each time the device is opened and the **close** routine is called only on the final close of the device.

Instead of having separate **read** and **write** routines, like character device drivers, each block device driver has a **strategy** routine. This routine is called with a pointer to a buffer header, known as the **buf** structure, which contains the I/O request parameter.

The **strategy** routine handles requests as buffers to be written or read from the device.



Entry Points for a Block Device Driver

ddconfig Entry Point

The configuration routine of a block device driver creates an entry in the device switch table for the block device driver. The device may support raw access. Raw access is character access to a block device. In this case, the driver's configure method must have created **/dev** special file entries for use in raw access. These special files retain the same major and minor numbers as their corresponding block device special files, but they have the letter *r* as a prefix, and the special files are created as character rather than as block.

For example, a block device named **/dev/hdisk0** that supports raw access also has a **/dev/rhdisk0** special file. The system calls the **read** and **write** routines of the raw device if

`/dev/rhdisk0` is opened. Other UNIX systems may not allocate the same major and minor numbers for both character and block devices.

ddopen and ddclose Entry Points

The AIX operating system supports only a few block devices in normal installation. These devices, such as hard disks and CD-ROM, are capable of random access and are opened by system services such as the buffer cache and paging subsystem. They are not to be opened directly by user space applications during normal system operations, but may be opened during maintenance by applications such as **fsck**.

The **ddopen** routine verifies that the device is a valid device and that it is online and available. It also performs any setup required by the driver, such as pinning of code and allocation of data structures.

Most of the block devices are attached to the SCSI adapter, and open the SCSI adapter device driver to communicate with the device. For more information on the SCSI subsystem, see the chapter on “Small Computer System Interface (SCSI) Subsystem” in *AIX Version 4.1 Kernel Extensions and Device Support Programming Concepts* and the chapter on SCSI Device Drivers in this book.

The device performs **ddclose** processing to release the resource. If the device is attached to the SCSI bus, refer to the *AIX Version 4.1 Kernel Extensions and Device Support Programming Concepts* book and the chapter on SCSI Device Drivers in this book for details.

ddstrategy Entry Point

The I/O requests to the physical device are accomplished through the **strategy** routine. The **strategy** routine provides a “strategy” for mapping I/O requests to the device so that it minimizes requests to the device and maximizes data transfer. When the **strategy** routine (**ddstrategy** device driver entry point) is called, a pointer to a buffer header or a chain of buffer headers specifies the request for device I/O. The **strategy** entry point is called in a user process context when the buffer cache does not contain the buffer requested by the user. The **strategy** routine, however, does not know about the user process.

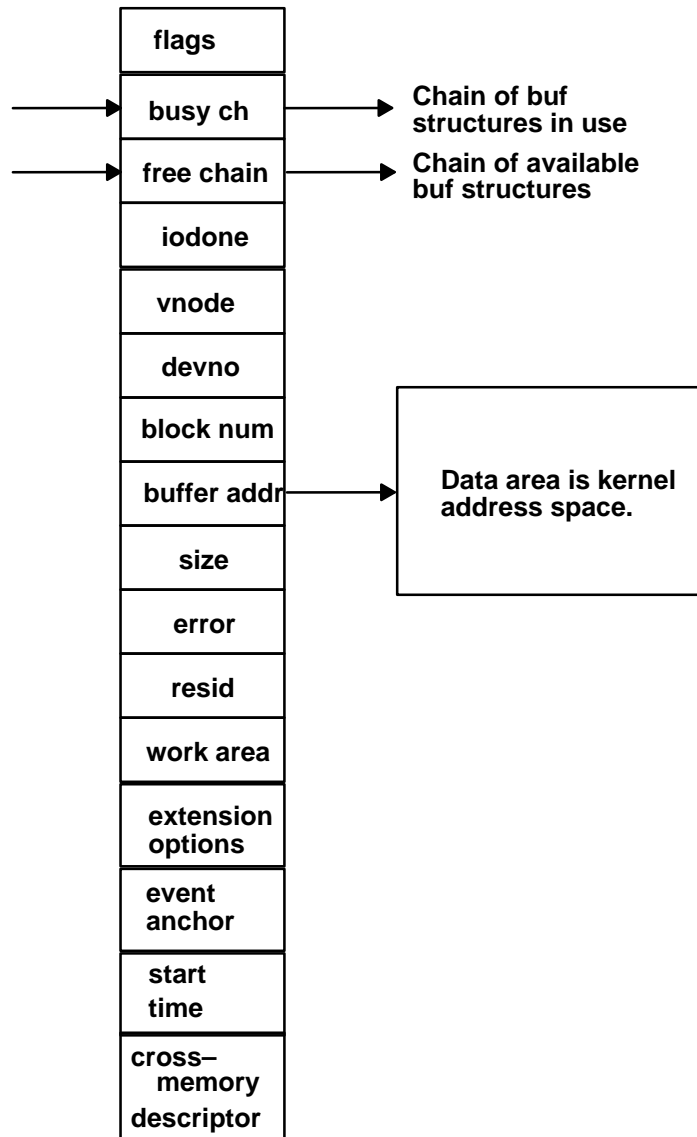
The buffer header contains the following information:

- The major and the minor numbers of the device.
- The description of the memory buffer to be used for the data transfer.
- The direction of the transfer.
- The transfer count.
- The block number on the device for which the transfer is targeted.
- The operation flag.

The **strategy** routine returns to the caller as soon as the buffer headers are queued to the appropriate device queue. Note that the **strategy** routine provides no return code to the caller and never waits for I/O completion before returning. This means that all requests are assumed valid in terms of parameters and that the request is asynchronous. Normal errors, such as out-of-range blocks, are caught but not returned directly as a return code.

The execution of the request completes some time later. The buffer structure contains fields for reporting the completion of the request.

A header contains all the information required to perform block I/O. The **buf** structure is shown in the `buf` Structure figure. It is the primary interface to the bottom half of block device drivers.



buf Structure

In AIX, the traditional **strategy** interface is extended as follows:

- The device driver **strategy** routine is called with a list of **buf** structures, chained using the *av_forw* and *av_back* pointers. The last entry in this list has a NULL *av_forw* pointer.
- When the operation is completed and the driver calls the **iodone** kernel service, the **b_iodone** function defined by the caller is scheduled to run as a software interrupt handler.

The **buf** structure and its associated data page must be pinned before calling the **strategy** routine. This is outlined in the `/usr/include/sys/buf.h` include file.

The **buf** structure contains the operation to be performed and status information to be returned to the caller, and is more like a message exchanged between a requestor and a service provider.

The caller is notified of I/O completion (or of an error associated with the request) by the device driver's call to the **iodone** kernel services. A residual count of the number of bytes requested but not transferred by the operation is placed in the *b_resid* field of the **buf**

structure by the device driver before the I/O is marked as complete for the buffer header. If all the requested bytes are transferred, this count is set to zero.

For more information on the specific fields of the **buf** structure, see “buf Structure” in *AIX Version 4.1 Technical Reference, Volume 5: Kernel and Subsystems*

Note: In AIX Version 4.1 a new flag, B_MPSAFE, has been added to the list of valid flags for the **b_flags** field of the **buf** structure definition. This flag is to indicate whether or not it is safe for the **iodone** processing to be performed on one processor or multiple processors. The **devstrat** kernel service will actually mark the **buf** structure with the B_MPSAFE_INITIAL flag once the B_MPSAFE flag has been sent. A device driver’s strategy routine, however, does not have to be concerned with this flag, since it is merely an instruction to the **iodone** system call about how to complete the processing of the **buf** structure.

Reordering Block I/O Requests

Multiple I/O requests can also be presented to the **strategy** routine, where the additional buffer headers can be chained to the first by using the *av_forw* pointers. While the device **strategy** routine is free to rearrange the buffers on its device queue with respect to the processing of single request, the ordering of the buffer headers provided in a chain to the **strategy** routine cannot be modified. The **strategy** routine also determines if the block number requested is valid for the device. In the case of a read-only operation, a block number at the end-of-media is not considered as an error, but no data is transferred.

For a write operation, if the block number is at the end-of-media, it is considered an error, the **B_ERROR** flag in the **buf** structure is set, and the *b_error* field contains the **ENXIO** value.

Categorizing Requests to the Start I/O Routine

To maintain the state of the device and its I/O requests, the device driver typically allocates a private data structure in the system memory associated with the device. The data structure contains the device status along with the device error information and the device queue pointers. Some device drivers maintain more than one queue of buffer headers. For example, one queue for the requests that are waiting for I/O start and another queue for the request that are currently in process.

For SCSI operations, the queueing process scans the pending queue for the requested device so that the number of SCSI operations is minimized. The requests are grouped by one of the following rules:

- Contiguous write operations
- Operations larger than maximum transfer size
- Operations requiring special processing

The coalesced (grouped) requests will be removed from the pending queue and placed in the **in_progress** queue so that a single command can be built to satisfy the requests. These requests are then queued and the routine to start I/O is called.

Starting Processing with the Start I/O Routine

The routine to start I/O checks to make sure that the device is not busy, and then scans the request queues in an attempt to find an operation to start.

First, the command stack is checked to see if a command needs to be restarted. Then the **in_progress** queue is checked to start any operations that have already been coalesced. Finally, the **pending** queue is checked. If it is not empty, the **coalesce** routine is called to group the operations into the **in_progress** queue.

When a request has been found and built, the adapter device driver is called by the strategy routine to begin processing the operation. While the queues are being scanned and an operation is in progress, the device busy flag is set. It is then reset if no request is found.

Once the I/O handling routine has completed an I/O transfer, it calls the **iodone** kernel service that determines if the indicated operation has completed successfully or if it has failed. If the operation is successful and complete, the next request is processed by the start I/O routine. If the operation has failed, your general failure processing routine is called in an attempt (such as retry) to clear the error.

dddump Entry Point

To support system dumps, a block device driver must supply the **dddump** entry point. It is called by the **devdump** kernel service. See Chapter 14, “Debugging Tools,” for more information on system dumps. See Chapter 8, “SCSI Device Drivers,” for more information on providing system dump support on SCSI devices.

Character Access to Block Device Drivers

While a character device driver can only be accessed by a character special file, most block device drivers provide both a block and a character special file. With this dual interface, a user can access the device in either block or character mode.

Note that the block device driver must have a **read** and a **write** entry point as well as a **strategy** entry point if it supports both character and block mode access. If it supports only block mode, it only needs to support a **strategy** entry point.

The diskette or hard disk device drivers are examples of the dual nature of block device drivers. The diskette is accessed by **/dev/fd0** for block mode and by **/dev/rfd0** for character (raw) mode. The hard disk is accessed by **/dev/hdisk0** for block mode and **/dev/rhdisk0** for character (raw) mode.

Raw I/O Processing

Raw I/O processing is a mechanism by which a block device driver has the ability to transfer data without using the I/O buffer cache. The raw I/O request is converted into a block and then sent to the device driver **strategy** entry point to be processed while the **read** and the **write** routines are typically waiting for the I/O completion.

When your device driver is configured, it contains entries for both **read** and **write** (raw access) and **strategy** (block access) routines. In addition, the configure method must set up the **/dev** entries for both special files.

If there is no buffer cache and you make the request directly, a different buffering facility is involved. You are providing a buffer passed in through the uio services. Therefore the **read** and **write** entry points are talking to a user process and translating the requests into **strategy** requests but still using **buf** structures. Because the **buf** structure contains a header that contains a pointer to the data area, it can be mapped to point to a user data area.

In fact, the user buffer can come out of the user data, text segments, shared memory segments, or the system segment. The different areas are defined in the **uio** and **iovec** structures.

The **read** and **write** routines of the raw device driver use the **uphysio** kernel service to map the **uio** areas into **buf** structures used by the strategy routines. After filling in the **buf** structure with data passed to it through the **uio** structure, the **uphysio** kernel service will call the block device driver’s strategy routine. The number of **buf** headers sent is

determined by the `buf_cnt` parameter passed to the **uphysio** kernel service. The **uphysio** kernel service returns only after all I/O has completed or after encountering an error. See "uphysio Kernel Service" in *AIX Version 4.1 Technical Reference, Volume 5: Kernel and Subsystems* for a more detailed discussion of this kernel service.

Note: Use care when accessing a block device through its character interface. Because the buffer cache is bypassed, the driver must be sure that no data currently exists in the kernel buffer cache. If another process did have data present in the cache, there is a high probability of data becoming inconsistent with data obtained through the character interface.

Block I/O Device Summary

A block I/O device contains a device name for its block device and its optional character device. Block devices support **strategy** routines, and possibly support **read** and **write** routines. The kernel cache speeds up access to data by allowing multiple processes to use the same data and keeping data that is referenced often in the cache. However, the cache is based on buffer sizes compatible with UNIX file system block sizes and is not efficient for applications that can use larger block sizes. To improve performance, a block device driver also provides raw or character access to block devices. More information is available in *AIX Version 4.1 Kernel Extensions and Device Support Programming Concepts*.

Chapter 8. SCSI Device Drivers

The AIX Small Computer Systems Interface (SCSI) subsystem has two parts:

- SCSI Device Driver
- SCSI Adapter Device Driver

The SCSI 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 SCSI device driver without having a detailed knowledge of the system hardware. You can look at the SCSI 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 SCSI device driver, because the SCSI adapter device driver is already provided in AIX.

The SCSI adapter device driver, or lower layer, is responsible only for the communications to and from the SCSI bus, and any error logging and recovery. The upper layer is responsible for all of the device-specific commands. The SCSI device driver should handle all commands directed towards its specific device by building the necessary sequence of I/O requests to the SCSI adapter device driver in order to properly communicate with the device.

These I/O requests contain the SCSI commands that are needed by the SCSI device. One important aspect to note is that the SCSI device driver cannot access any of the adapter resources and should never try to pass the SCSI commands directly to the adapter, since it has absolutely no knowledge of the meaning of those commands.

Device Driver Overview

The role of the SCSI device driver is to pass information between the operating system and the SCSI adapter device driver by accepting I/O requests and passing these requests to the SCSI adapter device driver. The device driver should accept either character or block I/O requests, build the necessary SCSI commands, and then issue these commands to the device through the SCSI adapter device driver.

The SCSI 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.

Adapter Device Driver Overview

Unlike most other device drivers, the SCSI 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 SCSI adapter diagnostics.

A SCSI device driver does not need to access the SCSI 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.

SCSI Adapter/Device Interface

The AIX SCSI 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 SCSI 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 SCSI adapter is accessed by the device driver through the **/dev/scsi#** 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.

For additional information on spanned and gathered write commands, see "Small Computer System Interface (SCSI) Subsystem" in *AIX Version 4.1 Kernel Extensions and Device Support Programming Concepts*

sc_buf Structure

The I/O requests made from the SCSI device driver to the SCSI adapter device driver are completed through the use of the **sc_buf** structure, which is defined in the **/usr/include/sys/scsi.h** header file. This structure, which is similar to the **buf** structure in other drivers, is passed between the two SCSI subsystem drivers through the **strategy** routine. The following is a brief description of the fields contained in the **sc_buf** structure:

struct buf bufstruct

This structure is a copy of the standard **buf** structure used for the I/O request that is defined in the **/usr/include/sys/buf.h** header file. Note that the **b_work** field in the **buf** structure is reserved for use by the SCSI adapter device driver.

struct buf *bp Contains a pointer to the original buffer structure used by the process calling the SCSI device driver. It can be a pointer to a list of SCSI spanned data transfer commands or it can contain a value of NULL. A NULL value indicates that no list exists and all required information is contained in the **bufstruct** field (see the preceding paragraph). A non-NULL value also requires the **resvd1** field of the **sc_buf** structure to be NULL.

uint resvd1 This field is usually set to NULL although it can contain a pointer to a **uio** structure which is used in gathered writes.

uint resvd2 Reserved (should be set to zero).

uint resvd3 Reserved (should be set to zero).

uint resvd4 Reserved (should be set to zero).

uint timeout_value

Contains the amount of time, in seconds, to be used in waiting for the completion of the command before it times out. A zero value indicates no timeout should be used.

uchar status_validity

This field can contain the `SC_SCSI_ERROR` bit flag that indicates that the `scsi_status` field return code is valid or the `SC_ADAPTER_ERROR` bit flag which indicates that the `general_card_status` return code is valid. There are three cases when considering the values of the return codes:

- If the `sc_buf.bufstruct.b_flag` field has the `B_ERROR` flag set, then the `status` field contains a valid **errno** value. If the `b_error` field has the value `ENXIO`, then the command needs to be restarted or the SCSI device driver canceled the request. If the `b_error` field has the value `EIO`, then the `status_validity` field indicates which status field, `scsi_status` or `general_card_status`, contains the error.

If the `status_validity` field is zero, examine the `sc_buf.bufstruct.b_resid` field for any possible error. Note that `b_resid` can be nonzero even if no error occurred. Evaluate the nonzero value carefully to ensure that it is a proper error value.

- If the `sc_buf.bufstruct.b_flag` field does not have the `B_ERROR` flag set, then no error is being reported. However, you must still examine the `b_resid` field to determine if an error has actually occurred. If an error has occurred, it is up to the SCSI device driver to recover since device queues are not stopped and future commands can still be sent to the adapter and driver.
- If the `sc_buf.bufstruct.b_flag` field has the `B_ERROR` flag set, then the device queue has been halted. To recover or continue after the error, the `sc_buf.flags` field must have the `SC_RESUME` bit set in the first `sc_buf` structure.

uchar scsi_status

This field is valid whenever the correct bit is set in the `status_validity` field. The `sc_buf.bufstruct.b_error` field should also contain the value `EIO` whenever the `scsi_status` field is valid. The various valid values are shown and defined in the `/usr/include/sys/scsi.h` header file.

uchar general_card_status

This field is valid whenever the correct bit is set in the `status_validity` field. The `sc_buf.bufstruct.b_error` field should also contain the value `EIO` whenever the `scsi_status` field is valid. The various valid values are shown and defined in the `/usr/include/sys/scsi.h` header file.

The `general_card_status` bit is set in the `status_validity` field whenever the SCSI adapter device driver encounters an unrecoverable error. Recovered errors are those that have been corrected and logged by the adapter device driver. The SCSI adapter device driver logs both bus and adapter-related errors.

If an error is detected after a command has reached a device, it is the responsibility of the device driver to attempt recovery and log the error. Of the values shown in `/usr/include/sys/scsi.h`, the device driver should handle:

- `SC_SCSI_BUS_FAULT`
- `SC_CMD_TIMEOUT`
- `SC_NO_DEVICE_RESPONSE`

The `SC SCSI BUS FAULT` error should be handled by the device driver and not the adapter device driver since this can be caused by a protocol or hardware failure.

uchar adap_q_status

This field is used to indicate that the adapter did not clear the device queue on a failure. The adapter will set this field to `SC DID NOT CLEAR Q` to indicate this condition. For example, this flag is returned if a check condition occurs while a command is being queued to the device.

uchar lun

This field should contain the LUN of the target device. Note that if the LUN is greater than 7, this field should contain the LUN value and the `lun` field in the `scsi_cmd` structure should be 0.

uint resvd7

Reserved (should be set to zero).

uchar q_tag_msg

This field is used when the SCSI device supports command tag queueing. It should be set to 0 if the device does not support queueing. See the `/usr/include/sys/scsi.h` file for a list and explanation of the valid values for this field.

uchar flags

This field contains various flags to instruct the adapter driver on how to process the transfer request. See the `/usr/include/sys/scsi.h` file for a list and explanation of the valid values for this field.

Adapter/Device Driver Intercommunication

In a typical request to the device driver, a call is first made to the device driver **strategy** routine, which takes care of any necessary queueing. This **strategy** routine then calls the device driver's **start** routine, which fills in the `sc_buf` structure and calls the adapter device driver's **strategy** routine through the `devstrat` kernel service.

The adapter's **strategy** routine validates all of the information contained in the `sc_buf` structure and also performs any necessary queueing of the transaction request. If no queueing is necessary, the adapter's **start** subroutine is called.

When an interrupt occurs, the SCSI adapter **interrupt** routine fills in the `status_validity` field and the appropriate `scsi_status` or `general_card_status` field of the `sc_buf` structure. The `bufstruct.b_resid` field is also filled in with the value of nontransferred bytes. The adapter's **interrupt** routine then passes this newly filled in `sc_buf` structure to the `iodone` kernel service which then signals the SCSI device driver's `iodone` subroutine. The adapter'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 `sc_buf` structure for any errors and attempt error recovery if necessary. The device driver should then dequeue the `sc_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.

SCSI Adapter Device Driver Routines

This section describes the following routines:

- **config**
- **open**
- **close**
- **openx**
- **strategy**

- **ioctl**

config

The **config** routine performs all of the processing needed to configure, unconfigure, and read Vital Product Data for the SCSI 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 (VPD) for the adapter, which contains information that is used to identify the serial number, change level, or part number of the adapter.

open

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

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

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

The **strategy** routine is the link between the device driver and the SCSI adapter device driver for all normal I/O requests. Whenever the SCSI device driver receives a call, it builds an **sc_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 SCSI commands required to carry out the request. When the command has completed, the SCSI device driver is notified through the **iodone** kernel service.

ioctl

The **ioctl** routine allows various diagnostic and nondiagnostic adapter operations. Operations include the following:

- **IOCINFO**
- **SCIOSTART**
- **SCIOEVENT**
- **SCIOSTOP**

- SCIOINQU
- SCIOSTUNIT
- SCIOTUR
- SCIOREAD
- SCIORESET
- SCIOHALT
- SCIODIAG
- SCIOTRAM
- SCIODNLD
- SCIOSTARTTGT
- SCIOSTOPTGT

SCSI Adapter ioctl Operations

This section describes the following ioctl operations:

- IOCINFO
- SCIOSTART
- SCIOSTOP
- SCIOINQU
- SCIOSTUNIT
- SCIOTUR
- SCIORESET
- SCIOHALT
- SCIODIAG
- SCIOTRAM
- SCIODNLD

IOCINFO

This operation allows a SCSI device driver to obtain important information about a SCSI adapter, including the card's SCSI ID and the maximum data transfer size in bytes. By knowing the maximum data transfer size, a SCSI 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_SCSI**. 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 **scsi** is the structure that applies to the adapter.

For example, the maximum transfer size value is contained in the variable *infostruct.un.scsi.max_transfer* and the card ID is contained in *infostruct.un.scsi.card_scsi_id*.

SCIOSTART

This operation opens a logical path to the SCSI device and causes the SCSI adapter device driver to allocate and initialize all of the data areas needed for the SCSI device. The SCIOSTOP operation should be issued when those data areas are no longer needed. This operation should be issued before any nondiagnostic operation except for IOCINFO. The following is a typical call:

```
rc = fp_ioctl(fp, SCIOSTART, idlun, NULL);
```

where *fp* is a pointer to a file structure and *idlun* is a type int value that contains the SCSI and Logical Unit Number (LUN) ID values of the device to be started. The least significant byte contains the LUN, the next least significant byte contains the SCSI ID, and the most significant two bytes should be set to zero.

A nonzero return value indicates an error has occurred and all operations to this SCSI/LUN pair should cease since the device is either already started or failed the start operation. Possible errno values are **EIO**, **EINVAL**, or **EACCES**.

- 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.
- EACCES** The adapter is not in normal mode.

SCIOSTOP

This operation closes a logical path to the SCSI device and causes the SCSI adapter device driver to deallocate all data areas that were allocated by the SCIOSTART operation. This operation should only be issued after a successful SCIOSTART operation to a device. The following is a typical call:

```
rc = fp_ioctl(fp, SCIOSTOP, idlun, NULL);
```

where *fp* is a pointer to a file structure and *idlun* is a type int value that contains the SCSI and LUN ID values of the device to be stopped. The least significant byte contains the LUN, the next least significant byte contains the SCSI ID, and the upper two bytes should be set to zero. A non-zero return value indicates an error has occurred. Possible errno values are **EIO** and **EINVAL**.

- EIO** An unrecoverable system error has occurred.
- EINVAL** The adapter was not in open mode.

This operation requires **SCIOSTART** to be run first.

SCIOINQU

This operation issues an inquiry command to a SCSI device and is used to aid in SCSI device configuration. The following is a typical call:

```
rc = ioctl(adapter, SCIOINQU, &inquiry_block);
```

where *adapter* is a file descriptor and *inquiry_block* is a **sc_inquiry** structure as defined in the **/usr/include/sys/scsi.h** header file. The SCSI ID and LUN should be placed in the **sc_inquiry** parameter block. The SC_ASYNC flag should not be set on the first call to this operation and should only be set if a bus fault has occurred. Possible errno values are **EIO**, **EFAULT**, **EINVAL**, **EACCES**, **ENOMEM**, **ETIMEDOUT**, **ENODEV**, and **ENOCCONNECT**.

- EIO** A system error has occurred. Consider retrying the operation several times, because another attempt may 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 may be successful.
- ENODEV** The device is not responding. Possibly no LUNs exist on the present SCSI ID.

ENOCCONNECT A bus fault has occurred and the operation should be retried with the SC_ASYNC flag set in the **sc_inquiry** structure. In the case of multiple retries, this flag should be set only on the last retry.

This operation requires **SCIOSTART** to be run first.

SCIOSTUNIT

This operation issues a start unit command to a SCSI device and is used to aid in SCSI device configuration. The following is a typical call:

```
rc = ioctl(adapter, SCIOSTUNIT, &start_block);
```

where *adapter* is a file descriptor and *start_block* is a **sc_startunit** structure as defined in the **/usr/include/sys/scsi.h** header file. The SCSI ID and LUN should be placed in the *sc_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 SC_ASYNC flag should not be set on the first call to this operation and should only be set if a bus fault has occurred. The *immed_flag* field allows overlapping start operations to several devices on the SCSI bus. 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 **SCIoTUR** operation can then be issued to check on completion of the operation on a particular device.

Note that when the SCSI adapter is allowed to issue simultaneous start operations, it is important that a delay of 10 seconds be allowed between successive **SCIOSTUNIT** operations to devices sharing a common power supply since damage to the system or devices can occur if this precaution is not followed. Possible error values are **EIO**, **EFAULT**, **EINVAL**, **EACCES**, **ENOMEM**, **ETIMEDOUT**, **ENODEV**, and **ENOCCONNECT**.

- | | |
|--------------------|---|
| EIO | A system error has occurred. Consider retrying the operation several times, because another attempt may 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 may be successful. |
| ENODEV | The device is not responding. Possibly no LUNs exist on the present SCSI ID. |
| ENOCCONNECT | A bus fault has occurred. Try the operation again with the SC_ASYNC flag set in the sc_inquiry structure. In the case of multiple retries, this flag should be set only on the last retry. |

This operation requires **SCIOSTART** to be run first.

SCIOTUR

This operation issues a SCSI Test Unit Ready command to an adapter and aids in SCSI device configuration. The following is a typical call:

```
rc = ioctl(adapter, SCIOTUR, &ready_struct);
```

where *adapter* is a file descriptor and *ready_struct* is a **sc_ready** structure as defined in the **/usr/include/sys/scsi.h** header file. The SCSI ID and LUN should be placed in the *sc_ready* parameter block. The **SC_ASYNC** flag should not be set on the first call to this operation and should only be set if a bus fault has occurred. The status of the device can be determined by evaluating the two output fields: *status_validity* and *scsi_status*. Possible *errno* values are **EIO**, **EFAULT**, **EINVAL**, **EACCES**, **ENOMEM**, **ETIMEDOUT**, **ENODEV**, and **ENOCCONNECT**.

EIO A system error has occurred. Consider retrying the operation several (around three) times, because another attempt may 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 the **SCIOTUR** 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 may be successful.

ENODEV The device is not responding and possibly no LUNs exist on the present SCSI ID.

ENOCCONNECT A bus fault has occurred and the operation should be retried with the **SC_ASYNC** flag set in the **sc_inquiry** structure. In the case of multiple retries, this flag should be set only on the last retry.

This operation requires **SCIOSTART** to be run first.

SCIORESET

This operation causes a SCSI device to release all reservations, clear all current commands, and return to an initial state by issuing a Bus Device Reset (BDR) to all LUNs associated with the specified SCSI ID. A SCSI reserve command should be issued after the **SCIORESET** operation to prevent other initiators from claiming the device. Note that because a certain amount of time exists between a reset and reserve command, it is still possible for another initiator to successfully reserve a particular device. The following is a typical call:

```
rc = fp_ioctl(fp, SCIORESET, idlun, NULL);
```

where *fp* is a pointer to a file structure and *idlun* is a type int value that contains the SCSI and LUN ID values of the device to be stopped. The least significant byte contains the LUN, the next least significant byte contains the SCSI ID, and the upper two bytes should be set to zero. A nonzero return value indicates an error has occurred. Possible errno values are **EIO**, **EINVAL**, **EACCES**, and **ETIMEDOUT**.

- 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.

This operation requires **SCIOSTART** to be run first.

SCIOHALT

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 SCSI adapter sends a SCSI abort message to the device and is usually used by the SCSI device driver to abort the current operation instead of allowing it to complete or time out.

After the **SCIOHALT** operation is sent, the device driver must set the SC_RESUME flag in the next **sc_buf** structure sent to the adapter device driver, or all subsequent **sc_buf** structures sent are ignored.

The SCSI adapter also performs normal error recovery procedures during this command which include issuing a SCSI bus reset in response to a SCSI bus hang. The following is a typical call:

```
rc = fp_ioctl(fp, SCIOHALT, idlun, NULL);
```

where *fp* is a pointer to a file structure and *idlun* is a type int value that contains the SCSI and LUN ID values of the device to be stopped. The least significant byte contains the LUN, the next least significant byte contains the SCSI ID, and the upper two bytes should be set to zero. A nonzero return value indicates an error has occurred. Possible errno values are **EIO**, **EINVAL**, **EACCES**, and **ETIMEDOUT**.

- 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.

This operation requires **SCIOSTART** to be run first.

SCIODIAG

This command is most commonly used by a SCSI adapter diagnostic program. The **SCIODIAG** operation allows the SCSI adapter to run various diagnostic commands, which include:

- Internal Diagnostics Test
- SCSI Wrap Test
- Read/Write Register Test
- POS Register Test

These diagnostics options are defined, along with the **sc_card_diag** structure which is used in the call, in the **/usr/include/sys/scsi.h** header file. Any errors detected by the diagnostics are returned in the same **sc_card_diag** structure. Refer to the header file for a definition of the returned values.

Whenever an error is detected, an EFAULT errno value is returned along with the appropriate error statuses in the **sc_card_diag** structure. When the ENOMSG value is received, no information is provided in the error status fields.

Because this operation attempts no retries or error recovery, no error logging is provided. The following is a typical call:

```
rc = ioctl(adapter, SCIODIAG, &diag_struct);
```

where *adapter* is a file descriptor and *diag_struct* is a **sc_card_diag** structure. Possible errno values are **EIO**, **EFAULT**, **EINVAL**, **EACCES**, **ENOMSG**, and **ETIMEDOUT**.

EIO	An unrecoverable system error has occurred.
EFAULT	A user process copy has failed or diagnostics has failed without completing all tests.
EINVAL	An invalid diagnostic command was passed.
EACCES	The adapter is not in diagnostics mode.
ENOMSG	The diagnostic command has completed with errors.
ETIMEDOUT	The operation did not complete before the time-out value was exceeded.

This operation requires the adapter be in diagnostic mode.

SCIOTRAM

This operation is designed to test SCSI adapter RAM. However, it is currently not supported and therefore, returns an errno value of ENXIO if called.

SCIODNLD

This operation downloads microcode to the SCSI adapter and is used in configuring the SCSI adapter. This command can also be used to query for the current version of the microcode from the adapter. The following is a typical call:

```
rc = ioctl(adapter, SCIODNLD, &dnld_struct);
```

where *adapter* is a file descriptor and *dnld_struct* is a **sc_download** structure. Possible errno values are **EIO**, **EFAULT**, **EINVAL**, **EACCES**, **ENOMEM**, and **ETIMEDOUT**.

EIO	An unrecoverable system error has occurred or there is a checksum error with the microcode. If the adapter has been opened in diagnostics mode, this error is logged in the system error log. If the adapter has onboard microcode, it may still function properly.
------------	---

- EFAULT** A user copy has failed or a severe I/O error has occurred during the download and all subsequent commands to this adapter should cease. If the adapter has been opened in diagnostics mode, this error is logged in the system error log.
- EINVAL** An invalid input parameter was passed.
- ENOMEM** A request for memory failed.
- ETIMEDOUT** The operation did not complete before the time-out value was exceeded.

SCSI Device Driver Routines

A SCSI device driver should contain the following routines:

- **config**
- **ioctl**
- **open**
- **close**
- **read**
- **write**

Additional routines that you may need are the **strategy**, **iodone**, and **dump** routines. A **dump** routine is needed if the device is to be used as the receiver of a system dump. A **strategy** routine is needed in the case of block I/O device drivers that handle lists of **buf** structures.

Because the AIX operating system allows device drivers to be paged out of memory, certain functions and data structures of the driver must be guaranteed to be in memory to avoid page faults while running in the interrupt environment. This is required because disabling interrupts to any level will prevent the system from paging any pages in or out.

Routines that are called from an interrupt handler run in the interrupt environment while those that are called from a kernel or user process run in the process environment. Routines that run in the process environment can be interrupted by those running in the interrupt environment and can be preempted by processes with a higher process priority.

Routines running in the interrupt environment can only be interrupted by those running at a higher interrupt priority or by exceptions. In order for interrupt environment routines to avoid causing page faults, all code and data accessed in the interrupt environment must be pinned. Kernel services that deal with the pinning of code are **pin**, **pinu**, and **pincode**.

A function name is passed to the **pincode** kernel service when using it to pin driver code. However, the entire module that contains the function is pinned, not just the function itself. In the case of large device drivers, this is a wasteful use of memory. To get around this restriction, split the driver into two halves:

- a top half
- a bottom half

The top half contains preemptable routines that run in the process environment while the bottom half contains routines that run in the interrupt environment.

You can compile the two halves so that they are cross-linked. This allows one half of the driver to be automatically loaded together with the other half when it is loaded. The bottom (pinned) half should not have any dependencies on the top (unpinned) half, so you must carefully plan the two halves to avoid this condition, while at the same time minimizing wasted memory.

Usually, top-half routines include the **config**, **ioctl**, **open**, **close**, **read**, and **write** routines. Routines that usually reside in the bottom half are the **strategy**, **dump**, and **iodone** routines. Put any routines that cannot be paged out in the bottom half of the device driver.

In the case of read-only devices, such as CD-ROM, you do not need a **write** routine.

Top-Half Routines

This section discusses the following routines:

- **config**
- **ioctl**
- **open**
- **close**
- **read**
- **write**

config

This routine is commonly used when configuring, unconfiguring, or changing attributes of the device. It is called by various configuration methods which include the configure, unconfigure, and change methods. This routine can also be called to return the Vital Product Data (VPD) of the device. The following is a typical call to this subroutine from a configuration method:

```
sysconfig(SYS_CFGDD, &cfg, sizeof(struct cfg_dd));
```

where **sysconfig** is a subroutine and *cfg* is a **cfg_dd** structure which is defined in **/usr/include/sys/sysconfig.h**. The `cmd` field of the **cfg_dd** structure should contain one of the following parameters:

- CFG_INIT
- CFG_TERM
- CFG_QVPD

If the routine is called with the CFG_INIT parameter, perform simple error checking to ensure that the driver, as well as the device, is in the correct state. Any required data structures should also be allocated and initialized with the values contained in the Device Dependent Structure (DDS). The DDS should have been correctly built by the configuration method and passed to the config routine through the **cfg_dd** structure which contains a field that points to the DDS. After all data structures have been properly initialized, the device switch table entry should be built and then entered using the **devswadd** kernel service.

If called with the CFG_TERM parameter, the routine should first check that the driver is in the correct state and that the device is closed before deallocating associated data structures, removing the driver from the device switch table, and changing the status of the device to CLOSED in the device driver.

If the routine is called with the CFG_QVPD parameter, the driver should fill in the passed **uio** structure with the appropriate information about the device and the driver and then return this structure to the calling routine.

ioctl

This routine handles any I/O access to the SCSI device other than that which is handled by the **read** and **write** routines. Any diagnostic capabilities and the standard **IOCINFO** operation should also be handled by this routine.

Typical I/O requests are those that reset the device, obtain certain configuration parameters, or perform various basic accesses to the device. This routine can also perform various diagnostic functions once the device is opened in diagnostic mode. The **IOCINFO** operation is used to return information specific to the device as well as statistics about the usage of the device.

open

This routine opens a SCSI device by initializing any data structures, pinning the bottom half of the driver code, registering and initializing and interrupt handlers, setting up any DMA channels, and preparing any needed timers.

If the current open is the first open to the device, any global data structures as well as the bottom half of the driver should be pinned using the **pincode** kernel service. The pinning of the code and global data structures should occur before the device is able to generate interrupts since the interrupt handler is usually placed in the bottom half of the driver.

The SCSI adapter should then be opened with the **fp_open** kernel service. Once this operation completes, a **SCIOSTART** should then be issued for the device. If a forced open was attempted on the device, then a bus device reset (BDR) must also be issued through the **SCIORESET** operation. If the open operation was a diagnostic open, then processing is usually complete at this point since all that is required for a diagnostic open is that the appropriate data structures are pinned and the adapter is opened and started for the device.

If the open is a normal open, then an **SCIOTUR** operation should be issued to test for the readiness of the SCSI device. An **IOCINFO** ioctl call should then be made to obtain any necessary operating parameters for the driver.

The open routine should serialize its operation using simple locks. This prevents more than one application from modifying driver data structures. Unlock the resource when the driver has completed its processing to allow other applications to open SCSI devices.

Note that if any of the commands during the open operation fail, the device driver should clean up before exiting with an error. This includes freeing any allocated memory, unpinning any pinned code, and releasing any locks.

close

This routine closes the SCSI device by deallocating any data structures associated with the device as well as unpinning any code or data. The adapter driver is also stopped and any timers are also released.

Serialize the **close** routine so that no more than one application can modify the data structures of a driver. Use simple lock kernel services to accomplish this.

If the device was opened in normal mode and the **SC_RETAIN_RESERVE** flag was not set, then a release must be performed on the device to allow other initiators to reserve the device. If the reservation is to be retained, then the release should not be performed.

A **SCIOSTOP** operation should now be issued to the adapter driver which should also be closed using the **fp_close** kernel service.

read

This routine reads data from a SCSI device and returns it to the calling process through a **uio** structure.

The driver must perform any parameter validation, command building, and cross memory descriptor processing that might be necessary before calling the SCSI adapter strategy routine through **devstrat**. For block device drivers, if the read routine is merely a raw interface, the **uphysio** routine should be called if the request falls on block boundaries. The **uphysio** routine will end up calling the device driver's strategy routine. If the request does not fall on a block boundary, the driver should then break up the request into blocks. Odd sized transfers can then be handled by using **devstrat** to call the strategy routine after the **sc_buf** structure has been correctly built.

write

This routine writes data received through a **uio** structure from the calling process to a SCSI device.

The write routine should be very similar to the read routine described above in the type of processing that is required.

Bottom-Half Routines

This section discusses the following routines:

- **strategy**
- **dump**

strategy

Note: This routine is only present in drivers for block SCSI devices. This routine accepts a linked list of **buf** structures, processes them, and determines the proper SCSI command to send to the device to perform the required operation. The **buf** structure is defined in the `/usr/include/sys/buf.h` header file.

This routine should also determine if requests made in successive **buf** structures can be consolidated into one SCSI command. The request is placed on the I/O request queue for future processing by the device driver. Once a request is ready to be sent out to the device, the adapter's strategy routine should be called via the **devstrat** kernel routine. The adapter's strategy routine takes a pointer to an **sc_buf** structure as its only parameter.

After the I/O operation has completed, the user-level calling routine's **iodone** routine is called through the **iodone** kernel service. The caller's **iodone** routine should have been passed in the `b_iodone` field of the **buf** structure.

dump

This routine is needed if the device is to be used as a possible recipient of a system dump. The dump routine should first disable interrupts to the INTIODONE level which are naturally re-enabled at the end of the routine. Also, there are several commands, passed through the `cmd` parameter, that the routine must handle.

For a DUMPINIT command, this routine should just call the adapter dump routine through the **devdump** system call with a DUMPINIT command.

For a DUMPSTART command, this routine should just call the adapter dump routine through the **devdump** system call with a DUMPSTART command.

For a DUMPQUERY command, this routine should just call the adapter dump routine through the **devdump** system call with a DUMPQUERY command and then fill in the **dmp_query** structure passed in through the `arg` parameter with the appropriate information.

For a DUMPWRITE command, this routine should fill out an **sc_buf** structure for each **iovec** structure passed in through the `uio` parameter. A SCSI WRITE command should be constructed into each of the **sc_buf** structures and the adapter's dump routine should then be called through the **devdump** system call with a DUMPWRITE command. This should be repeated until all the **iovec** structures have been processed.

For a DUMPEND command, this routine should just call the adapter dump routine through the **devdump** system call with a DUMPEND command.

For a DUMPTERM command, this routine should just call the adapter dump routine through the **devdump** system call with a DUMPTERM command.

PVIDs

If a SCSI device is intended to contain a JFS filesystem, it must first meet two requirements. The first is that even though the device is not a hard disk, the driver must respond as such when it is issued an IOCINFO **ioctl** call. The device type must be either DD_DISK or DD_SCDISK as defined in the header file **/usr/include/sys/devinfo.h**.

The second requirement is that the device must contain a Physical Volume Identifier (PVID) so that it can be uniquely distinguished from other media on the system that contain filesystems. The PVID is used by the system to keep track of media even if it is moved from one SCSI ID to another or from one adapter to another. This allows the Logical Volume Manager (LVM) to maintain the consistency of volume groups (VGs) and logical volumes (LVs) that span several physical storage devices.

Normally, once a PVID is written, it remains with the device until it is overwritten or somehow damaged. This implies that a PVID only needs to be written once for a newly added device. New IBM hard disks will have a PVID assigned to them the first time they are configured on a RISC System/6000 by the configure method for hard disks. The system accomplishes this by first performing a read of the IPL record area from each disk detected on the system. If this area contains no PVID (the PVID value is 0), one is created and written out to the IPL record area. The PVID is created by the following scheme:

```
#define makehex(x) "0123456789abcdef"[x&15]

struct unique_id          unique_id;
struct utsname            uname_buf;
long                      machine_id;
struct timestruc_t        cur_time;
int                        i;
char                      pvidstr[33];
char                      bdevice[64];
int                        fd;
IPL_REC                   clearipl;
IPL_REC                   ipl_rec;
off_t                     offset;

bzero((caddr_t) &unique_id, sizeof (struct unique_id));

if (gettimer (TIMEOFDAY, &cur_time))
    exit (-1);

if (uname(&uname_buf))
    exit (-1);

machine_id = 0;
sscanf(uname_buf.machine, "%8x", &machine_id);

/* Note that words 3 and 4 remain 0 */
unique_id->word1 = machine_id;
unique_id->word2 = cur_time.tv_sec*1000 + cur_time.tv_nsec/1000000;

for(i=0;i<32;i++) {
    if (i&1)
        pvidstr[i] = makehex(unique_id[i/2]);
    else
        pvidstr[i] = makehex(unique_id[i/2]>>4);
}
pvidstr[32] = '\0';

/* lname is the name of the disk */
sprintf(bdevice, "/dev/r%s", lname);
fd = open(bdevice, O_RDWR);
```

```

    if (fd < 0)
        exit (-1);

offset = lseek(fd, PSN_IPL_REC, 0);
    if (offset < 0)
        exit (-1);

    if (read(fd, &ipl_rec, sizeof(ipl_rec)) < 0)
        exit (-1);

    if (ipl_rec.IPL_record_id != (unsigned int) IPLRECID) {
        /*
         * Boot record does not exist on disk yet.
         */
        ipl_rec = clearipl;
        ipl_rec.IPL_record_id = IPLRECID;
    }
    ipl_rec.pv_id = pvid;

offset = lseek(fd, PSN_IPL_REC, 0);
    if (offset < 0)
        exit (-1);

    if (write(fd, &ipl_rec, sizeof(ipl_rec)) < 0)
        exit (-1);

    close (fd);

```

Once the PVID has been created and stored to disk, it should also be added into the CuAt database for the disk under the `pvid` attribute.

SCSI Device Attributes

The following is a list of standard attributes that apply to the SCSI devices that are supported on the RISC System/6000. This list contains attributes most SCSI devices should contain. It is not an exhaustive list since each device may require additional attributes, depending on the implementation of each device.

- model_name** The name of the SCSI device that is returned in the inquiry string when an **SCIOINQU** ioctl call is made to the device. It is usually a 16-byte character string.

- maxlun** The maximum allowed value of the Logical Unit Number (LUN) for this device. Currently, all supported IBM SCSI devices have a **maxlun** value of zero. This indicates that only that particular device is allowed to occupy its SCSI ID and instructs the SCSI adapter to cease searching for other devices at other LUNs on the current SCSI ID.

- pvid** This attribute applies only to SCSI disks and contains the PVID value of the disk which is stored in the boot record of the disk. The value is usually a 32-bit character string and its default value is "none" as stored in the Predefined Attribute (**PdAt**) object class. Once a disk is assigned a PVID, create a Customized Attribute (**CuAt**) object class entry that contains the proper PVID value of the disk.

SCSI Configuration Methods

The configuration of SCSI devices is similar to the configuration of other devices on the system. Please refer to the Device Driver Configuration chapter elsewhere in this document

for a detailed explanation of configuration methods and how they should be written. However, there is an additional requirement that configuration methods for SCSI devices be written to take care of PVIDs (Physical Volume Identifiers) if necessary.

The configure method should try to read the PVID from the disk after the disk has been spun up and an INQUIRY command has been issued. If the disk has no PVID, then one should be created and written out if the media is allowed to contain an LV (logical volume) or filesystem. For information on how to accomplish this, see "PVIDs" on page 8-17.

If a non-NULL PVID is read, either of the following situations may exist:

- The device is already known to the system.
- The device is being moved from another RISC System/6000.

To determine which situation exists, the config method should scan the CuAt database for a previous record of this PVID. If none is found, then the disk should be assumed to be one being moved from another system.

If a non-NULL PVID is read, then the config method should verify that the PVID stored in the database is for the same logical name as the disk currently being configured. If so, then the disk is in the same position as it was when it was last configured and the config method does not need to perform any more PVID processing.

If a non-NULL PVID is read but its matching database entry does not match the logical name of the device being currently configured, this is an indication that the device has been moved from another location on the same system. If this occurs, calling the current parent SCSI adapter's config method will update the CuDv to correctly pair the device's logical name with its PVID.

Chapter 9. Writing a Virtual File System

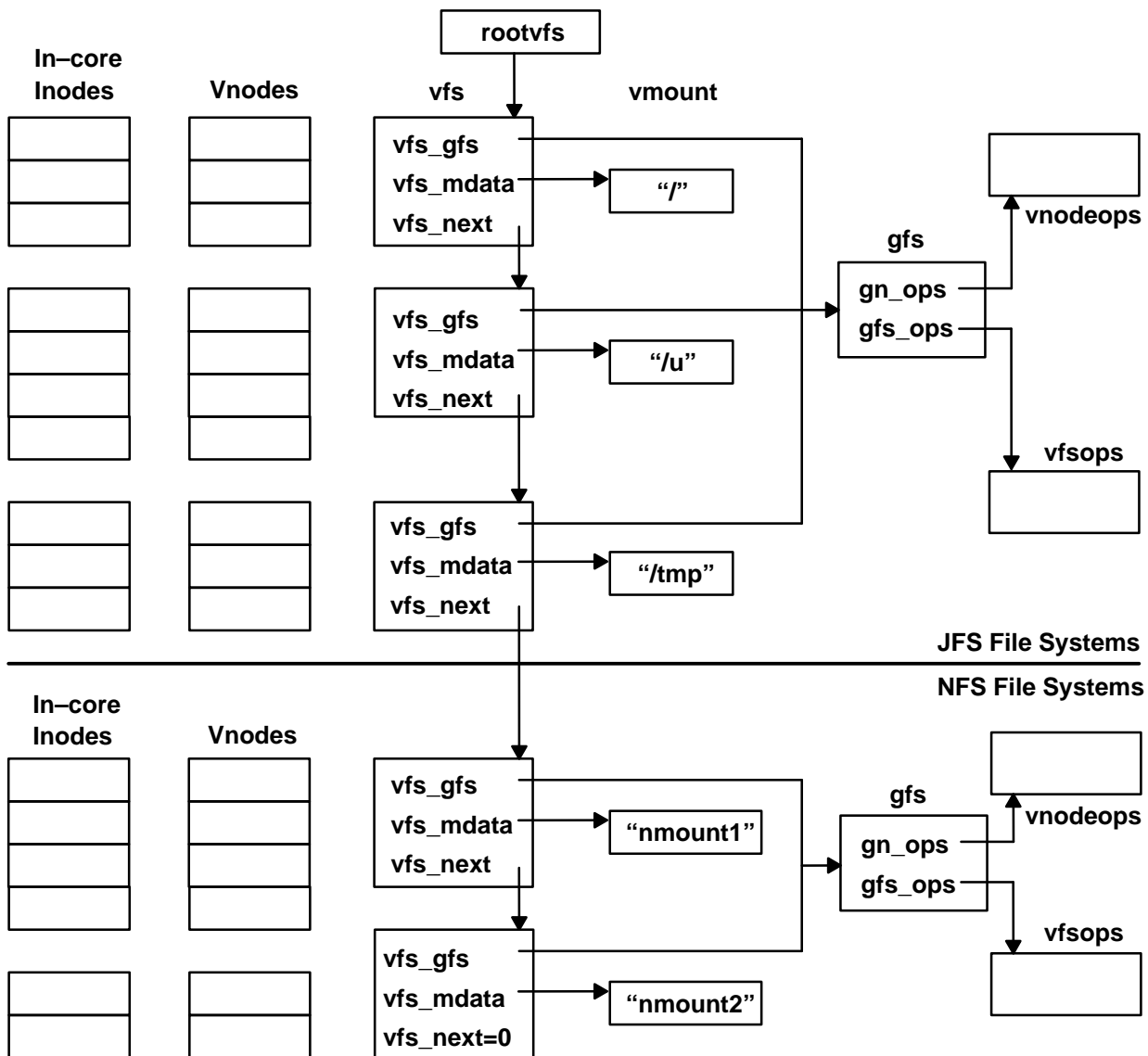
In addition to Journaled File System (JFS), Network File System (NFS), and the CD-ROM file system types included within AIX, it is possible to write your own Virtual File System (VFS). This may be desirable if you want to incorporate a new concept such as a distributed file system, or to use a new device such as a WORM drive or tape jukebox, or simply for performance reasons where a large amount of a specific type of data needs to be managed in an unusual manner.

This chapter is intended as a continuation to the virtual file system information in *AIX Version 4.1 Kernel Extensions and Device Support Programming Concepts*

Multiple File System Types within the Kernel

Each type of file system is represented within the kernel by a *struct gfs*. Each mounted file system is represented by a *struct vfs* which contains a pointer to the appropriate *struct gfs*. These structures are the basis for file system related operations.

The **vfs** structures are in a linked list, regardless of file system type, with the kernel variable *rootvfs* pointing to the first **vfs** structure in the list. The following **vfs** and **gfs** Structures figure shows the relationship between **vfs** and **gfs** structures.



vfs and gfs Structures

The **gfs** structures each contain a pointer to a *struct vnodeops* and a *struct vfsops*. These structures contain a list of functions which have a standardized interface by which system calls can invoke file related (with *vnodeops*), and file system related (with *vfsops*) operations.

You add new file system types to the kernel by loading the kernel extension into the kernel using the **sysconfig(KLOAD,...)** subroutine and then invoking the config entry point of the new kernel extension. The kernel extension contains the virtual file system dependent functions.

This in turn creates a **vnops** structure and a **vsops** structure, and initializes these structures with addresses of the functions within the extension. A temporary **gfs** structure is created, and the **gfsadd** kernel service is used to insert the gfs record into the kernel's array of **gfs** structures. At this point, the virtual file system type is usable, and file systems of the new type can be mounted.

Data Structures within a Virtual File System

As already described, a **gfs**, **vnops**, and **vsops** structure are created each time a new file system kernel extension is added to the kernel.

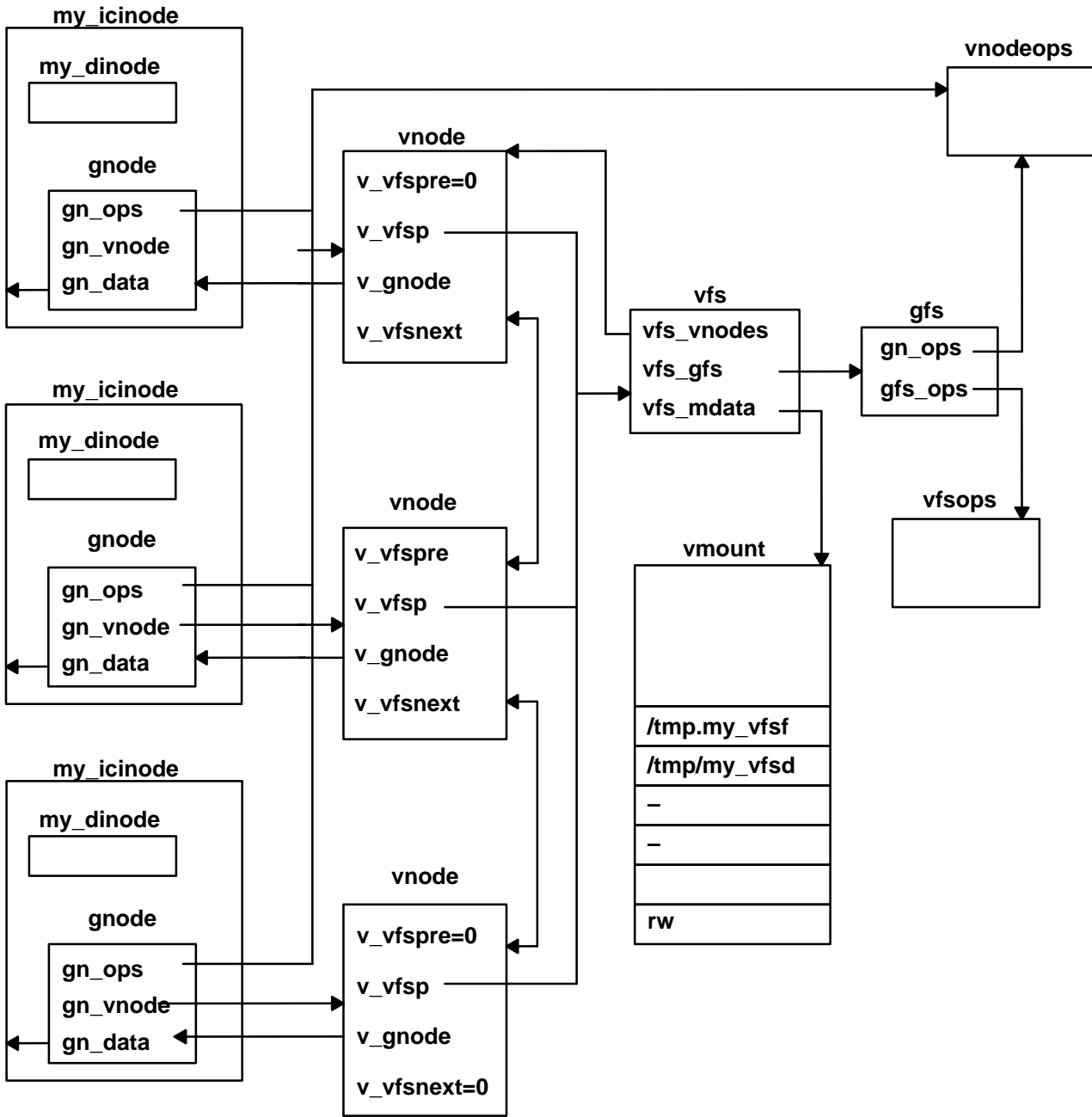
When new file systems are mounted, a **vfs** and **vmount** structure are created. The **vmount** structure contains specifics of the mount request, such as the object being mounted, and the stub over which it is being mounted. The **vfs** structure is the connecting structure which links the vnodes (representing files) with the vmount information, and the **gfs** structure.

The mount helper creates the **vmount** structure, and calls the **vmount** subroutine. The **vmount** subroutine then creates the **vfs** structure, partially populates it, and invokes the file system dependent **vfs_mount** subroutine which completes the **vfs** structure, and performs any operations required internally by the particular file system implementation. See "Mount Helper" on page 9-17 for more information about the mount helper.

Whenever a file is accessed, it is represented by a **vnode** structure, and an in-core **inode** structure, which also contains a **gnode** structure, and possibly a copy of the disk i-node. The **vnode** and **gnode** structures are file system independent, and the kernel can access them directly. However, the kernel has no information about the internal storage of the data files that the vnodes represent. To perform actions on a data file, the kernel passes a vnode pointer to the relevant vnode operation listed in the **vnops** structure.

The following File System Data Structures figure shows the relationships of the various data structures within a mounted sample file system `my_fs`. Note that in this file system type, there is a copy of the disk i-node (`my_dinode`) within each in-core i-node (`my_icinode`).

This figure also shows the way that the in-core i-node contains the gnode, which in turn contains pointers both to the containing in-core i-node (gn_data), and a direct link back to the **vnnodeops** structure (gn_ops) for efficiency.



File System Data Structures

The following sections contain descriptions of each of the file-system-independent structures.

gfs Structure

The **gfs** structure is defined in the **sys/gfs.h** header file:

```
struct gfs {
    struct vfsops          *gfs_ops;
    struct vnodeops        *gn_ops;
    int    gfs_type;       /* type of gfs (from vmount.h) */
    char    gfs_name[16];  /* name of vfs (eg. "jfs","nfs",..)*/
    int    (*gfs_init)();  /* ( gfsp ) - if ! NULL, */
                                /* called once to init gfs */
    int    gfs_flags;     /* flags for gfs capabilities */
    caddr_t gfs_data;     /* ptr to gfs's private config data*/
    int    (*gfs_rinit)();
    int    gfs_hold       /* count of mounts of the ... */
}
}
```

There is one **gfs** structure for each type of virtual file system currently installed on the machine. For each **gfs** entry, there may be any number of **vfs** entries.

The **gfs** structures are stored within a global array accessible only by the kernel. The **gfs** entries are inserted with the **gfsadd** kernel service. Only one **gfs** entry with a given **gfs_type** can be inserted into the array. Generally, **gfs** entries are added by the **CFG_INIT** section of the configuration code of the file system kernel extension.

The **gfs** entries are removed with the **gfsdel** kernel service. This is usually done within the **CFG_TERM** section of the configuration code of the file system kernel extension.

The operating system uses the **gfs** entries as an access point to the virtual file system functions on a type-by-type basis. There is no direct link from a **gfs** entry to all of the **vfs** entries of a particular **gfs** type. The file system code generally uses the **gfs** structure as a pointer to the **vnodeops** structure and the **vfsops** structure for a particular type of file system, although the **gnodes** also contain a pointer to the **vnodeops** structure for their type of file system.

vfs structure

The **vfs** structure is defined in the **sys/vfs.h** header file:

```
struct vfs {
    struct vfs    *vfs_next;      /* vfs's are a linked list */
    struct gfs    *vfs_gfs;      /* ptr to gfs of vfs */
    struct vnode  *vfs_mntd;     /* pointer to mounted vnode, */
                                /* the root of this vfs */
    struct vnode  *vfs_mntdover; /* pointer to mounted-over */
                                /* vnode */
    struct vnode  *vfs_vnodes;   /* all vnodes in this vfs */
    int           vfs_count;     /* number of users of this vfs */
    caddr_t       vfs_data;     /* private data area pointer */
    unsigned int  vfs_number;    /* serial number to help */
                                /* distinguish between */
                                /* different mounts of the */
                                /* same object */
    int           vfs_bsize;     /* native block size */
    short         vfs_rsvd1;     /* Reserved */
    unsigned short vfs_rsvd2;    /* Reserved */
    struct vmount *vfs_mdata;    /* record of mount arguments */
};
```

There is one **vfs** structure for each file system currently mounted.

New **vfs** structures are created with the **vmount** subroutine. This subroutine calls the **vfs_mount** subroutine found within the **vfsops** structure for the particular virtual file system type.

The **vfs** entries are removed with the **uvmount** subroutine. This subroutine calls the **vfs_umount** subroutine from the **vfsops** structure for the virtual file system type.

The **vfs** structure is central to each mounted file system. It provides access to the **vnodes** currently loaded for the file system, mount information through the **vfs_mdata** pointer, and provides a path back to the **gfs** structure and its file system specific subroutines through the **vfs_gfs** pointer.

The **vfs** structures are a linked list with the first **vfs** entry addressed by the **rootvfs** variable which is private to the kernel. The **vfs_mntd** pointer points to the **vnode** within the file system which generally represents the root directory of the file system. The **vfs_mntdover** pointer points to a **vnode** within another file system, also usually representing a directory, which indicates where the file system is mounted. In this sense, the **vfs_mntd** pointer corresponds to the object within the **vmount** structure referenced by the **vfs_mdata** pointer, and the **vfs_mntdover** pointer corresponds to the stub within the **vmount** structure referenced by the **vfs_mdata** pointer.

Refer to the “Mount Helper” section on page 9-17 for details of the **vmount** structure.

vnode structure

The **vnode** structure is defined in the **sys/vnode.h** header file:

```
struct vnode {
    ushort    v_flag;           /* see definitions below          */
    ulong     v_count;         /* the use count of this vnode    */
    int       v_vfsgen;       /* generation number for the vfs  */
    Simple_lock v_lock;       /* lock on the structure          */
    struct vfs *v_vfsp;       /* pointer to the vfs of this vnode */
    struct vfs *v_mvfsp;      /* pointer to vfs which was mounted over this */
                                /* vnode; NULL if no mount has occurred */
    struct gnode *v_gnode;    /* ptr to implementation gnode    */
    struct vnode *v_next;     /* ptr to other vnodes that share same gnode */
    struct vnode *v_vfsnext; /* ptr to next vnode on list off of vfs */
    struct vnode *v_vfsprev; /* ptr to prev vnode on list off of vfs */
    union v_data {
        void *          _v_socket;      /* vnode associated data */
        struct vnode *  _v_pfsvnode;    /* vnode in pfs for spec */
    } _v_data;
    char * v_audit;          /* ptr to audit object          */
};
```

Vnodes are the primary handles by which the operating system references files. Most **vnode** operations are passed a pointer to a **vnode** as their first parameter.

Vnodes are created by the **vfs**-specific code when needed, using the **vn_get** kernel service.

Vnodes are deleted with the **vn_free** kernel service.

Vnodes are the result of a path resolution. Each time an object (file) within a file system is located (even if it is not opened), a **vnode** for that object is located (if already in existence), or created. Naturally, **vnodes** for each directory searched to resolve the path are also created, or referenced. **Vnodes** are also created for files as the files are created.

The **vnode** structure provides the link between the **vfs** structure and the **gnode** structure. There are two **vnodes** for one **gnode** only in the case of file-over-file mounts. In this case, the **gnode** refers to the first related **vnode**, and the other **vnodes** for that **gnode** are linked using the **v_next** field.

gnode structure

The **gnode** structure is defined in the **sys/vnode.h** header file:

```
struct gnode {
    enum vtype gn_type;           /* type of object: VDIR,VREG,... */
    short  gn_flags;             /* attributes of object          */
    ulong  gn_seg;              /* segment into which file is mapped */
    long   gn_mwrcnt;           /* count of map for write        */
    long   gn_mrdcnt;           /* count of map for read         */
    long   gn_rdcnt;           /* total opens for read          */
    long   gn_wrcnt;           /* total opens for write          */
    long   gn_excnt;           /* total opens for exec           */
    long   gn_rshcnt;           /* total opens for read share     */
    struct vnodeops *gn_ops;
    struct vnode *gn_vnode; /* ptr to list of vnodes per this gnode */
    dev_t  gn_rdev;           /* for devices, their "dev_t" */
    chan_t gn_chan;           /* for devices, their "chan", minor's minor */

    Simple_lock  gn_reclk_lock; /* lock for filocks list        */
    int          gn_reclk_event; /* event list for file locking  */
    struct filock *gn_filocks; /* locked region list           */

    caddr_t gn_data;           /* ptr to private data (usually contiguous) */
};
```

A gnode refers directly to a file (regular, directory, special, and so on), and is usually embedded within a file system implementation specific in-core i-node.

Gnodes are created as needed by file system specific code at the same time as creating in-core i-nodes. This is normally immediately followed by a call to the **vn_get** kernel service to create a matching vnode.

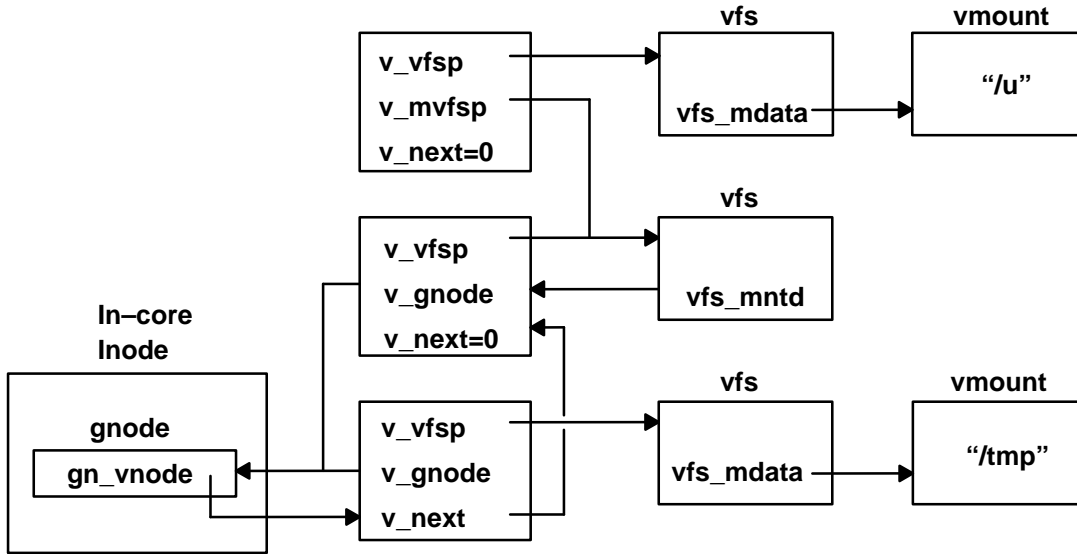
The gnode structure is usually deleted either when the file it refers to is deleted, or when its in-core i-node is removed to make room for an in-core i-node representing a more recently accessed file.

In early UNIX virtual file system implementations, the vnode represented the file system independent information about a file, and the in-core i-node was totally file system implementation specific. In AIX, the gnode is a part of the in-core i-node which is uniform, but the remainder of the in-core i-node remains implementation specific.

File-Over-File Mounts

It is possible to have more than one vnode referring to a gnode. This occurs when a file is mounted over another file. In this case, the **v_mvfs** field within the vnode of the file which is being mounted over is set to point to a newly created **vfs** structure. The new **vfs** structure represents a file system with one file, whose vnode points to the gnode of the file which is mounted.

The following File-Over-File Mounts figure shows the situation where a file in `/tmp` is mounted over a file in `/u`. Note that the `gnode` points to one of the `vnodes` which is linked to the next `vnode` for that `gnode` with the `v_next` field.



File-Over-File Mounts

Components of a Third-Party Virtual File System

The following are the basic components that you must provide to implement a virtual file system for the AIX operating system:

1. The kernel extension containing initialization functions, `vfs` operations, `vnode` operations, and virtual memory operations.
2. A file system helper or some other mechanism to create file systems on media.
3. A mount helper or a specialized mount command.
4. A configuration program to load the kernel extension into the kernel.
5. A software installation package (preferably an `installp` package) to install the software on a customer's computer, and modify `/etc/vfs`.

Creating the Virtual File System Kernel Extension

The virtual file system kernel extension is similar to a driver in that it runs in kernel mode and has a configuration entry point. Considerations for compiling and linking are also similar.

The following is an example of Makefile entry for building a virtual file system kernel extension:

```
CFLAGS = -O -D_AIX -D_SUN -D_KERNEL -DKERNEL -I../.. -I. -qxref -qlist
CDEFS  = -U_STR_ -U_MATH_ -D_KERNEL -D_AIX -D_IBMR2      \
        -DMACHINE=_IBMR2 '-DMACHNAME="R2_System"'        \
        -DNLS -D_NLS -DMSG -Daiws -Dunix

all:    my_vfs

my_vfs: my_vfsops.o my_vnodeops.o
        /bin/ld -b"binder:/usr/lib/bind glink:/lib/glink.o" \
        -s -D0 -H512 -T512                                \
        -lcsys -lsys                                       \
        -b"I:/lib/kernex.exp I:/lib/syscalls.exp E:myfs_t.exp
map:sym.myfs_t"                                           \
        -e myfs_config                                     \
        my_vfsops.o my_vnodeops.o                         \
        -o my_vfs

my_vnodeops.o: my_vnodeops.c myfs.h
        ${CC} ${CFLAGS} ${CDEFS} -c my_vnodeops.c 2> $.err

my_vfsops.o:   my_vfsops.c myfs.h
        ${CC} ${CFLAGS} ${CDEFS} -c my_vfsops.c 2> $.err
```

Note: To bring in the appropriate version of `printf` and have debugging statements appear on the terminal, do *not* use `cc` for the linking phase of the operation.

Entry Points within the File System Kernel Extension

This section describes the **config**, **init**, and **rootinit** entry points.

config

```
int myfs_config (int cmd, struct uio *uiop);
```

The primary entry point within the file system kernel extension is the config entry point. The configuration program calls this entry point once the kernel extension is loaded. The configuration routine to which the config entry point refers usually consists of a case statement which switches based on the command passed in. The two cases which should be taken into account are **CFG_INIT** and **CFG_TERM**. However, you can add more commands.

The **CFG_INIT** branch should:

1. Check that the extension has not already been initialized.
2. Initialize any private data structures.
3. Create and add an entry into the **gfs** table using the **gfsadd** kernel service. This calls the **init** entry point from the **gfs** structure.
4. Register the page fault handler (**strategy** routine) using the **vm_mount** kernel service.
5. Pin any code or data within the file system which must not be paged. For example, pin interrupt handlers or the strategy routine.

The **CFG_TERM** branch should:

1. Verify that it is safe (appropriate) to terminate the extension. For example, verify that the file system has been unmounted.
2. Remove the **gfs** struct entry using the **gfsdel** kernel service.
3. Unpin the areas pinned during the **CFG_INIT**.

init

```
int myfs_init (struct gfs *gfsp);
```

This routine is called by the **gfsadd** kernel service immediately after creating the **gfs** structure. The address of the function is passed into **gfsadd** within the **gfs_init** field of the **gfs** structure. Any functions that must be performed after the **gfs** entry is created but before users can start to use the file system should be performed here.

The most common operation within this function is to initialize file system dependent structures. The function may also spawn a pager process using the **creatp** and **initp** kernel services.

rootinit

This function is only called for the root file system. The entry point is passed to the **gfsadd** kernel service in the **gfs_rinit** field of the **gfs** structure. The most common operations are to initialize the root **vmount** and **vfs** structures, and to open the root device.

VFS Operations within the File System Kernel Extension

The **vfs** operations are those operations which affect an entire file system, such as:

- **mount**
- **unmount**
- **sync**

The addresses of these functions are stored within a **vfsops** structure that is pointed to by the **gfs_ops** field of the **gfs** structure for a particular virtual file system implementation.

For specific details of the **vfs** operations requirements, see the *AIX Version 4.1 Kernel Extensions and Device Support Programming Concepts*.

The following are descriptions of the **vfs** operations for a typical virtual file system.

vfs_mount

```
int myfs_mount(struct vfs *vfsp, struct ucred *crp);
```

This routine is called by the **vmount** subroutine, and performs the following actions in addition to any internal work required:

1. Extract important data passed by the mount helper program in the **vfsp->vfs_mdata** field.
2. Obtain the vnode of the object being mounted with the **lookupvp** kernel service.
3. Check that the object being mounted is appropriate.
4. Initialize the following fields within the **vmount** structure (pointed to by **vfsp->vfs_mdata**):
 - **vmt_fsid**
 - **vmt_vfsnumber**
 - **vmt_time**
5. If applicable, call **vm_mount** to register the file system's strategy routine.

6. Create a vnode (and matching in-core i-node) for the root directory of the file system, and initialize the `v_mvfsp` field of the stub vnode (pointed to by `vfsp->vfs_mntdover`) with its address.
7. Use the `vn_rele` vnode operation to release the object's vnode.

vfs_unmount

```
int myfs_unmount(struct vfs *vfsp, int flags, struct ucred *crp);
```

This routine is called by the **umount** subroutine. This routine performs the following actions:

1. Free up any resources used by this mount, including vnodes, using the `vn_free` vnode operation.
2. Set the `v_mvfsp` field of the vnode for the stub to NULL.
3. Set `vfsp->vfs_mntd` to NULL.
4. Call the `vm_umount` kernel service to unregister the file system's strategy routine.
5. Call the `vfsrele` kernel service to release the system-created resources.

vfs_root

```
int myfs_root(struct vfs *vfsp, struct vnode **vpp,
              struct ucred *crp);
```

This routine is frequently called when locating files within a file system. Essentially, it finds, or creates a vnode (associated with an in-core i-node) for the root of the file system (that is, the file containing the root directory). Before returning success, the routine should check that the file system is mounted, and calls the `vn_hold` vnode operation for the vnode it is about to return.

vfs_statfs

```
int myfs_statfs(struct vfs *vfsp, struct statfs *statfsp,
                struct ucred *crp);
```

The `vfs_statfs` routine returns basic file system statistics in the form of a `statfs` structure. See the `sys/statfs.h` header file for the fields of the `statfs` structure. The methods to obtain this information are implementation dependent.

vfs_sync

```
int myfs_sync();
```

Every 60 seconds, and whenever the `sync` command is run, the `sync` function is called. This in turn calls the `vfs_sync` function *once* for each different file system type. Note that a pointer to a `gfs` structure is passed in.

The actual use of this routine is totally implementation dependent, but once completed each file system of the same `gfs` type should be updated on secondary storage (assuming there is any) to the point where it is consistent.

vfs_vget

See File Systems Operations in *AIX Version 4.1 Technical Reference, Volume 5: Kernel and Subsystems*.

vfs_cntl

See File Systems Operations in *AIX Version 4.1 Technical Reference, Volume 5: Kernel and Subsystems*.

Vnode Operations within the File System Kernel Extension

The vnode operations are the operations that work on specific files (represented by vnodes) within a file system. This section covers the required vnode operations.

See File Systems Operations in *AIX Version 4.1 Technical Reference, Volume 5: Kernel and Subsystems* for details on the requirements for the vnode operations.

vn_hold

```
int myfs_hold (struct vnode *vp);
```

This routine simply increments the `v_count` field (the usage count) of the vnode.

vn_rele

```
int myfs_rele (struct vnode *vp);
```

This routine decrements the usage count of the vnode. If the usage count is reduced to zero, remove the vnode from memory using the **vn_free** vnode operation.

It is practical to release the related in-core i-node at this point, assuming (as is normal) that this is the only vnode referencing the in-core i-node.

Check that the file system was being unmounted but could not complete because there were still open files (`vfsp->vfs_flag & VFS_UNMOUNTING`). If this was the last vnode for the file system (`vfsp->count == NULL`), call the **vfsrele** kernel service to release the `vfs` entry and complete the unmount.

vn_getattr

```
int myfs_getattr (struct vnode *vp, struct vattr *ubuf,  
                 struct ucred *crp);
```

This routine obtains the attributes of the file referenced by `vp` from the related gnode and in-core i-node which is in a file-system specific format, and returns the information in a **vattr** structure.

The following shows the **vattr** structure as defined in the **sys/vattr.h** header file:

```
struct vattr {
    enum vtype va_type; /* from gnode.gn_type*/
    mode_t va_mode; /* access modes from in-core i-node*/
    uid_t va_uid; /* user id from in-core i-node*/
    gid_t va_gid; /* group id from in-core i-node*/
    union {
        dev_t _va_dev;
        long _va_fsid; /* vfs_mdata->vmt_fsid.fsid_dev*/
    } _vattr_union;
    long va_serialno; /* i-node number from in-core i-node*/
    short va_nlink; /* number of links to fil from in-core
                    i-node*/

    long va_size; /* file size in bytes from in-core i-node*/
    long va_blocksize; /* block size for file system */
    long va_blocks; /* blocks reserved for file
                    from in-core i-node*/
    struct timeval va_atime; /* last file access time
                             from in-core i-node*/
    struct timeval va_mtime; /* last file modification
                             time from in-core i-node*/
    struct timeval va_ctime; /* last disk i-node modification
                             time from in-core i-node*/

    dev_t va_rdev; /* from gnode.gn_rdev*/
    long va_nid; /* use unamex(), field nid*/
    chan_t va_chan; /* from gnode.gn_chan*/
    char *va_acl; /* Access Control List*/
    int va_aclsiz; /* size of ACL*/
    int va_gen; /* inode generation number */
};
```

vn_open

```
int myfs_open (struct vnode *vp, int flags, int ext,
              caddr_t vinfo, struct ucred *crp);
```

The **vn_open** routine should check the validity of the operation, in particular any file-system dependent tests, such as the requested open mode conflicting with the mode in which the file is already opened by another process.

The file can be bound in at this point, or the binding can be delayed until the first read or write to the file. See the “Virtual Memory Operations” section on page 9-15 for information on binding.

This routine should also increment the appropriate usage counter within the gnode (*gn_rdcnt*, *gn_wrcnt*, or *gn_excnt*).

vn_close

```
int myfs_close (struct vnode *vp, int flags, caddr_t vinfo,
               struct ucred *crp);
```

The **close** routine calls **vn_close** routine. Any usage counters incremented by the **vn_open** call should be decremented. Do *not* remove the vnode within **vn_close**, this is done within **vn_rele**.

vn_strategy

```
int myfs_strategy (struct buf *bp);
```

The **strategy** routine is called by the virtual memory manager when a page of memory is referenced that is mapped to a file, but the memory has not been paged in. Usually, this routine simply places the buffer on the list of pending work for the pager, and wakes the pager.

The buffer pointer passed in will have the following values in its fields:

b_bcount	Always PAGESIZE (4K).
b_blkno	The 512-byte block number of the offset within the file (0, 8, 16, and so on).
b_baddr	The 4K offset within the mapped segment. This rolls over to zero once a 256-MB segment is crossed.
av_forw	Non-null if there are multiple requests to process. This will happen if the virtual memory manager attempts to perform read ahead operations.
b_flags	Set to B_READ for a read operation, B_WRITE for a write operation and (B_READ B_PFSTORE) for a read before a write operation.

vn_rdwr

```
int myfs_rdwr (struct vnode *vp, enum uio_rw op, int flags,
              struct uio *uiop, int ext, caddr_t vinfo,
              struct vattp *vattp, struct ucred *crp);
```

This routine performs file (not directory) reads or writes. Follow these procedures:

1. Verify that the file type is correct.
2. Return success if the transfer size is zero.
3. Return EINVAL if the start offset is negative.
4. Verify that the file is mapped. If it is not mapped in, map it using the **vms_create** kernel service.
5. Use **vm_move** to copy the data from the source to the destination, reducing the transfer length on reads that go past the end of the file.

Note: The **vm_move** amounts to a copy which can page fault.

vn_lookup

```
int myfs_lookup(struct vnode *dp, struct vnode **vpp,
               char *name, int flags, struct ucred *crp);
```

This routine finds the file *name* which is the basename of a file in the directory *dp*, and returns a vnode pointer for the file. If the file is not found, set *vpp* to NULL.

vn_readdir

```
int myfs_readdir (struct vnode *vp, struct uio *uiop,
                 struct ucred *crp);
```

This routine is called to read directories, which it returns in file-system independent format using a **dirent** structure. Even **read** and **fread** call this routine if the file is a directory, and the directory contents are reconstituted (imperfectly) from the structures.

Read directory entries from the directory starting with the first entry at or beyond *uiop->uio_offset*. Entries are read into an internal buffer while the buffer size is less than the *uio->resid* passed in (which is 4K), or until the end of the directory is reached. You do not

need to place empty directory entries in the buffer. The buffer should be filled in with as many entries as possible.

The buffer is then copied back into user space using the **uiomove** kernel service which updates `uio->uio_resid`. The **vn_readdir** routine then updates `uiop->uio_offset` to be the offset within the directory of the next unread entry.

If the routine is called with no directory entries beyond the `uio_offset`, the routine should return with `uio_resid` untouched.

The following listing shows the fields of the **dirent** structure as found in the **sys/dir.h** header file:

```
#define _D_NAME_MAX 255
struct dirent {
    ulong_t      d_offset;    /* offset within directory file
                               of next directory entry */
    ino_t        d_ino;       /* i-node number of entry */
    ushort_t     d_reclen;    /* Use DIRSIZ(struct dirent *)
                               after d_namlen is initialized */
    ushort_t     d_namlen;    /* strlen(d_name field)*/
    char         d_name[_D_NAME_MAX+1];
    /* name must be no longer than this including the '\0' */
};
```

Virtual Memory Operations

Once a file is opened, it is bound into an address range. This address range is not in either the user or kernel space, but is inserted into the kernel address space on a temporary basis by **vms_create**.

Accesses to this range may cause page faults which cause the registered **strategy** routine to be called. The **strategy** routine then arranges for the information to be paged in or out without itself page faulting. For this reason, pin the **strategy** routine and any data areas it uses.

A separate process can be started to handle page faults as requested by the **strategy** routine, and this process is under no obligation to avoid page faulting.

Binding the File to an Address Range

To bind a file to an address range, call the **vms_create** kernel service. This reserves part of the 52-bit virtual memory address space of the machine, and associates it with the gnode passed in. The `vmid` passed back is effectively a segment register value, which should be placed in the `gn_seg` field of the **gnode** structure. This is typically done by the file system's open routine or by the read/write function (upon the first I/O attempt).

On subsequent reads or writes to the file, the **vn_rdwr** subroutine passes this `vmid` to **vm_move**, along with transfer length, `uio` pointer and so on. Then **vm_move** places the address space of the file within the kernel address space long enough to perform the transfer.

For more information on **vms_create**, refer to kernel services in *AIX Version 4.1 Technical Reference, Volume 5: Kernel and Subsystems*.

The Page Handling Process

The pager for a virtual file system may run as a separate process, a process spawned by the **vfs_init** subroutine. The routine uses the **e_sleep_thread** kernel service on a list of **buf** structures to which the **vn_strategy** subroutine appends further entries. Once the list contains at least one request, the pager parses the **buf** structures, and in many virtual file system implementations, attempts to order them depending upon the situation.

For instance, buffers may be separated into requests for different media so that one disk does not need to wait until another disk has completed a series of transfers. Or the requests may be ordered to optimize disk seek sequencing.

Paging a Block In or Out

The actual technique used to page a block in or out varies greatly between virtual file systems. In a situation where the file system resides directly on media, such as in the CD-ROM file system, the **buf** structures are modified to contain the actual device block number, and device IDs, then the **devstrat(bufp)** kernel service is called by the page fault handler. This calls the **strategy** routine within the device driver to perform the transfer.

Once the data is transferred, call the **iodone(bp)** kernel service to notify the virtual memory manager that the memory space is now valid.

File System Helper

The file system helper for a virtual file system type is listed in the **/etc/vfs** file. This program is called by AIX system management programs, usually to operate on an unmounted file system. Examples include:

- fsck** Calls the file system helper to perform an implementation-dependent check of the contents of a file system.
- mkfs** Invokes the file system helper to create a file system.
- chfs** Uses the file system helper to extend a file system.

An example of an entry in the **/etc/filesystems** file follows:

```
/mymnt :
      dev           = /dev/fd0
      vfs           = myfs
      mount         = true
      options       = rw
      type          = myfs
      nodename      = myfs
```

For the preceding example, the command:

```
mkfs /mymnt
```

translates to a call to the file system helper with the following parameters:

```
fshop_make 5 5 7 -ip 0 \
name=/tmp/my_vfsd,label=/tmp/my_vfsd,dev=/tmp/my_vfsf
```

Note: The **fshop_make** is *not* the name of the file system helper. The subcommand **argv[0]** is set to a name other than the program name using the **execlp** subroutine.

We suggest that third-party virtual file system writers use **argv[1]**, which is a command from the **fshelp.h** file such as **FSHOP_MAKE**.

The command:

```
fsck -v myfs
```

translates to the following call:

```
fshop_check 1 3 6 -ifp 0 device=/dev/hd1,mounted
```

The basic subcommands are:

argv[A_NAME (=0)]

The command in `char *` format.

argv[A_OP (=1)]

The command as an entry from the **fshelp.h** file.

argv[A_FSFID (=2)]

File descriptor number for the file system.

argv[A_COMFD (=3)]

File descriptor number for pipe.

argv[A_MODE (=4)]

Mode flags. For example, `i = FSHMOD_INTERACT_FLAG`.

argv[A_DEBG (=5)]

Debug level.

argv[A_FLGS (=6)]

Operation-dependent flags.

The functions performed by any file system helper are not fixed. Implement only those appropriate to a virtual file system. Furthermore, it is acceptable for some functions to be omitted, and implemented as separate programs, although the AIX commands will not be able to invoke them.

Mount Helper

The format of the mount helper is similar to that of the file system helper. The mounting program may be in any format, but for it to be invoked by the AIX **mount** command it should be listed as the mount helper in the **/etc/vfs** file, and conform to the following invocation format:

argv[0]	The command path
argv[1]	Command (M=mount, U=unmount)
argv[2]	Debug Level
argv[3]	Nodename
argv[4]	Object (the file system being mounted)
argv[5]	Stub (the directory that the file system is mounted over)
argv[6]	Options (for example, "rw" = read/write)

The basic operations of the mount helper are:

1. Checks that the operation is valid.
2. Allocates space for a **vmount** structure with enough space at the end for the additional parameters. See the **sys/vmount.h** header file for information on the **vmount** structure.
3. Populates the **vmount** structure.
4. Calls **vmount (vmountptr, size_of_vmount_struct)**.

The contents of the **vmount** structure are defined as follows:

```
#define VMT_OBJECT 0 /* I index of object name */
#define VMT_STUB 1 /* I index of mounted over stub name */
#define VMT_HOST 2 /* I index of (short) hostname */
#define VMT_HOSTNAME 3 /* I index of (long) hostname */
#define VMT_INFO 4 /* I index of binary vfs specific info */
/* includes network address, opts, etc*/
#define VMT_ARGS 5 /* I index of text of vfs specific args*/
#define VMT_LASTINDEX 5 /* I the last in the array of structs */
struct vmount {
    ulong vmt_revision; /* revision level, currently 1 */
    ulong vmt_length; /* length of structure and data */
    fsid_t vmt_fsid; /* Do not set */
    int vmt_vfsnumber; /* Do not set */
    time_t vmt_time; /* Do not set */
    ulong vmt_timepad; /* Do not set */
    int vmt_flags; /* general mount flags */
    int vmt_gfstype; /* type of gfs, e.g. MNT_JFS */
    struct vmt_data {
        short vmt_off; /* I offset of data, word aligned */
        short vmt_size; /* I actual size of data in bytes */
    } vmt_data[VMT_LASTINDEX + 1];
};
```

Note: The `vmt_off` field is the offset from the start of the **vmount** structure to the start of the data element.

The following is an example of a hexadecimal dump of a **vmount** structure:

```
0000: 00000001 0000006c 00000000 00000000 .... ..1 ....
0010: 00000000 00000000 00000000 00000000 .... ....
0020: 00000008 003c000e 004c000e 005c0002 .... .<.. .L.. \..
0030: 00600002 00640004 00680003 2f746d70 .`. .d.. .h.. /tmp
0040: 2f64656d 6f766673 66000000 2f746d70 /dem ovfs f... /tmp
0050: 2f64656d 6f766673 64000000 2d000000 /dem ovfs d... -...
0060: 2d000000 00000000 72770000 00000000 -... .... rw.. ....
```

Virtual File System Configuration Program

You must load the virtual file system kernel extension in much the same way as a driver. This operation can be performed by a separate program that is invoked by an rc script, or preferably, as a rule in the `Config_Rules` database. It can also be invoked by the `mount` helper when the first file system of a particular type is mounted.

Load the kernel extension with the **sysconfig** subroutine:

```
sysconfig (SYS_KLOAD,...)
```

Then use **sysconfig** to call the configuration entry point of the extension:

```
sysconfig (SYS_CFGKMOD,...)
```

The following is an example of using the **sysconfig** subroutine:

```
struct cfg_load load;
struct cfg_kmod kmod;
.
.
.
load.path = kern_ext_name;
if( sysconfig( SYS_KLOAD, &load, sizeof(load) ) == -1 ){
    fprintf( stderr, "Unable to SYS_KLOAD %s, errno = %d\n",
            load.path, errno );
    exit(1);
}
    fprintf( stderr, "loaded at 0x%08x\n", load.kmid );

/* initialize extension */
kmod.kmid = load.kmid;
kmod.cmd = CFG_INIT;
kmod.mdiptr = (caddr_t) &kmod.kmid;
kmod.mdilen = sizeof( kmod.kmid );

if( sysconfig( SYS_CFGKMOD, &kmod, sizeof(kmod)) == -1 ){
    fprintf( stderr,
            "Unable to configure module, errno=%d\n",
            errno);
    exit(1);
}
```

During development, and for some file system implementations, you may want to unload the kernel extension. To do this, use the following calls:

```
kmod.cmd = CFG_TERM;
sysconfig (SYS_CFGKMOD, &kmod,...);
sysconfig (SYS_KUNLOAD, &load,...);
```

Software Installation Package

You should package your software in the form of an installp module, because this is the form in which customers expect to receive software for AIX Version 4.1. See "Software Product Packaging" in *AIX Version 4.1 General Programming Concepts* [Volume 1: Writing Programs](#) for more information.

In addition to installing the software in the customer's system, the installation program also:

1. Appends an entry into the **/etc/vfs** file for the new file system type. For example:

```
my_fs      8  /etc/helpers/my_fsmnhelp  /etc/helpers/my_fshelper
```

2. (Optional) Inserts a new rule into the Config_Rules database to invoke the kernel extension installation program during system startup.
3. (Optional) Creates an initial entry in **/etc/filesystems**.

/etc/vfs

The **/etc/vfs** file describes the currently installed virtual file systems. Entries include:

- Comment lines, which start with a #.
- The general control line which defines the default local, and (optionally) the default remote virtual file system. This is typically:

```
%defaultvfs      jfs      nfs
```

- VFS entries that consist of:

name The name of the vfs type.

type The number of the vfs type.

mount-helper The name of the mount helper, or "none". This is with respect to **/etc/helpers** unless it starts with a / (slash).

file system-helper

Same conventions as mount-helper.

Use white space to separate these fields.

For a user file system, use numbers from MNT_USRVFS (=8) to MNT_AIXLAST (=15) for the type. These values are defined in the **sys/vmount.h** header file. IBM reserves the right to assign file system numbers between zero and MNT_USRVFS – 1 (=7).

Once the example file system "my_fs" is inserted, the end of **/etc/vfs** contains:

```
%defaultvfs      jfs      nfs
#
cdrfs             5         none     none
jfs               3         none     /etc/helpers/v3fshelper
nfs               2         none     /etc/helpers/nfsmnthelp
                 none     remote
my_fs             8         none     /etc/helpers/my_fsmnthelp
                 none     /etc/helpers/my_fshelper
```

Loading the File System Kernel Extension during System Startup

The easiest way to load the kernel extension into the kernel during system startup is to insert a new rule into the Config_Rules database. To do this, create a file with the new rule in it. Our sample file is named **my_fs_rule.add**:

```
Config_Rules:
  phase = 2
  seq = 0
  rule = "/etc/methods/installmy_fs -a /etc/drivers/my_fs"
```

To install the rule, use the following ODM commands:

```
odmdelete -q "rule LIKE '*my_fs*'" -o Config_Rules
odmadd my_fs_rule.add
```

The **odmdelete** command deletes old Config_Rules entries which may have been inserted for this virtual file system. If you use the **odmadd** command twice without **odmdelete**, there will be two rules in the database, and the extension installation program will be invoked twice.

Do *not* use this technique during development, as a faulty kernel extension could prevent a successful system startup. Instead, add the rule at the end of the development cycle when the extension is stable.

Glossary

- Disk i-node** A disk i-node is an i-node that resides on a disk (or some other secondary storage). Generally, this i-node contains all of the information about a file which also resides on the disk except its name (which is in one or more directories), and its contents (which are stored in the data area reserved for the file). This information includes file type, size, permissions, access times, and data location information (assuming the file is a data file, not a special file). Disk i-nodes are referred to by their i-node number which is unique within a particular file system.
- gfs entry** A **gfs** structure exists for each TYPE of virtual file system that is being used within the system. For example, there is one **gfs** structure for the Journaled File System (JFS) regardless of the number of JFS file systems currently mounted, and if there are any NFS mounts, there is one **gfs** structure for NFS. The **gfs** structure provides the system with a set of pointers to the standardized functions for a file system type.
- gnode** A gnode is a data structure containing information about a file that is currently (or has recently) being referenced. The gnode is effectively an extension of the vnode, and contains the information which is directly associated with the file, independent from the vfs by which the file is being accessed. There can only be one gnode within a system for a particular file at any one time. The gnode is a system-defined structure (that is, vfs independent), contained within an in-core i-node which is vfs dependent.
- in-core i-node** An in-core i-node is a data structure maintained by a file system kernel extension which represents a file on disk. The structure of the in-core i-node is implementation dependent; however, it contains a gnode structure that is implementation independent. It is common that the in-core i-node also contain a copy of the disk i-node for the file it represents.
- virtual file system (vfs)**
A virtual file system represents a conventional mounted file system. The purpose of having virtual file systems is so that the kernel, and system calls need not know anything about the file system internals, but rather they can perform file system operations, and file operations through a defined interface that is independent of the data storage format or media.
- vnode** A vnode is a data structure representing a file within a virtual file system. Accesses by system calls to a file are performed using a vnode pointer which provides an implementation independent interface. Each vnode contains a pointer to the vfs which contains it, and to the gnode which contains limited information about the file, and a pointer to the in-core i-node that contains the implementation specific data about the file.

Chapter 10. STREAMS-Based TTY Subsystem Interface

This chapter describes the programming interface of the tty subsystem. This interface can be used by a tty programmer to write a specific module (such as a converter or a line discipline), or a specific driver.

The following information is provided:

- The description of the tty modules and drivers and especially the protocol used in the exchange of STREAMS messages.
- The way the tty subsystem is integrated in a multiprocessor environment.
- The list of the IOCTL operations relating to the tty subsystem.
- The tty data structures extracted from the **str_tty.h** file where tty constants, functions and messages are defined.

Note: Because the tty subsystem is based on STREAMS, a good knowledge of STREAMS concepts is highly recommended. See Related Information on page 10-28.

Overview

The tty subsystem is based on STREAMS. The stack structure of a tty stream is made up of the following modules:

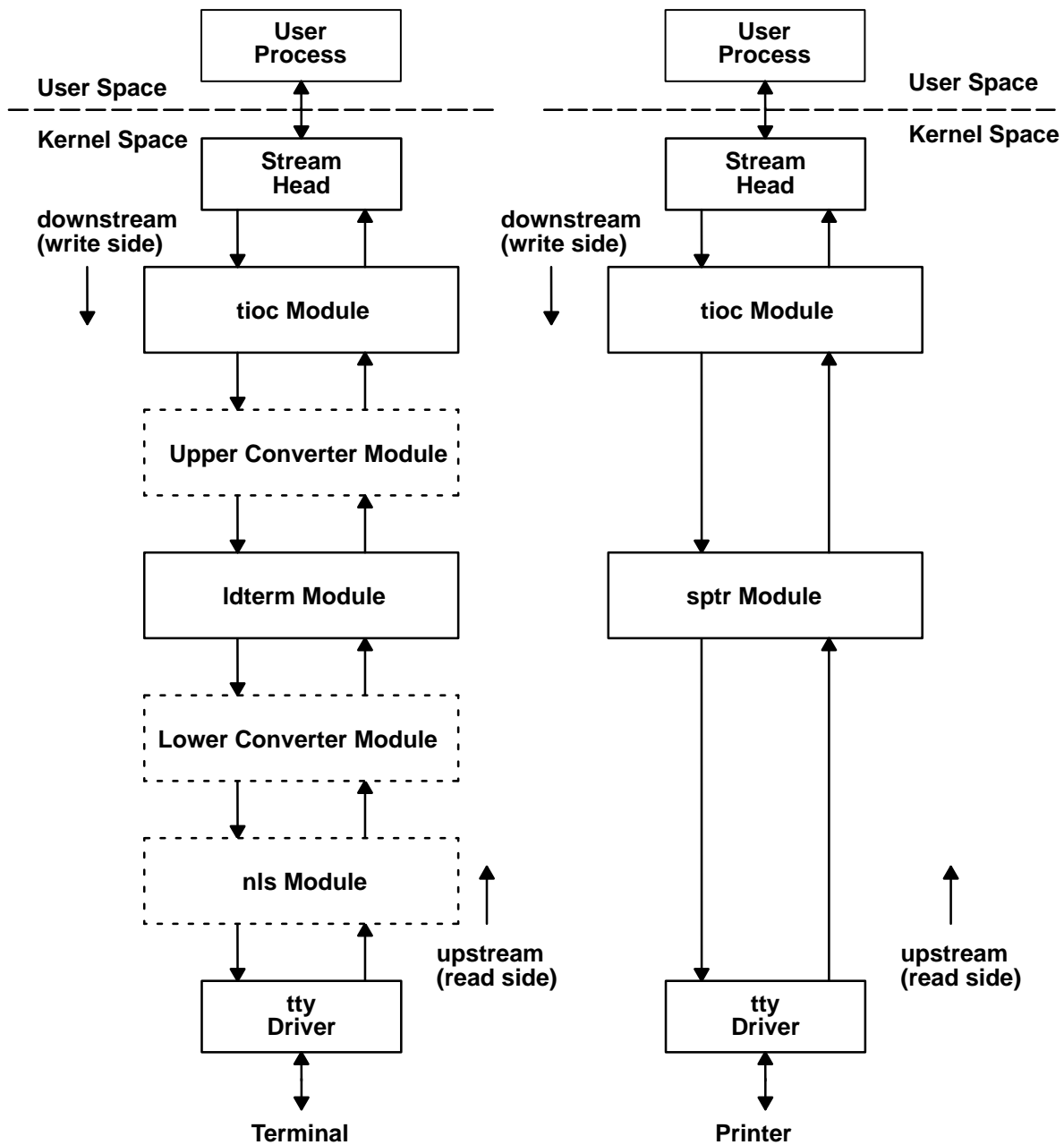
- The stream head, which processes the user's requests. It is common to all tty devices, regardless of the line discipline or tty driver in use.
- The **tioc** module. It is provided to facilitate processing of the transparent ioctl operations.
- The line discipline (**ldterm**, **sptr** or **slip**).
- The stream end. It is a tty driver (**cxma**, **lft**, **lion**, **rs**, or **sf**) or pseudo-driver (**pty**).

The tty stream may also contain:

- An optional character mapping module (**nls**). It is a converter module pushed above the tty driver to support input and output data mapping. The **nls** module uses the input and output map files provided by the **setmaps** command.
- A pair of optional converters to convert upstream and downstream data. An upper converter module is pushed above the line discipline and a lower converter module is pushed below the line discipline. (For example, **uc_sjis** and **lc_sjis** are the upper and lower converters used to convert IBM-932 code set into or from the EUC code set handled by the **ldterm** line discipline.)

Note: The optional modules are not described in this chapter. See "The TTY Subsystem" in *AIX Version 4.1 General Programming Concepts, Volume 1: Writing Programs* for more information.

The following figure shows the tty stream structure for a terminal and for a serial printer. The line discipline for a terminal is **ldterm**. The line discipline for a serial printer is **sptr** and the stream structure is simpler (the converters, and **nls** modules are not needed.)



TTY Terminal Stream

TTY Serial Printer Stream

User modules can also be added or can replace the standard tty modules or drivers to support other specific functionalities.

The means of communication within a stream are messages. Messages are sent through a stream by calls to routines of each queue (write-side or read-side queue) in the stream. Messages can be generated by a driver, a module, or by the stream head. The messages are exchanged between modules or drivers in conformance with protocol rules which are described in the following sections.

Messages that are passed from the stream head toward the driver are said to travel *downstream* (also called *write side*). Similarly, messages passed in the other direction, from the driver to the stream head, travel *upstream* (also called *read side*).

For more general information about tty subsystem and STREAMS, refer to *AIX Version 4.1 General Programming Concepts* *Volume 1: Writing Programs* and *AIX Version 4.1 Communications Programming Concepts*.

Stream Head

The stream head is the interface between the stream and the user application. The stream head processes the user's requests; it is common to all tty devices regardless of the line discipline or tty driver in use.

The stream head performs the following functions:

- It allocates the stream as the controlling terminal if none is already allocated.
- It handles the process group and the session associated with the controlling terminal.
- It handles the job control.
- It processes the **M_PCSIG** and **M_HANGUP** messages coming from the line discipline and generates the appropriate signals to the appropriate user process. It also handles the message type **M_ERROR** from any downstream modules.
- It handles "trusted path" for security purposes. This functionality is shared between the stream head and the driver.
- It handles the "revoke" function which kills all waiting processes attached to the specified file descriptor.

The stream head processes various messages. In particular, the **M_SETOPTS** message is sent by the line discipline module to alter some characteristics of the stream head. The **SO_ISTTY** flag contained in an **M_SETOPTS** message indicates to the stream head that the stream is established for a terminal.

The IOCTLs listed here are directly processed by the stream head:

TIOCSPGRP (or **TXSPGRP**)

Sets the process group identifier for the tty.

TIOCGPGRP (or **TXGPGRP**)

Gets the process group identifier for the tty.

TIOCGSID

Gets the session identifier for the tty.

TXISATTY

Answers if this stream is a tty or not.

TCTRUST

Sets or resets the trust flag in the stream head.

TCQTRUST

Gets the state of the tty (trusted or not).

TCSAK

Sets or resets the state for secure attention key (SAK).

TCQSAK

Identifies the process of a user asking for the SAK sequence.

TIOCCONS

Sets the console redirection active or not.

TCXONC

Generates the following message types downstream based on the input parameter:

Input Parameter	Message Type
M_START	TCCON

	M_STOP	TCCOFF
	M_STARTI	TCION
	M_STOPI	TCIOFF
TCFLSH	Generates a M_FLUSH message type based on input parameter type (FLUSHR , FLUSHW).	

TIOC Module

The IOCTL system calls from applications that are addressed to modules or drivers can be processed according to two STREAMS mechanisms:

- The first mechanism allows the processing of ioctls issued using the **I_STR** call. The associated structure (**strioc** defined in the **stropts.h** file) contains the IOCTL command and the number of bytes of data. The **I_STR** calls are handled by the stream head and forwarded downstream.
- The second mechanism allows the processing of ioctls other than **I_STR**. These IOCTLs are known as transparent IOCTLs. (The transparent mechanism transparently supports applications developed prior to the introduction of STREAMS.)

If an IOCTL is not processed in the stream head, the stream head creates an **M_IOCTL** message which includes the ioctl arguments, and sends this message downstream for processing by a specific module or driver. **M_IOCTL** messages containing a transparent IOCTL have a **TRANSPARENT** indicator in their associated **iocblk** structure defined in the **stream.h** file. The **M_IOCTL** messages can be followed by one or more **M_DATA** blocks.

The **tioc** module is provided in the tty subsystem to facilitate processing of the transparent IOCTLs in the lower modules. It has two functions:

1. Identification of the transparent IOCTLs that can be processed on a tty stream:

The **tioc** module maintains a table containing a list of ioctls commands with the number of bytes to copy, and a type (see the **tioc_reply** structure in “Open Routine”). The **tioc** module recognizes a predefined list of ioctls. In addition, it asks the lower modules or drivers for their specific IOCTLs and adds them to its list. This protocol is described in “Open Routine.”

2. Management of the data transfers from or to the stream head for the transparent ioctls:

Since a module or driver has no user context, it has to request the stream head to perform data transfers between user and kernel environments. The **tioc** module is responsible for the transfer of data from or to the user space. This avoids messages transfers through all modules. If the ioctl processing requires data to be transferred in from user space, **tioc** issues an **M_COPYIN** message. In the same way, **tioc** issues an **M_COPYOUT** message to transfer out any data back to user space. The **tioc** module may also intermix **M_COPYIN** and **M_COPYOUT** messages in any order, if both input and output transfers are required for an ioctl. The protocol used to perform the **M_COPYIN** and **M_COPYOUT** messages is described in “Copy in Data for an IOCTL” and “Copy out Data for an IOCTL”.

Open Routine

The open routine has to initialize the IOCTL commands table of the **tioc** module. For that, it sends downstream an **M_CTL** message containing the **TIOC_REQUEST** command.

The **M_CTL** message has to come downstream to the stream end. Then the drivers or modules send an **M_CTL** message upstream containing the **TIOC_REPLY** command. Each module on read side adds an **M_DATA** message to the **M_CTL** message if required. The

M_DATA message contains a list of **tioc_reply** structures which define the specific IOCTLS to be handled by **tioc**.

The **tioc_reply** structure, defined in the **str_tty.h** file, is the following:

```
struct tioc_reply {
    int tioc_cmd;           /* command */
    int tioc_size;         /* number of bytes to copy */
    int tioc_type;         /* type of ioctl */
};

/*
 * STREAM tioc module tioc_type
 */
#define TTYPE_NOCOPY      0    /* don't need any copies */
#define TTYPE_COPYIN     1    /* need a M_COPYIN */
#define TTYPE_COPYOUT    2    /* need a M_COPYOUT */
#define TTYPE_COPYINOUT  3    /* need both M_COPYIN and M_COPYOUT */
#define TTYPE_IMMEDIATE  4    /* use immediate value */
```

The following examples show how IOCTL types are used.

The IOCTLS that do not require data transfer to or from STREAMS (for example, **TIOCSDTR**) are declared as **TTYPE_NOCOPY**.

The IOCTLS that use immediate parameter values (for example, **TCXONC**) are declared as **TTYPE_IMMEDIATE** and are called as follows:

```
ioctl(fd, TCXONC, TCOON);
```

The IOCTLS that request information from the STREAMS and require COPYOUT operations (for example, **TIOCGETA**) are declared as **TTYPE_COPYOUT** and are called as follows:

```
ioctl(fd, TIOCGETA, termios_ptr);
```

In the preceding sample, **termios_ptr** is the address of a **termios** structure to be obtained from the tty.

The ioctls that send information to the STREAMS and require M_COPYIN operations (for example, **TIOCSETA**) are declared as **TTYPE_COPYIN** and are called as follows:

```
ioctl(fd, TIOCSETA, termios_ptr);
```

In the preceding sample, **termios_ptr** is the address of a **termios** structure to be set in the tty.

Copy in Data for an IOCTL

When an **M_IOCTL** message arrives from the stream head, the write-side put routine recognizes an IOCTL which requires data to be copied in, and it also knows the size of required data.

So, the write-side put routine immediately sends an **M_COPYIN** message with this size to the stream head. The stream head processes a **copyin** function and sends an **M_IOCTLDATA** message downstream containing the data and the result of the **copyin** process.

If the **copyin** was successful, the write-side put routine sends an **M_IOCTL** message, linked to an **M_DATA** message containing the data, downstream. The **M_IOCTL** message is sent downstream until a module can process it and send an **M_IOCACK** up to the stream head.

If the **copyin** was not successful, the write-side put routine sends an **M_IOCNAK** message to the stream head.

Copy out Data for an IOCTL

When an **M_IOCTL** message arrives from the stream head, the write-side put routine recognizes an IOCTL which requires data to be copied out.

The **tioc** module sends the **M_IOCTL** message downstream and the IOCTL is processed by the appropriate module. If the ioctl processing is successful, this module sends an **M_IOCACK** upstream containing data to copy out.

The read-side put routine sends an **M_COPYOUT** message, linked to an **M_DATA** message, to the stream head. The stream head processes a **copyout** function and sends downstream an **M_IOCDATA** message containing the result of the **copyout** process.

The write-side put routine replies to this message by an upstream message, depending on the success of the **copyout**. If it was successful, it sends an **M_IOCACK** message. If not, it sends an **M_IOCNAK** message.

LDTERM Module

The **ldterm** module is a key part of the STREAMS-based tty subsystem. It supplies the line discipline for terminal devices. This line discipline is POSIX compliant and also supports the System V and BSD interfaces. The **ldterm** module provides the terminal interface functions specified in the **termios.h** header file. (The **termios.h** header file is described in *AIX Version 4.1 Files Reference*). The **ldterm** module also handles EUC and multibyte characters.

The **ldterm** module processes various types of STREAMS messages. The messages processed by this module are listed in Messages Summary. Any other message received by **ldterm** is passed downstream or upstream unchanged.

Open Routine

When first called, the open routine allocates space for the **ldterm** internal structure, and also sends an **M_SETOPTS** message upstream. This message includes a **stroptions** structure part defined in the **stream.h** file, which contains options that inform the stream head how to use this stream.

The open routine allocates space for the **termios** structure, which contains the flags used to control the terminal. The flags are defined in the **termios.h** file.

- **c_iflag** defines input modes.
- **c_oflag** defines output modes.
- **c_cflag** defines hardware control modes.
- **c_lflag** defines terminal functions handled by **ldterm**.
- **c_cc** defines the control characters.

The open routine initializes the flags with default values.

When **ldterm** is pushed during stream initialization, it sends some **M_CTL** messages downstream that query the driver for the current flags (**TIOCGETA** command) and the flags the driver may process (**MC_CANONQUERY** command). The driver may modify the flags processed by **ldterm** with its response to **MC_CANONQUERY**.

Close Routine

The close routine sends an **M_SETOPTS** message upstream to undo stream head changes done on the first open.

The **ldterm** module also sends **M_START** and **M_STARTI** messages downstream to undo the effect of any previous **M_STOP** and **M_STOPI** messages.

Finally, the close routine frees all the outstanding buffers allocated by the **ldterm** module.

Read-Side Put Routine

The **ldterm** read-side put routine processes the following STREAMS messages coming from downstream modules or driver:

M_FLUSH

This is a request to flush the read-side queue of all its data messages and all data being accumulated. The read-side put routine processes the request and forwards the message upstream.

M_BREAK

The **M_BREAK** message provides the following status information:

`break_interrupt`, `parity_error`, `framing_error` or `overrun`.

The read-side put routine reads the status, processes the event and discards the message.

M_DATA

The read-side put routine processes the **M_DATA** message and performs various actions according to the characters encountered in the data and the setting of the **termios** flags:

- The read-side put routine generates echo characters which are sent downstream in **M_DATA** messages.
- **ldterm** can control the output flow of data: if the **IXON** flag is set, the read-side put routine processes START (VSTART) and STOP (VSTOP) characters and sends **M_START** and **M_STOP** messages downstream.
- **ldterm** can control the input flow of data: if the **IXOFF** flag is set and input is to be stopped or started, the read-side put routine generates **M_STOPI** and **M_STARTI** messages downstream.
- If the **ISIG** flag is active, the read-side put routine manages signals characters. It sends **M_PCSIG** messages upstream when signal characters are encountered and then discards these characters.
- At the logical termination of input, the read-side put routine sends the currently buffered characters upstream to the stream head. The logical termination of input depends on the state of the **ICANON** flag:
 - If **ICANON** is set, **ldterm** is in canonical input mode. In this case, the input logically terminates at the end of a line of input. The canonical line termination characters are NEWLINE, EOF, EOL, and EOL2.
 - If **ICANON** is not set, **ldterm** is in noncanonical or raw input mode. In this case, the input terminates when at least VMIN bytes are present in the input message buffer or when the timer specified by VTIME expires.

M_IOCACK

The **M_IOCACK** message signals a positive acknowledgment of a previous **M_IOCTL** message.

- If the **M_IOCACK** message acknowledges the **TIOCGETA**, **TIOCSETA**, **TIOCSETAW**, or **TIOCSETAF** commands, the **termios** structure is updated as specified in the commands.
- If the **M_IOCACK** message acknowledges switching the current canonical mode (**-ICANON** to **ICANON**, or **ICANON** to **-ICANON**), the read-side put routine sends an **M_SETOPTS** message upstream to notify the stream head of the change.

- If the message acknowledges a **TIOCOUTQ** command, the required number of bytes are added to the reply value in the **M_IOCACK** message.
- In all other cases the message is sent upstream.

M_CTL

The **M_CTL** messages received on the read-side and processed by **ldterm** are sent by the driver for different reasons:

1. The **M_CTL** message can be sent to communicate changes in the driver's state:

If the **CLOCAL** flag of **ldterm** is not set, and the carrier state has just made a transition from **on** to **off**, the read-side put routine sends an **M_HANGUP** message upstream to inform the stream head that the terminal connection was broken.

2. The **M_CTL** message can be sent to answer to a previous request or command from **ldterm**. The following commands contained in an **M_CTL** message are processed:

TIOCGETA The driver sends this command either as a response to an inquiry for current settings or to reflect an asynchronous change in the flags of its **termios** structure. The read-side put routine copies the **termios** structure from the attached **M_DATA** message block into its internal **termios** structure. Then, it frees the **M_CTL** message.

MC_NO_CANON

The input canonical processing normally performed on **M_DATA** messages is disabled and those messages are passed upstream unmodified; this is used by modules or drivers that perform their own input processing (For example a pseudo terminal in TIOCREMOTE mode connected to a program that performs the input processing).

MC_DO_CANON

All input processing performed on **M_DATA** messages is enabled.

MC_PART_CANON

The driver sends this message to notify **ldterm** that it handles some part of the input processing itself (for example, flow control). An **M_DATA** message containing a **termios** structure is expected to be attached to the original **M_CTL** message. The **ldterm** module will examine the **c_iflag**, **c_oflag**, and **c_lflag** fields of the **termios** structure and will process only those flags which have not been turned ON.

TIOCGETMODEM

The driver sends this message to communicate the state of its flag **modem carrier on**. The associated **M_DATA** message contains a value of 1 (one) to indicate the carrier is **on**, or a value of 0 (zero) to indicate the carrier is **off**. This information is used to update the **ldterm** module state. If the **CLOCAL** flag of **ldterm** is not set, and the carrier state has just made a transition from **on** to **off**, the read-side put routine sends an **M_HANGUP** message upstream to inform the stream head that the terminal connection was broken.

When the command is processed the **M_CTL** message is freed.

3. The **M_CTL** message can be sent to answer to the **TIOC_REQUEST** coming from the **tioc** module with a **TIOC_REPLY**. (See "TIOC Module" on page 10-4.) This **M_CTL** message is forwarded upstream.

Write-Side Put Routine: Immediate Processing

The **ldterm** write-side put routine immediately processes the following STREAMS messages:

M_FLUSH

The write-side put routine flushes the write-side queue and discards any buffered output data. Then, it forwards the message downstream.

M_DATA

If the write-side queue is empty, the write-side put routine processes the **M_DATA** message. Else, it queues the **M_DATA** message to the write-side queue for later processing by the write-side service routine.

M_IOCTL

The write-side put routine validates the format of the **M_IOCTL** message and checks for known IOCTL command:

- If the message format is invalid, it turns the **M_IOCTL** message into an **M_IOCNAK** message, and returns it upstream.
- If the IOCTL command is not recognized, it forwards the **M_IOCTL** message downstream for processing by other modules.
- If the IOCTL is recognized, the write-side put routine determines if the command must be processed in the proper sequence relative to **M_DATA** messages. If so, it queues the **M_IOCTL** message to the write-side queue for later processing. The commands that require processing in sequence are:
TIOCSETAW, TIOCSETAF, TCSETAW, TCSETAF, TCSBRK, TCSETXW, TCSETXF, and TCSBREAK.
Otherwise, the write-side put routine processes the command immediately.

M_READ

The **M_READ** message is processed only if the **ldterm** module is in noncanonical input mode.

The **M_READ** message is sent by the stream head to notify downstream modules when an application has issued a read request and there is not enough data queued at the stream head to satisfy the request. The message contains the number of characters requested by the application.

If VTIME is positive, the write-side put routine starts an input timer. When the timer expires, it sends all buffered input upstream.

M_START, M_STOP, M_STARTI, M_STOPI

Some ioctl commands (**TCXONC, TCOOF**) are issued by the stream head to block the flow of data. These blocked data are stored in the modules queues.

The stream head sends an **M_IOCTL** message with the **TCOON** command when it wants to unblock the data. However this message can also be blocked if the waiting data entirely fill the queues.

To solve this blocking state, the stream head sends the following high priority messages:

- **M_START** to restart output of data
- **M_STOP** to stop output of data
- **M_STARTI** to restart input of data

- **M_STOPI** to stop input of data.

The **ldterm** write-side put routine updates internal state fields and may forward these messages downstream.

Write-Side Service Routine: Delayed Processing

The write-side service routine processes messages that may be delayed due to STREAMS flow control or to ioctls requiring sequential processing. The write-side service routine is called by the scheduler.

M_DATA

The write-service service routine processes the data according to the flags in the **termios** structure. It sends the processed characters downstream to the driver when the write-side queue fills up and all of the data is processed.

M_IOCTL

Some ioctl commands must wait until output drains before they are processed. **M_IOCTL** messages containing these commands are queued on the write-side queue so that the write service routine processes them in the correct sequence relative to preceding data. The commands that require processing in sequence are:

TIOCSETAW, TIOCSETAF, TCSETAW, TCSETAF, TCSBRK, TCSETXW, TCSETXF, and TCSBREAK.

Multibyte Processing

The **ldterm** module handles the extended UNIX code (EUC) character encoding scheme. This encoding scheme enables the module to process multibyte characters as well as single-byte characters. It correctly handles backspacing and tabulation expansion for multibyte characters.

By default, multibyte processing by **ldterm** is turned off. When **ldterm** receives an **EUC_WSET** IOCTL call that sets character width on screen to a value greater than one, it enables multibyte processing.

When multibyte processing is turned on, the **ldterm** module automatically calls EUC routines as necessary.

Messages Summary

Messages include read-side messages and write-side messages.

Read-Side Messages

Messages processed by **ldterm**:

M_BREAK, M_CTL, M_DATA, M_FLUSH, M_HANGUP, M_IOCACK.

Messages sent by **ldterm** upstream:

M_CTL, M_DATA, M_ERROR, M_FLUSH, M_HANGUP, M_IOCACK, M_IOCNAK, M_PCSIG, M_SETOPTS.

Write-Side Messages

Messages processed by **ldterm**:

M_CTL, M_DATA, M_FLUSH, M_IOCTL, M_READ, M_START, M_STARTI, M_STOP, M_STOPI, M_NOTIFY.

Messages sent by **ldterm** downstream:

M_BREAK, M_CTL, M_DATA, M_DELAY, M_FLUSH, M_IOCTL, M_READ, M_STOP, M_START, M_STOPI, M_STARTI.

SPTR Module

The serial printer module (**sptr**) supplies a specific line discipline for serial printers. It is mainly used by the spool subsystem.

The user interface (spool/application) is specified in the **lp** special file. (The **lp** special file is described in *AIX Version 4.1 Files Reference*)

The **sptr** module processes various types of STREAMS messages. The messages processed by this module are listed in Messages Summary. Any other message received by **sptr** is passed downstream or upstream unchanged.

Open Routine

When first called the open routine allocates space for the **sptr** structure and sends an **M_SETOPTS** message upstream to the stream head.

When **sptr** is pushed on the stream, it sends an **M_CTL** message to the driver (containing the **TIOCMGET** command) to obtain the status and the CTS modem control signal.

Read-Side Put Routine

The read-side put routine processes the following STREAMS messages:

M_FLUSH

The read-side put routine flushes the read-side queue, then forwards the message upstream. The flushing is dependent on the input parameter (FLUSHR or FLUSHW).

M_DATA

The read-side put routine stores the **M_DATA** message for next **M_READ** message.

M_PCPROTO

The driver sends this message containing the **LPWRITE_ACK** command to indicate to **sptr** that all the data it sent in the previous **M_PROTO** message was transmitted to the line.

M_IOCACK

The **M_IOCACK** message signals the positive acknowledgment of a previous **M_IOCTL** message.

M_CTL

The **M_CTL** message is sent by the driver to communicate changes in the driver's state, or to reply to a previous **M_CTL** message.

The following commands contained in an **M_CTL** message are processed:

TIOCMGET The **sptr** module registers the CTS state in its private data.

TIOC_REPLY The **sptr** module adds information concerning the specific IOCTL commands in the message and sends it upstream.

cts_on or **cts_off**

This message is sent by the driver when the CTS signal changes. **sptr** uses it to get information about the printer connection and in turn takes appropriate actions.

Write-Side Put Routine

The write-side put routine immediately processes the following STREAMS messages: (Messages not listed here are simply forwarded downstream.)

M_FLUSH

The write-side put routine flushes the write-side queue and forwards the message downstream. The flushing is dependent on the input parameter (FLUSHR or FLUSHW).

M_DATA

If the write-side queue is not empty, the write-side put routine queues the message to this queue for later processing. The message will be processed by the write-side service routine when called by the scheduler.

If the write-side queue is empty, the write-side put routine processes the message immediately: it formats the data if needed, and sends an **M_PROTO** message with the data downstream to the driver.

M_IOCTL

The write-side put routine processes the following IOCTL commands:

IOCINFO, LPRGET, LPRSET, LPRMODG, LPRMODS, LPRGTOV, LPRSTOV, LPQUERY, LPRGETA, LPRSETA, LPWRITE_REQ, TCGETA, TCSETA, TCSETAW, and TCSETAF.

The **sptr** module will reply to **LPWRITE_REQ** when it receives the write completion message (**M_PCPROTO** message) from the driver or if an error condition arrives (timeout, disconnection of the printer).

Other IOCTL commands are sent downstream unchanged.

M_READ

This message is sent by the stream head as a data request. The write-side put routine returns the required number of data previously stored on the read side.

Messages Summary

Messages include read-side messages and write-side messages.

Read-Side Messages

Messages processed by **sptr**:

M_CTL, M_DATA, M_FLUSH, M_IOCACK, M_PCPROTO.

Message sent by **sptr** upstream:

M_FLUSH, M_IOCACK, M_IOCNAK, M_SETOPTS.

Write-Side Messages

Message processed by **sptr**:

M_DATA, M_FLUSH, M_IOCTL, M_READ.

Messages sent by **sptr** downstream:

M_CTL, M_DATA, M_FLUSH, M_IOCTL, M_PROTO

SLIP Module

The Serial Line Internet Protocol (**slip**) line discipline enables the TCP/IP protocol layer to use the serial lines as network interfaces.

The **slip** module is used to read internet protocol (IP) packets from a tty port, and to write IP packets to a tty port for an IP network interface.

The **slip** module reads raw data character by character from a tty port until it can assemble an IP packet and forward it for the IP network interface. In the same way, packets written to a network interface and corresponding to a tty device (**slip** interface) are written to the tty port for transmission to a remote system.

SLIP Applications

Two **slip** applications are provided: the **slattach** command and the **sliplogin** command. Both attach a tty device to a network interface.

Both the **slattach** and **sliplogin** commands configure the stream head for a tty in the following fashion:

The tty port is opened, and set to “raw” mode, with all echoing disabled. Then all active disciplines are popped from the stream for the tty. The **slip** line discipline is pushed using the **strload** command. For more information on **slip** configuration, see the **/etc/rc.net.serial** shell script.

The stream stack with the **slip** module is reduced to:

- The stream head
- The **slip** module
- The driver.

SLIP Routines

The **slip** module has open and close routines, and put and service routines for the read-side and the write-side.

The read-side put routine handles **M_FLUSH**, **M_CTL** and **M_DATA** messages sent by the driver. Unrecognized messages are just passed upstream.

The write-side put routine handles **M_FLUSH**, **M_IOCTL** and **M_DATA** messages sent by the stream head. Unrecognized messages are just passed downstream.

TTY Drivers

tty drivers directly control the hardware (asynchronous devices) or pseudo-hardware (pty devices). They perform the actual input and output to the adapter. The tty subsystem consists of the following drivers:

Name	Function
cxma	128-port adapter
lft	Low function terminal emulation
lion	64-port adapter
pty	Pseudo-terminal
rs	Native serial ports, 8-port and 16-port adapters.
sf	Native serial ports, in ISA bus hardware.

This section first explains how to provide a configuration routine for a tty driver. Then the open, close, write and read processings are described, as well as the open and pacing disciplines which are managed by the driver.

STREAMS tty device drivers have interrupt entry points at the hardware device interface, and have direct entry points only for the **open**, **close**, and **sysconfig** system calls. These entry points are accessed via STREAMS. The stream head translates **write** and **ioctl** calls into messages and sends them downstream to be processed by the driver write put routine.

Drivers Configuration Routine

In order to support dynamic loading, unloading, configuring, and unconfiguring, each tty driver must provide a configuration routine. This routine is called each time the tty driver is referenced in a load or unload operation.

Unlike an AIX non-STREAMS-based character device driver, the configuration entry point of a tty driver does not add an entry to the device switch table. It simply declares itself to STREAMS by calling the **str_install** utility as shown in following example. The **str_install** utility performs the internal operations necessary to add or remove the tty driver from the STREAMS internal tables.

Example

The following example is a minimal configuration routine for a tty driver called **ttyd**. Device specific configuration and initialization logic can be added as necessary. The **ttyd_config** entry point defines and initializes the **strconf_t** structure required by the **str_install** utility. In this example the **ttyd_config** operation retrieves the argument passed through the pointer specified by the **uiop** parameter. The major number is required for tty drivers and is retrieved from the **dev** parameter.

```
/* ttyd driver example:
/* BEGINNING. */

#include <sys/device.h>          /* for the CFG_* constants */
#include <strconf.h>            /* for the STR_* constants */

static struct module_info ttydm_info = {
    DRIVER_ID, "ttyd", 0, INFPSZ, 512, 256
};

static struct qinit ttyd_rinit = {
    NULL, ttydrsrv, ttydopen, ttydclose, NULL, &ttydm_info, NULL
};

static struct qinit ttyd_winit = {
    ttydwput, ttydwsrv, NULL, NULL, NULL, &ttydm_info, NULL
};

static struct streamtab ttydinfo = {
    &ttyd_rinit, &ttyd_winit, NULL, NULL
};

ttyd_config(dev, cmd, uiop)
    dev_t    dev;
    int      cmd;
    struct   uio *uiop;

{
    struct   ttyd_dds    tmp_dds;
    static   strconf_t   ttyd_conf = {
        "ttyd", &ttydinfo, STR_NEW_OPEN, -1
    };
};
```

```

        if (uimove(&tmp_dds,sizeof(struct ttyd_dds),UIO_WRITE, uiop))
            return EFAULT;

        ttyd_conf.sc_major = major(dev);
        ttyd_conf.sc_sqllevel = SQLVL_QUEUEPAIR; /* for example */

        switch (cmd) {
        case CFG_INIT: return str_install(STR_LOAD_DEV, &ttyd_conf);
        case CFG_TERM: return str_install(STR_UNLOAD_DEV, &ttyd_conf);
        default: return EINVAL;
        }
    }

/*ttyd driver example:
/* END.*/

```

In this example, all the **tttyd** driver entry points are declared in the **qinit** structures that will be handled by STREAMS. These entry points are:

On the read-side:

- an **open** routine
- a **close** routine
- a **service** routine.

On the write-side:

- a **put** routine
- a **service** routine.

All these entry points are standard STREAMS interface entry points. The only direct entry point to the driver is **tttyd_config**.

Open Disciplines

Open disciplines specify the protocol to establish a connection

Open disciplines are defined at configuration time. They can be **dtropen** (default value) or **wtopen**.

Pacing Disciplines

Pacing disciplines (or flow disciplines) control the flow of input and output data. The flow control can be software or hardware, depending on driver and adapter capabilities. The default modes are **xon** for a terminal and **dtr** for a printer.

Software Flow Control

In this case, special characters (START and STOP) indicate when the flow of data has to be stopped or resumed. **IXON**, **IXOFF** and **IXANY** flags of **termios** structure enable start and stop output or input control.

Hardware Flow Control

In this case, the flow of data is suspended or resumed by toggling an EIA modem control signal. The hardware flow control modes are: **dtr**, **rts**.

Open and Close Routines

The open routine is called whenever a device is opened. The open routine uses the open discipline specified at configuration time, and pins the private data structure.

The close routine unpins the private data structure.

During open and close, the driver is in user context such that it is able to issue sleeps and allocate dynamic memory.

The open and close routines are normally serialized by the STREAMS. Only one close can be active at a time per major/minor device pair. In some cases, the open routine can also sleep waiting for the specified conditions, such as modem status, to be met.

Write-Side Put Routine

The write-side put routine processes the following received messages:

M_BREAK, **M_CTL**, **M_DATA**, **M_DELAY**, **M_FLUSH**, **M_IOCTL**, **M_PROTO**, **M_START**, **M_STARTI**, **M_STOP**, **M_STOPI**.

Because the tty driver is the lowest module on the STREAMS stack, all other messages have to be freed by the driver.

M_DATA

Data are to be sent to output.

M_DELAY

Data output is suspended for a delay passed as a parameter.

M_IOCTL

The **M_IOCTL** message contains the IOCTL command to be processed by the driver. When the processing ends, the driver sends an **M_IOCACK** or **M_IOCNAK** message upstream, depending on the result of the IOCTL processing. For those IOCTLs which are not to be processed by the driver, the driver sends an **M_IOCNAK** message upstream. However some IOCTL commands are ignored by the driver and cause only an acknowledgement of the command.

M_FLUSH

If the request is done for the write side, the driver directly flushes its write-side queue. If the request is for the read side, the driver immediately sends the same message upstream and finally flushes its read-side queue.

M_CTL

The **M_CTL** message is generally sent by **ldterm**, **sptr** or **tioc** modules on their **open** processing to obtain information from the driver. It must be processed immediately. The message contains one of the following commands:

MC_CANONQUERY, **TIOCGETA**, **TIOCGETMODEM**, **TIOC_REQUEST**, **TIOCMGET**, **TXTTYNAME**.

If these commands need to be processed by the driver, the driver sends upstream the same message containing the requested information in the data part of the message.

If these commands do not need to be processed by the driver, then the driver just frees the message.

In reply to the **TIOC_REQUEST** command the driver sends the **M_CTL** message upstream replacing **TIOC_REQUEST** by **TIOC_REPLY**, and adding the list of its specific ioctl commands, if any.

In reply to the **MC_CANONQUERY** command, the driver answers **MC_CANON_PART** if needed.

TXTTYNAME is used to initialize the line discipline module name.

TIOCGETA is used to initialize the line discipline with the default **termios** settings.

M_STOP, M_START

These requests stop and restart output.

M_STOPI, M_STARTI

These requests stop and restart input.

M_BREAK

This message is sent by the line discipline to the driver to request the transmission of a **BREAK** on the device if it supports the break condition.

If the integer pointed by the `b_rptr` part of this message is 1, then the driver sets the break condition; otherwise the driver clears the break condition.

M_PROTO

This message is sent by the **sptr** module and contains the **LPWRITE_REQ** command. The **M_DATA** message associated with the **M_PROTO** message contains data to transmit to the line. The driver must acknowledge this message with a **M_PCPROTO** message with parameter **LPWRITE_ACK** when it has transmitted all the data to the line.

Read-Side Processing

The driver has no read-side put routine because it is the last module on the stream. However the driver has a read-side service routine which is scheduled by **STREAMS**.

The following messages are sent upstream by the driver:

M_BREAK, **M_CTL**, **M_DATA**, **M_PCPROTO**, **M_PCSIG**.

M_DATA

When the driver is ready to send data or other information to the user process, it does not wake up the process. It stores the received characters in **M_DATA** messages which are queued. These messages are sent later by the read-side service routine to the stream corresponding to the driver line.

M_BREAK

The driver sends this message upstream to provide the line discipline with the following status information:

`break_interrupt`, `parity_error`, `framing_error` or `overrun`.

M_CTL

The driver sends this message upstream either to communicate changes in the modem status (a transition of the carrier state from `on` to `off` for example), or to answer to a previous request (**TIOCGETA** from module **ldterm**, for example).

M_PCSIG

If the **SAK** recognition is set, the driver sends this message to signal the reception of the **SAK** sequence.

M_PCPROTO

The driver sends this message containing the **LPWRITE_ACK** command to indicate to **sptr** that all the data it sent in the previous **M_PROTO** message was transmitted to the line.

Interface with the TIOC Module

On opening, the **tioc** module sends an **M_CTL** message downstream containing the **TIOC_REQUEST** command. Then the **tioc** module waits for an **M_CTL** message with a **TIOC_REPLY** command containing the downstream modules or driver specific ioctls. The **tioc** module will update its ioctls table according to these specific ioctls. (See “TIOC Module” on page 10-4.)

Example

The **rs** driver has two transparent IOCTLs (**RS_SETA** and **RS_GETA**) which require data transfers from and into user space. The **tioc** module will perform these transfers on behalf of the **rs** driver. For that, **rs** defines two **tioc_reply** structures that will be sent to **tioc** module at its open time in reply to a **TIOC_REQUEST** included in an **M_CTL** message.

```
/* tioc_reply structures array. */
static struct tioc_reply
srs_tioc_reply[] = {
    { RS_SETA, sizeof(struct rs_info), TTYPE_COPYIN },
    { RS_GETA, sizeof(struct rs_info), TTYPE_COPYOUT },
};
```

In the write put routine of the **rs** driver, if the message to process is an **M_CTL** containing a **TIOC_REQUEST** command, the following code is executed:

```
/* mp is the M_CTL to process,
   iocp is mp->b_rptr: pointer to an iocblk struct describing M_CTL
   mpl will contain the 2 tioc_reply structures.
   q is the driver's write-side queue.
*/

case TIOC_REQUEST:

int reply_size = 2 * sizeof(struct tioc_reply);
iocp->ioc_cmd = TIOC_REPLY;
if (!(mpl = allocb(reply_size, BPRI_MED)))
    break; /* just reply with the same message, next RS_SETA and
           RS_GETA will arrive transparent and will fail */
iocp->ioc_count = reply_size;
bcopy(srs_tioc_reply, mpl->b_rptr, reply_size);
mpl->b_wptr = mpl->b_rptr + reply_size;
mp->b_cont = mpl;
qreply(q,mp); /* send the M_CTL upstream */
break;
}
```

The **M_CTL** message containing the **TIOC_REPLY** command is sent upstream and other modules will add **M_DATA** messages to **M_CTL** if necessary.

Interface with the LDTERM Module

The driver answers to the following **ldterm** commands which are included in **M_CTL** messages:

TIOCGETA The **ldterm** module asks for current **termios** structure settings.

TIOCGETMODEM

The **ldterm** module asks for the modem carrier state.

MC_CANONQUERY

The **ldterm** module negotiates which **termios** structure flags are handled by the driver.

Another specific interface between the tty driver and the line discipline module (**ldterm**, **sptr**) module is the possibility for the driver to send special information:

- A modem status change
- A break interrupt
- A parity error
- A framing error.

When the driver detects a modem status change, it sends upstream an **M_CTL** message with a status pointed by the `b_rptr` part of the message. This status can be:

`cts_on, cts_off, dsr_on, dsr_off, ri_on, ri_off, cd_on, cd_off.`

When the driver detects an error, it sends an **M_BREAK** message upstream with a status pointed by the `b_rptr` part of the message. This status can be one of the values:

`break_interrupt, framing_error, parity_error, overran.`

Interface with the SPTR Module

When the driver has sent on the line all the data associated with an **M_PROTO** message (with the **LPWRITE_ACK** command), it sends **sptr** an **M_PCPROTO** message (with the **LPWRITE_REQ** command) in reply.

The TTY Subsystem in a Multiprocessor Environment

Note: Information supplied in this section requires knowledge of the STREAMS synchronization, and of the device drivers in a multiprocessor environment. See Related Information on page 10-28.

On a multiprocessor system, a program can run on any processor and can migrate between processors. A program which consists of multiple threads can run on several processors at the same time. This creates a problem of concurrent access to global data.

However, the STREAMS-based tty subsystem takes advantage of the synchronization provided by STREAMS. Most of tty STREAMS modules and drivers are configured with the queue pair level synchronization (**SQLVL_QUEUEPAIR**). This ensures the serialization of operations done on read and write sides, without explicit locking by the module.

In the tty subsystem, the problem of maintaining data consistency in a multiprocessor environment is different in the driver and in the other stream modules.

TTY Modules Other Than Driver

If the module doesn't have global data shared by all the modules instances, the queue pair level synchronization (**SQLVL_QUEUEPAIR**) is enough to ensure the module is multiprocessor-safe. Nevertheless, the queue level synchronization (**SQLVL_QUEUE**) can be used for better throughput if there is no shared data between the read and write sides of the module, or if the put and service routines of the module guarantee the consistency of accesses to such data. The kernel provides a set of locking services and atomic primitives for that purpose.

If the module has global data, it must be protected with locks, such as simple locks. Use the **disable_lock** and **unlock_enable** kernel services to safely protect data on a multiprocessor system.

Drivers

There are three basic types of critical sections for drivers:

- thread-thread: critical sections shared between threads
- thread-interrupt: critical sections shared between threads and interrupt handlers
- interrupt-interrupt: critical sections shared between interrupt handlers.

STREAMS resolves the first critical section problem, except for concurrent multiple opens. The serialization of open and close, and the synchronization between the last close and first open is performed by STREAMS, and no specific lock is needed in the driver.

The other critical sections must be managed by the driver.

Depending on the STREAMS synchronization level selected, the STREAMS will also ensure the serialization of the put and service routines for the driver.

For the **rs** driver for example, only one interrupt per adapter can be handled at a time. But off-level routines which may be called for the same port on different processors, need to be serialized. To do so, the driver uses simple locks and interrupt priority masking.

Drivers which are not multiprocessor-safe can rely on the possibility of funnelling provided by the kernel and STREAMS.

Special Cases

These include callback functions and driver configuration routines.

Callback Functions

There are callback functions (for the **timeout** or **bufcall** utilities) that need to be protected against interrupts. This protection can be ensured by STREAMS. In this case the flag **STR_QSAFETY** will be specified in the **str_install** utility. It is also possible to use locks to protect the callback functions.

Driver Configuration Routine

The minimal configuration routine of the driver is not called by STREAMS. Consequently it is not protected by STREAMS and must ensure its protection as a non-STREAMS driver does.

An example of a driver configuration routine is provided on page 10-14.

IOCTL Support and Origin

The following table indicates for each IOCTL:

- Its origin: AIX, AT&T, BSD or SVID (AIX means that the IOCTL is specific to AIX system).
- In which part of the stream tty subsystem (stream head, **ldterm**, **sptr**, **nls**, or driver) the IOCTL is processed. An asterisk (*) in a column indicates where the ioctl is processed.

The last column (comment) gives an additional information for some IOCTLs:

- “STREAMS” indicates that the ioctl is processed by the STREAMS framework.
- “TIOCGETA”, “TIOCSETA”, “TIOCSETAF” and “TIOCSETAW” indicate the internal names of the TCGETS, TCSETS, TCSETSF and TCSETSW ioctls respectively.
- “rs”, “sf”, “lion”, and “cxma” indicate that the IOCTL can only be used by the specified driver.
- “pty” indicates that the IOCTL can only be used by a pseudo-terminal driver.

IOCTL	origin	stream head	ldterm	sptr	nls	driver	comment
CXMA_GETA	AIX					*	cxma
CXMA_SETA	AIX					*	cxma
CXMA_SETAW	AIX					*	cxma
CXMA_SETAF	AIX					*	cxma
CXMA_KME	AIX					*	cxma
CXMA_GETFLOW	AIX					*	cxma
CXMA_SETFLOW	AIX					*	cxma
CXMA_GETAFLOW	AIX					*	cxma
CXMA_SETAFLOW	AIX					*	cxma
CXMA_RESET	AIX					*	cxma
EUC_WGET	AT&T		*				
EUC_WSET	AT&T		*				
FIOASYNC	BSD	*					STREAMS
FIONREAD	BSD	*					STREAMS
LI_GETVT	AIX					*	lion
LI_SETVT	AIX					*	lion
LI_GETXP	AIX					*	lion
LI_SETXP	AIX					*	lion
LI_SLPI	AIX					*	lion
LI_DSLP	AIX					*	lion
LI_SLPO	AIX					*	lion
LI_PRES	AIX					*	lion, cxma
LI_DRAM	AIX					*	lion, cxma
LI_GETTBC	AIX					*	lion

IOCTL	origin	stream head	ldterm	sptr	nls	driver	comment
LI_SETTBC	AIX					*	lion
LPRGET	AIX			*			
LPRSET	AIX			*			
LPRMODG	AIX			*			
LPRMODS	AIX			*			
LPRGOTV	AIX			*			
LPRSTOV	AIX			*			
LPQUERY	AIX			*			
LPRGETA	AIX			*			
LPRSETA	AIX			*			
LPWRITE_REQ	AIX			*			
RS_GETA	AIX					*	rs, sf
RS_SETA	AIX					*	rs, sf
TCFLSH	SVID	*					
TCGETA	SVID		*	*			
TCGETS	SVID		*			*	TIOCGETA
TCGETX	SVID		*			*	
TCGMAP	AIX				*		
TCKEP	AIX	*					STREAMS
TCLOOP	AIX					*	
TCQSAK	AIX					*	
TCQTRUST	AIX	*					
TCSAK	AIX					*	
TCSBRK	SVID					*	
TCSBREAK	AIX					*	
TCSETA	SVID		*	*		*	
TCSETAF	SVID		*	*		*	
TCSETAW	SVID		*	*		*	
TCSETS	SVID		*			*	TIOCSETA
TCSETSF	SVID		*			*	TIOCSETAF
TCSETSW	SVID		*			*	TIOCSETAW
TCSETX	SVID		*			*	
TCSETXF	SVID		*			*	
TCSETXW	SVID		*			*	
TCSMAP	AIX				*		
TCTRUST	AIX	*					
TCXONC	SVID	*					
TIOCCBRK	BSD		*				

IOCTL	origin	stream head	ldterm	sptr	nls	driver	comment
TIOCCDTR	BSD					*	
TIOCCONS	AIX	*					
TIOCEXCL	BSD					*	pty
TIOCFLUSH	BSD		*				
TIOCGETC	BSD		*				
TIOCGETD	BSD		*				
TIOCGETP	BSD		*				
TIOCGLTC	BSD		*				
TIOCGPGRP	SVID	*					
TIOCGSID	SVID	*					
TIOCGWINSZ	BSD		*				
TIOCHPCL	BSD		*				
TIOCIXCL	BSD					*	pty
TIOCLBIC	BSD		*				
TIOCLBIS	BSD		*				
TIOCLGET	BSD		*				
TIOCLSET	BSD		*				
TIOCMBIC	SVID					*	
TIOCMBIS	SVID					*	
TIOCMGET	SVID		*			*	
TIOCMSET	SVID		*			*	
TIOCOUTQ	BSD		*			*	
TIOCPKT	BSD					*	pty
TIOCREMOTE	AT&T					*	pty
TIOCSBRK	BSD		*				
TIOCSDTR	BSD					*	
TIOCSETC	BSD		*				
TIOCSETD	BSD		*				
TIOCSETN	BSD		*				
TIOCSETP	BSD		*				
TIOCSLTC	BSD		*				
TIOCSPPGRP	SVID	*					
TIOCSTART	BSD	*					
TIOCSTI	BSD		*				
TIOCSTOP	BSD	*					
TIOCSWINSZ	BSD		*			*	pty
TIOCUCNTL	BSD					*	pty
TXGPGRP	AIX	*					

IOCTL	origin	stream head	ldterm	sptr	nls	driver	comment
TXISATTY	AIX	*					
TXSETIHOOG	AIX		*				
TXSETOHOOG	AIX		*				
TXSPGRP	AIX	*					
TXTTYNAME	AIX					*	

TTY Data Structures

The following is an extract from the `usr/include/sys/str_tty.h` file. This file defines all the constants, functions, structures and types of messages that are used by the tty modules and drivers for the exchange of messages. Additional comments explain how some of the structures are used.

```
#ifndef _H_STR_TTY
#define _H_STR_TTY

#include <sys/termio.h>
#include <sys/stream.h>
#ifdef _KERNEL
#include <sys/errno.h>
#include <sys/ioctl.h> /* for TTNAMEMAX definition */
#include <termios.h>
#include <sys/trchkid.h>
#include <sys/atomic_op.h>

/* Macro definition for atomic operation on sysinfo fields. */
#define sysinfo_add(x,y) fetch_and_add(&(x), (y))
```

TIOC Module

```
/* Commands of the M_CTL message at open */

#define TIOC_REQUEST _IO('J', 0x91) /* request for ioctls */
#define TIOC_REPLY _IO('J', 0x92) /* reply for ioctls */

/* Data structures for TIOC_REPLY */

struct tioc_reply {
    int tioc_cmd; /* ioctl command */
    int tioc_size; /* number of bytes to copy */
    int tioc_type; /* type of ioctl */
};

/* Values for tioc_type */
```

tioc_type indicates if the IOCTL requires data to be copied into or from user space.

```
#define TTYPE_NOCOPY 0 /* don't need any copy */
#define TTYPE_COPYIN 1 /* need an M_COPYIN */
#define TTYPE_COPYOUT 2 /* need an M_COPYOUT */
#define TTYPE_COPYINOUT 3 /* need both M_COPYIN and M_COPYOUT */
#define TTYPE_IMMEDIATE 4 /* use immediate value */
```

TTY Commands Associated with M_CTL Messages

```
#define TIOCGETMODEM    _IO('J', 0xa0) /* get the modem state from driver */
#define MC_CANONQUERY   _IO('J', 0xa1) /* query the termios state */
#define MC_NO_CANON     _IO('J', 0xa2) /* set pty's remote mode */
#define MC_DO_CANON     _IO('J', 0xa3) /* reset pty's remote mode */
#define MC_PART_CANON   _IO('J', 0xa4) /* oflag,iflag and lflag of termios*/
/* are handled by driver or ldterm */
```

M_PCPROTO Commands

```
#define LPWR              ('l' << 8)
#define LPWRITE_ACK      (LPWR | 31) /* command in M_PCPROTO message */

typedef int      OSR_STATUS;
```

Status Information

enum status gives general status definitions for drivers and line discipline in the case of parity and framing error or break_interrupt from the adapter, or in the case of modem status changes. The status information is sent by the driver in an **M_BREAK** message.

```
enum status {
    good_char, overrun, parity_error, framing_error, break_interrupt,
    cts_on, cts_off, dsr_on, dsr_off, ri_on, ri_off, cd_on, cd_off };
```

The possible values for the status enumeration are:

- good_char** A valid character was received.
- overrun** Characters were not removed from the hardware in a timely fashion and some data was lost by the hardware.
- parity_error** Character was received with improper parity. The character is passed as received by the hardware.
- framing_error** Character was received with a framing error. This usually indicates the number of bits per character is set incorrectly or the baud rate is not set correctly. The character is passed as received by the hardware.
- break_interrupt** The hardware detected a break condition. The break condition is different for various adapters and physical link layers. For asynchronous communications, the break condition is usually defined as a spacing condition on the line for more than one total character time.
- cts_on** The clear to send signal, **cts**, made a low to high transition.
- cts_off** The clear to send signal made a high to low transition.
- dsr_on** The data set ready signal, **dsr**, made a low to high transition.
- dsr_off** The data set ready signal made a high to low transition.
- ri_on** The ring indicate signal, **ri**, made a low to high transition.
- ri_off** The ring indicate signal has made a high to low transition.
- cd_on** The data carrier detect signal, **cd**, made a low to high transition.
- cd_off** The data carrier detect signal made a high to low transition.

TTY Trace Support

The main tty hook IDs are defined in the **sys/trchkid.h** common header file. The tty subhooks are defined below. The tty hooks and subhooks are used in the **Return** and **Enter** macros defined below.

```

/* TTY subhooks identifiers */

#define TTY_CONFIG      0x01
#define TTY_OPEN        0x02
#define TTY_CLOSE       0x03
#define TTY_WPUT         0x04
#define TTY_RPUT         0x05
#define TTY_WSRV         0x06
#define TTY_RSRV         0x07
#define TTY_REVOKE      0x08 /* for stream head */
#define TTY_IOCTL        0x09 /* for ioctls */
#define TTY_PROC          0x0a /* for drivers */
#define TTY_SERVICE      0x0b /* for drivers */
#define TTY_SLIH          0x0c /* for drivers */
#define TTY_OFFL          0x0d /* for drivers */
#define TTY_LAST          0x0e /* can be used for any specific entry*/

```

The parameters for the **Return** and **Enter** macros are:

w	TTY hookid TTY subhookid
dev	dev(type dev_t)
ptr	address of the private data (q->q_ptr) of each module or driver
a, b, c	specific parameters for each subhook
retval	return value

The parameters for each subhook are:

TTY_CONFIG	a=command, no dev, no ptr.
TTY_OPEN	a=oflag, b=sflag
TTY_CLOSE	a=flag
TTY_WPUT	a=@msg, b=message type
TTY_RPUT	as TTY_WPUT
TTY_WSRV	a=q_count
TTY_RSRV	as TTY_WSRV
TTY_REVOKE	a=flag
TTY_IOCTL	a=ioctl command
TTY_PROC	a=cmd, b=arg
TTY_SERVICE	a=service command, b=arg
TTY_SLIH	no dev, ptr=@struct intr, a=adapter type
TTY_OFFL	ptr=@struct intr

Enter and Return Macros

The **Enter** and **Return** macros use the tty hooks and subhooks.

```

#define Enter(w, dev, ptr, a, b, c) \
    dev_t   DEV; \
    int     PTR; \
    int     Flag = 0; \
    int     W; \
    if (TRC_ISON(0)) { \
        DEV = (dev); \
        PTR = (ptr); \
        Flag = 1; \
        W = w; \
        TRCHKGT(W, DEV, PTR, a, b, c); \
    }

```

```

#define Return(retval) {
    int    RET = (retval);
    int    Line = __LINE__;
    if (Flag)
        TRCHKGT(((W)|0x80), DEV, PTR, RET, Line, 0);\
    return(RET);
}

#define Returnv(retval) {
    int    Line = __LINE__;
    if (Flag)
        TRCHKGT(((W)|0x80), DEV, PTR, 0, Line, 0);\
}

#define Data(xxx, a, b, c)
    (Flag ? TRCHKGT(((W)|0x40), DEV, PTR, a, b, c) : 0)

```

TTY IOCTL Internal Names

The following definitions indicate the internal names of some IOCTLs.

```

/* tty ioctls commands internal names */

#define TIOCGETA    TCGETS
#define TIOCSETA    TCSETS
#define TIOCSETAW   TSETSW
#define TIOCSETAF   TCSETSF

```

STREAMS TTY Modules and Drivers DDS

```

/* Maximum device name length */

#define DEV_NAME_LN    16

/* DDS types used at configuration time */

enum dds_type {
    LC_SJIS_DDS,    /* Which DDS type */
    LDTERM_DDS,    /* sjis lower converter module*/
    LION_ADAP_DDS, /* ldterm module */
    LION_LINE_DDS, /* 64-port driver (for adapter) */
    NLS_DDS,       /* 64-port driver (for lines) */
    PTY_DDS,       /* nls module */
    RS_ADAP_DDS,   /* pty module */
    RS_LINE_DDS,   /* Native 8 and 16-port driver (for adapters) */
    SPTR_DDS,      /* Native 8 and 16-port driver (for lines) */
    TIOC_DDS,      /* sptr module */
    UC_SJIS_DDS,   /* tioc module */
    CXMA_ADAP_DDS /* sjis upper converter module */
    CXMA_LINE_DDS /* 128-port driver (for adapter) */
    CXMA_LINE_DDS /* 128-port driver (for lines) */
};

```


STREAMS TTY Modules and Drivers Names

```
enum module_names {
    tioc,          /* transparent ioctl modulename */
    ldterm,       /* line discipline module name */
    pty,          /* pseudo tty driver module name */
    uc_sjis,      /* upper converter sjis module name */
    lc_sjis,      /* lower converter sjis module name */
    nls,          /* mapping discipline module name */
    sptr,         /* serial line printer discipline name */
    rs,           /* rs driver name */
    lion,         /* lion driver name */
    cxma         /* 128 port driver name */
};

#endif /* _H_STR_TTY */
```

Related Information

Discussion of Multiprocessing (MP) Serialization in *Serialization Services*, on page 5-8, and *MP-Safe Coding Example*, on page 5-10.

The TTY Subsystem in *AIX Version 4.1 General Programming Concepts*, *Volume 1: Writing Programs*.

STREAMS in *AIX Version 4.1 Communications Programming Concepts*

Understanding Multiprocessor-Safe Device Drivers in *AIX Version 4.1 Kernel Extensions and Device Support Programming Concepts*

UNIX System V Release 4, Programmer's Guide: STREAMS, Englewood Cliffs, N.J.: Prentice-Hall, Inc., 1990.

Chapter 11. Implementing Graphical Input and 2D Graphics Device Drivers

This chapter provides technical guidance for developers who want to add functionality into the XServer on AIXwindows. This additional functionality may include graphics adapters, input devices and dynamically loaded extensions.

This chapter has several sections:

- “Porting to the AIXwindows X Server: Overview,” on page 11-1, provides a description of the basics of porting to the AIXwindows X Server.
- “Porting 2D Graphics Adapters,” on page 11-2, outlines the general procedure to develop a graphics adapter device driver and the loadable DDX interface used by the X Server.
- “Porting Input Devices,” on page 11-48, describes how to add a new input device to the AIXwindows X Server via the X11 Input Extension.
- “Building a Dynamically Loadable Module,” on page 11-62, is a brief description of how to develop a generic X extension load module.

The information provided in this chapter does not attempt to educate readers about the X Window System or X programming. Rather, it describes the tasks required to implement the above-mentioned functionality in the AIX system. For general information about the X Window System programming see the list of related information at the end of this chapter.

Porting to the AIXwindows X Server: Overview

You can add functionality to the AIXwindows X server through the creation of dynamically loadable modules. Each load module must provide the following basic functionality:

- An entry point definition
- A set of routines to initialize the appropriate data structure provided by the X server
- X Consortium and vendor-written implementation-specific routines

The load modules can interact with the X Window System in a variety of ways. The AIXwindows X server currently supports the following functionality:

- Porting the X server to run on a different 2D graphics adapter
- Writing load modules to the standard X11 Input Extension
- Creating your own dynamically loadable X Window System extension

The load modules for 2D graphics adapters and extension input devices require that entries be defined and configured in the AIX Object Data Manager (ODM) databases. This allows the X server to determine where the load modules for the specific devices reside.

Information about 2D graphics adapters is stored in the ODM Graphics Adapter Interface (GAI) database, and information about extension input devices is stored in the ODM XINPUT database.

Current database entries can be examined by setting the environment variable **ODMDIR=/usr/lib/objrepos** and running the command `odmget GAI` or `odmget XINPUT`. These entries will be explained in greater detail in the respective sections of this chapter.

The dynamically loaded extensions are normally loaded via a per-user basis. If the X server is started with a `-x ext_name` flag, then the X server will look in the file **`/usr/lpp/X11/bin/dynamic_ext`** for the load module that corresponds to `ext_name`. There are also static extensions, which can be found in **`/usr/lpp/X11/bin/static_ext`**, but they are always loaded into the X server, and are thus not recommended.

The information provided here assumes an ability on the part of the developer to program generic actions without further instruction. Refer to the following sections to learn more about extension-specific tasks:

- Porting 2D Graphics Adapters (on page 11-2)
- Porting Input Devices (on page 11-48)

The subroutines used in these extensions are summarized in “List of X Server Porting Subroutines” on page 11-64.

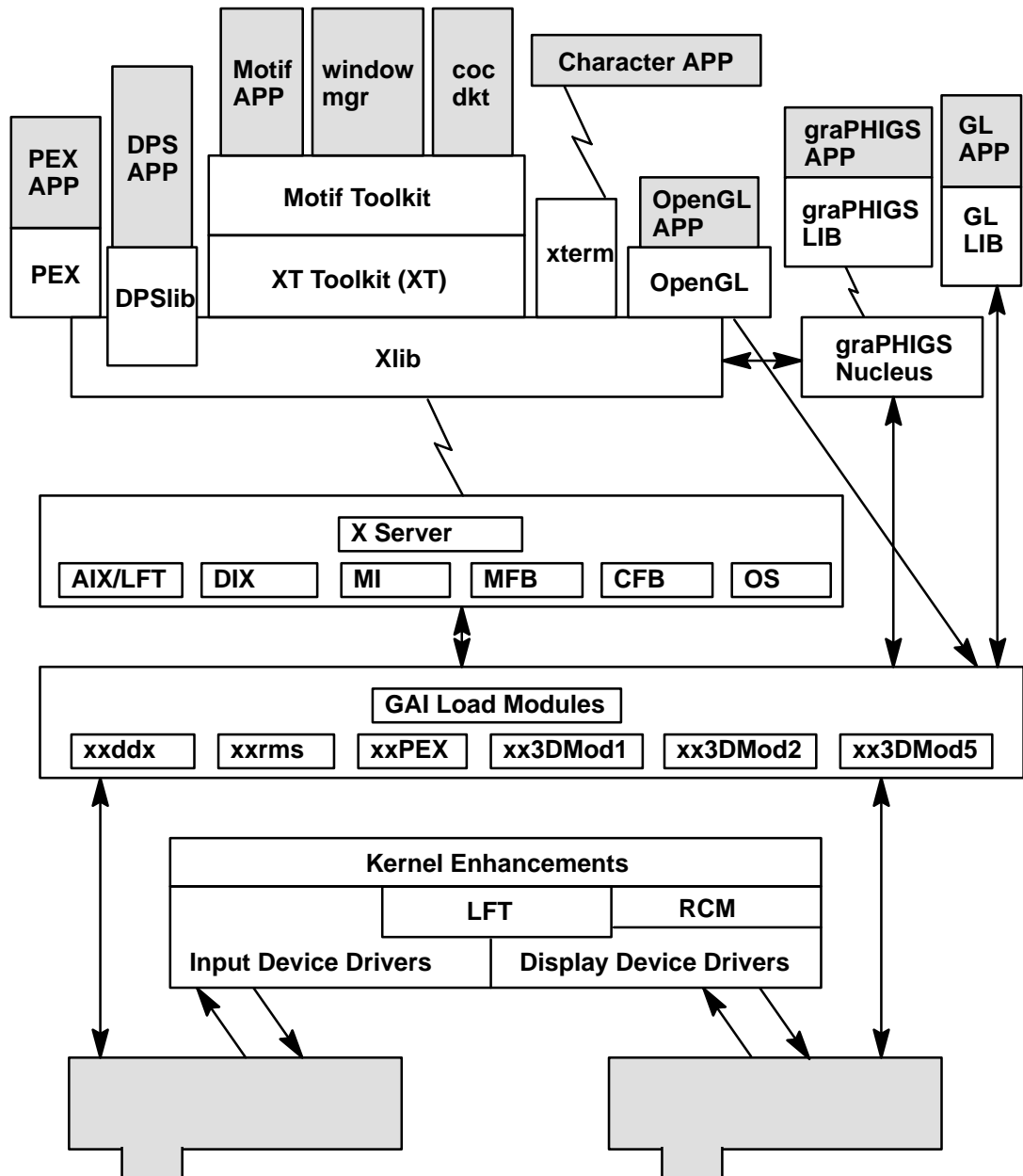
Porting 2D Graphics Adapters

This section describes the method for porting 2D graphics adapters into the X server on the AIXwindows. The information is divided as follows:

- Graphics Architecture Overview
 - Graphics Adapter Interface (GAI) Display Subsystem
 - X Server
- Low-Level Display Driver
 - Display Device Driver
 - Display Device Driver Subroutines
 - LFT Interface Routines
 - Display Driver Structure Descriptions
- Device Dependent Driver (DDX)
 - GAI 2D Adapter Load Modules
 - X Server Initialization Routines
 - Device Dependent Initialization Subroutines
 - Adapter Access and the `aixgsc` System Call
 - Minimum RMS for 2D Adapters
 - Configuring the 2D Adapter into the ODM Database

Graphics Adapter Interface (GAI) Display Subsystem

The Architecture of the Display Subsystem figure illustrates the display subsystem architecture.



Architecture of the Display Subsystem

The display subsystem is the set of software entities that provides all display related input and output to applications and to the AIX Version 4.1 kernel. Its capabilities range from simple 2D graphics to support for 3D models which enable lighting and shading.

A Graphics Adapter Interface (GAI) defines an interface between the user application programming interface (API) and the device-specific code.

The principal components of the architecture include:

- X Server
- graPHIGS Nucleus
- Kernel Enhancements
- GAI Load Modules

The shaded objects in the architecture figure show input and graphics hardware devices and user applications. These are not components of the display subsystem.

The components of the display subsystem receive their external inputs from X clients and from other kernel subsystems. User inputs are processed by the kernel and routed through the X server, which sends them appropriately to X clients.

Applications use programming interfaces made available through the various libraries of the display subsystem. Each programming interface, (API), defines a graphical model. These graphical models evolve in many forums, including the American National Standards Institute, the International Standards Organization, the X Consortium, the GL Architecture Review Board (ARB), and various highly successful industry *de-facto* standards.

The graphical models supported by the display subsystem are listed below:

- graPHIGS
- X Window System
- Display PostScript
- Graphics Library of Silicon Graphics, Inc. (GL)
- OpenGL Version 1.0
- PEX 5.1

Display Subsystem Definitions

The display subsystem is the set of software units that provides all display-related input and output to applications and to the AIX Version 4.1 kernel.

List of Component Types

The types of components are defined below:

Component Type	Definition
Adapter	With reference to input devices, an adapter specifies the device driver used to process user events from one of the input devices. This could be the keyboard/sound, mouse, tablet, dials, or lighted programmable function keys.
Application	Specifies a program that uses the display subsystem. The use of this term with respect to the display subsystem is analogous to the standard definition of application.
Client	Specifies an application that uses the services of X Window System. In the terminology of the display subsystem, however, this concept has been extended. A client is an application that uses a server. These clients are typically different processes and may reside on different platforms.
Display Device Driver	Specifies the set of subroutines that control read or write access to a

display adapter. The display drivers provide a common programming interface to other components in the display subsystem, and account for the particular function and implementation within the specific device. Thus, each display driver contains code unique for that device. A set of subroutines is not considered a display driver unless it:

- Contains software that is applicable only to the function or control of a particular device.

OR

- Issues I/O level commands to specifically manipulate the device hardware.
- Manages the graphics context and graphics resources on behalf of servers.

Library	Specifies a programming interface consisting of a functionally related set of subroutines that are typically grouped as one module and that provide the formal application programming interface into a collection of software.
Nucleus	Specifies a server. The term nucleus is used principally to distinguish an instance of a server process from the X server. The principal user of nucleus is the graPHIGS API, via the graPHIGS nucleus.
Server	Specifies the component of X Window System that receives requests from the Xlib library and that returns events to the Xlib library. This concept has been modified and extended in the display subsystem. A server can also mean a process executing in a platform in which rendering occurs and that is directly attached to the display adapter. This server accepts X protocol requests and returns X protocol events, and as such is the X server. However, other graphical models use the server to obtain screen resources such as window geometries or colormaps, so the X server definition is now extended to become the resource server.
Toolkit	Specifies a library of calls that allows applications to manipulate data presentation using a different level of abstraction than that offered by an underlying library. The display subsystem expansion of this concept of a toolkit is the concept that it in turn calls upon another library. The toolkit often takes its name from the underlying library. Frequently, toolkits are written to be window system independent.

List of Component Names

The following list introduces the components and provides an description of the function of each. Where appropriate, the general category of component is described rather than each instance. Some of the components listed (such as DPS, GL and graPHIGS) are packaged into separate LPPs and must be ordered as such.

Component	Description
DPS	Specifies Display PostScript, a program product from Adobe Systems, Inc. This program implements the language PostScript and its extensions, as applicable, on a display system instead of on printers. DPS uses 2D graphical functions and complex character fonts, and operates as an extension of the X server.
GAI	Specifies the Graphics Adapter Interface, including the interface to and the set of libraries that form the interface for graphical output to the display

adapters. The GAI libraries implement 2D and 3D graphics functions as well as permit direct access to the hardware. GAI components include:

- GAI Resource Management Support
- GAI 2D Drawing Library
- GAI 3D Model 1 Drawing Library
- GAI 3D Model 2 Drawing Library
- GAI 3D Model 5 Drawing Library

- GL** Specifies Graphics Library. This is the library of graphical functions defined and implemented by Silicon Graphics, Incorporated, under the product name GL. This library is a set of 3D functions, based on a graphical model differing from the PHIGS graphical model. GL is the second graphical model used for 3D. Thus, it is often referred to as the 3D-M2 or the GAI 3D Model 2 Drawing Library.
- graPHIGS** Specifies the graPHIGS application programming interface (API), which is IBM's implementation of PHIGS. This program includes extensions to and deviations from the PHIGS standard. The term PHIGS is a reference to the graPHIGS API, unless explicitly stated to the contrary.
- The graPHIGS API is the first 3D graphical model. It is referred to as 3D-M1 or the GAI 3D Model 1 Drawing Library.
- OpenGL** Specifies OpenGL, a network transparent API for developing applications using 3D graphics. OpenGL is derived from the proprietary SGI GL API.
- PEX** Specifies the 3D Extension to X. The core X protocol provides for basic 2D graphics functionality. The PEX extension adds 3D graphics at about the same functional level as PHIGS and PHIGS PLUS, including features such as lighting, shading, server-side stored structures, and advanced primitives.
- PHIGS** Specifies Programmer's Hierarchical Interactive Graphics System, an accepted international standard (ISO 9592) for the definition, display, and modification of 2D or 3D graphical data, the additional manipulation of geometrically related objects, and the definition and modification of the relationships between the data and the objects. The relationships and the data are stored in a hierarchical data store.
- LFT** Specifies low-function terminal. The LFT is a STREAMS-based simple character oriented tty-like terminal emulator for full screen processing. LFT is a low-cost, low-function enablement feature intended to be used only during system startup, installation and standalone diagnostics. It is not expected to be used in steady-state processing. It supports all the display adapters and keyboards available in AIX Version 4.1, but does not support the mouse, tablet, or any other input devices.

X Window System

Specifies the core of the display subsystem. The X Window System supplies many functions to other components of the subsystem, including resource management, window allocation, a transparent method for communications between platforms, and 2D drawing operations. It has a well-defined mechanism for extending its services. It operates using two X components, the Xlib library and the X server that operate as follows:

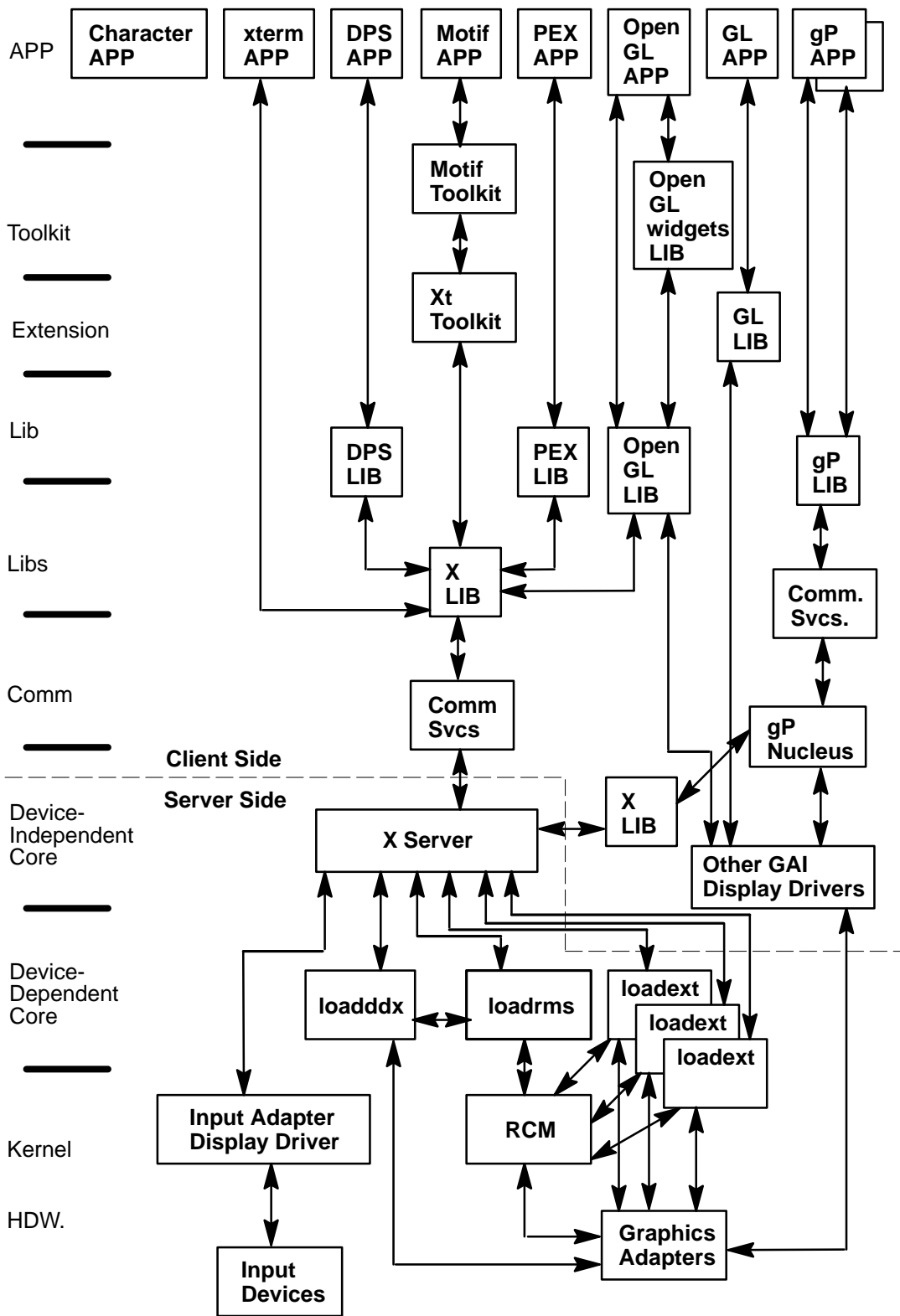
- | | |
|------|--|
| Xlib | Resides in a client and provides the API to users of X Window System. Xlib passes the requests to the X server |
|------|--|

and events to the application. It interacts with the X server using a protocol through an AIX or shared memory socket.

X Server

The X server resides in the X executable and is a component within the display subsystem core. The X server consists of two parts. A part that is device independent (dix) that interprets requests from the Xlib, schedules client activity, manages the return of events and input to the Xlib Library, and performs other generic actions.

The X server also consists of a device-dependent part. The device-dependent part renders the 2D graphics operations defined by X Window System for the specific display adapter. The **loadddx** GAI Load Modules implement this interface.



Functional Block Diagram of Display Subsystem

Application Programming Interface (API)

All types and combinations of applications may exist within a platform. Many of them communicate with the user and display adapter in the server by means of the X Window System protocol over a communications link. Others communicate by means of a graPHIGS protocol over (perhaps a different) communications link.

A list of the generic application types include:

- Character-based applications
- Terminal emulator applications (such as `xterm`)
- Toolkit applications (Motif and Xt)
- Display PostScript applications
- X Window System applications
- OpenGL applications
- GL 3.2 applications
- graPHIGS applications
- PEX applications

Applications *bind* to their appropriate library using any of the AIX system services available to them. The API component provides a set of libraries and toolkits for use by application developers.

API Names		
Common Name	Library Name	Description
Xlib	libX11.a	Low-level X11 interface. Interface to communication services.
Xt	libXt.a	Toolkit intrinsics provided by X Window System.
Xext	libXext.a	X Window System Extension libraries. Provides API for the shape, cursor, colormap, DPS, Direct Access.
Input Library	libXi.a	Provides API for the standard X11 Input Extension.
Motif toolkit	libXm.a	Provides Motif widget and API support.
Motif Resource Manager	libMrm.a	Motif Resource Manager.
OpenGL	libGL.a	Provides OpenGL support.
OpenGLwidgets	libXGLW.a	Provides support for OpenGL widgets.
OpenGL utilities	libGLu.a	Provides a set of utilities to be used with OpenGL.
GL3.2	libg1.a	Provides support for GL 3.2. Sometimes the math library is also needed with GL 3.2.
PEXlib	libPEX5.a	PEXlib API.

API Names		
Common Name	Library Name	Description
PEX PHIGS	libphigs.a	PHIGS API to PEX. This is sometimes called PEX_SI PHIGS.
graPHIGS lib	libgP.a	graPHIGS shell (API bindings and interface to the Nucleus)

The API and library architecture is shown in the Client Side portion of the Functional Block Diagram of Display Subsystem figure, on page 11-8.

All communications begin with AIX sockets. In local or standalone platforms, the sockets are local domain or in shared memory. Between distributed platforms, the sockets are TCP/IP sockets. The type of communication or transport medium is transparent to the application.

X Server

The X server plays an extremely important role within the display subsystem architecture. The X server is the central arbiter and provider of graphical resources. It can be considered the resource server of the display subsystem. The family of X server extensions rely upon the X server. All input flows through the X server. Other components, such as the graPHIGS nucleus or one of the 3D libraries, request basic graphics resources from the X server. The X server's role in the display subsystem architecture is shown in the Functional Block Diagram of Display Subsystem figure, on page 11-8.

The subcomponents of the server are:

- Device Independent Core
- GAI Load Modules—loadddx, loadrms
- X Extensions

The X server communicates with the applications to present data to the screen. It has an input and an output path. *Input* is defined as the direction of data flowing from an input device, through the server, and out through the communications services to the application. *Output* is defined as the direction of data flowing from an application, through the communications services, through the server and device drivers, and onto the display adapter.

The device-independent portion of the X server is made up of modified code from the X Consortium. Specifically the dix (device-independent x), mi (machine independent), mfb (monochrome frame buffer), cfb (color frame buffer) and os (operating system) directories are included. The structure and subroutines have been modified to provide for windows that are being accessed directly by one of the 3D GAI models. In other cases, modifications have been made to optimize for performance. However, in all cases, the programmer interfaces have remained the same to facilitate porting from other platforms.

The device-dependent portion of the X server is provided by the 2D GAI load modules (loadddx and loadrms). The X Window System has a regular architecture to support the definition of extensions. New and third-party extensions will be permitted and assisted.

The following extensions are supported in AIXwindows:

X Window System Extensions		
Common Name	Official Name	Source
DPS	Adobe-DPS-Extension DPSExtension	Adobe
PEX	X3D-PEX	MIT Standard
OpenGL	OpenGL	OpenGL Architecture Review Board
cursor	aixCursorExtension	IBM
colormap blink	xColormapExtension	IBM
X Input	XInputExtension	MIT Standard
Shape	SHAPE Non Rectangular Window (Shape) Extension	MIT Standard
Screen Saver	SCREEN-SAVER	MIT Draft Standard

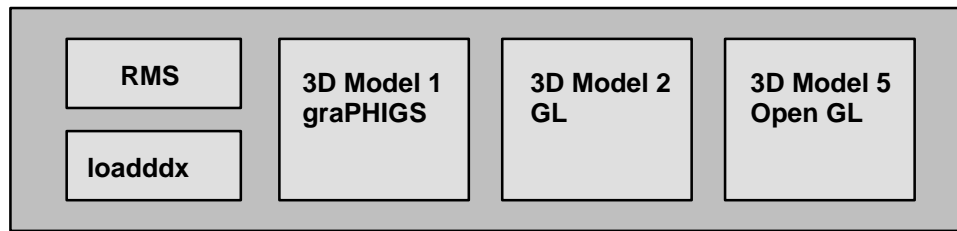
The X Window System extensions have subcomponents in both the X server and API display system components.

GAI Load Modules

One of the fundamental components of the display subsystem is the set of graphics services provided X servers, the graPHIGS nucleus or one of the 3D libraries. These services are contained within the GAI load modules. The GAI load module is bound to the X server, graPHIGS nucleus or 3D library upon its initialization. It provides a library of graphics function to the loading entity.

Except for the RMS load modules, the GAI load modules each provide the device-specific code for one of the GAI models, as illustrated in the Block Diagram of the GAI Load Modules figure. RMS provides a set of support routines that manage hardware resources common to all models, such as windows and clipping regions, but not all models are supported on all adapters. Consequently, there is no need for full function RMS on this class of adapter.

GAI Models



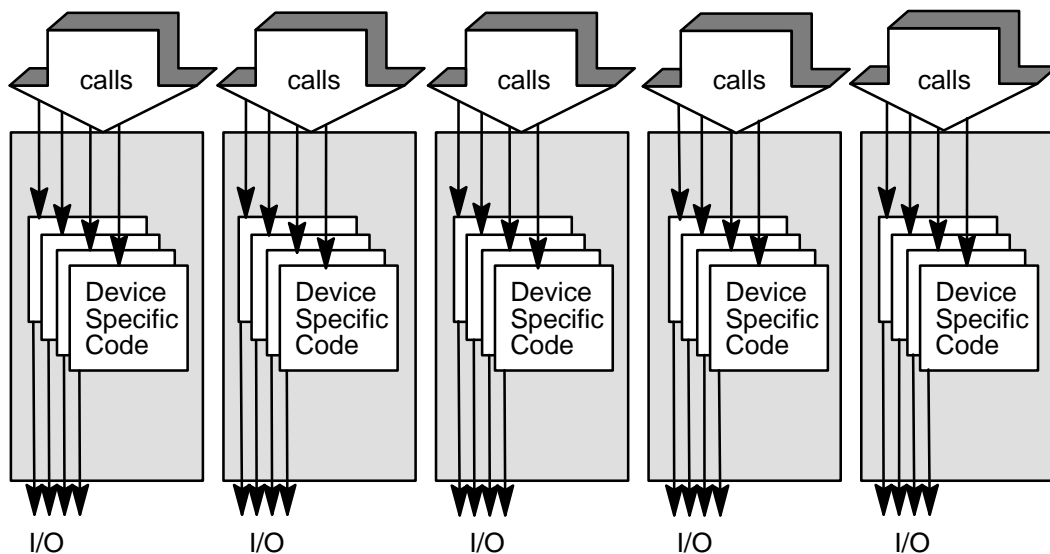
LOADDDX

RMS

3DMOD1

3DMOD2

3DMOD5



Block Diagram of the GAI Load Modules

The GAI load module directly controls the display adapters. When called, it passes direct I/O operations to the display adapters.

The GAI load module also provides functions that assist the X server and other servers in implementing a protocol to directly access the hardware and in managing the resources of the display adapters. It also processes inputs from the display adapters, particularly interrupts and pick events.

The GAI load modules are as follows:

loadddx The loadddx load modules implement the porting layer of the X server. Its interfaces can be found in *Strategies for Porting the X 11 Sample Server*. There is one loadddx for each graphics adapter supported. “Graphics Adapter Interface (GAI) 2D Adapter Load Modules”, on page 11-35, provides techniques for the development of a loadddx.

loadrms The loadrms load modules implement the resource management support for the GAI models. This includes device specific management of window geometries, windows, monitors, system queues, and clipping regions. This resource management is most important for the 3D models accessing hardware directly. “Minimum Resource Management Subsystem (RMS) for 2D Adapters”, on page 11-45, discusses the minimum RMS support needed to implement a 2D graphics adapter.

Note: The following load modules are necessary to provide certain types of 3D support to adapters. They are **not** required for 2D adapters, but are listed here only for informational purposes.

- load3dm1** The GAI 3D Model 1 load modules provides the set of functions required to support a PHIGS-like 3D graphics model within the display subsystem. The support functions have been tailored to complement the architecture of the graPHIGS Nucleus. The functions supported fall into the broad categories of context support and rendering support. The 3dm1 load module provides the support necessary for multiple graPHIGS nuclei to access a display adapter.
- load3dm2** The API within the display subsystem that supports immediate mode graphics is GL. To support this API, GAI 3D Model 2 load module has been defined. It provides direct support for high-function display adapters being driven by GL applications.
- load3dm5** The GAI 3D Model 5 load modules provide the rendering functions for supporting the OpenGL API. These load modules configure the rendering libraries depending on the amount of hardware support available for the API to use. This may range from full hardware support to full software support or some combination of the two. The 3D Model 5 was designed to fully support both direct (server bypass) and indirect (protocol to X server) rendering.

Kernel Components of the Display Subsystem

The kernel components are as follows:

- LFT** The LFT implements a simple character-oriented tty-like terminal emulation. LFT is not intended to be used in a steady state RISC System/6000 environment. Customers are expected to use a graphical interface for all processing except system startup, installation and stand-alone diagnostics
- RCM** The rendering context manager (RCM) permits time sharing of the adapter between graphics processes. This sharing of the adapter is based upon system requests. It also permits space sharing of the adapter, in the sense of having the displayable area of the adapter. The RCM provides the environment for multiple graphics processes independently accessing display hardware. Each of these processes requests that the graphics adapter be in a particular state; the state is called the rendering context for that graphics process. It protects each graphics process from other graphics processes.
- For 2D adapters only displaying X Window System clients, the RCM is not necessary.

Display Device Driver

This material covers configuration information that is unique to graphics and display device driver functions. Use it as a guide to the unique aspects of your adapter, the device driver you write, and how the X server and the LFT subsystem use your device.

Note: This section is preliminary and a prerequisite for writing the ddx portion of your adapter. It is recommended that this section be completed first before moving to the other sections.

LFT Overview

The LFT (low-function terminal) is a simple character oriented, tty-like terminal emulation. It replaces and simplifies the HFT (high-function terminal) of previous releases of AIX, no longer supporting virtual terminals or hotkeying. The LFT supports all graphics adapters and keyboards supported by AIX, but it does *not* support any other types of input devices such as mouse or tablet; these have their own drivers.

The LFT is configured during phase two of the boot process. The LFT will be configured only if there are graphics adapters already available. You must first configure your own adapter and then allow your device to be used by the LFT.

Configuration and ODM Object Classes

Most of the configuration information in the ODM related to devices tells the system management routines and other devices about your device. For graphics adapters some of the information is for graphics adapters only while other information is general about any device.

You must add the following objects for your device and information for other object classes so they can find your device. For example, the LFT and X server need to know if your device is available.

PdDv

The **PdDv** object class is the Predefined Device object class. It contains information about the device and the driver it uses, the configuration methods, class, subclass, and type information of the device.

PdAt

The **PdAt** object class contains information about various attributes of your driver that are unique to graphics. There are also **PdAt** attributes that are used by system management routines and the LFT subsystem.

Some of the attributes used for a graphics adapter are:

display_id	The display_id attribute takes a value of the form 0x04xx0000, where:
04	Fixed.
xx	An adapter specific ID. 0x00 – 0x7F are for IBM use only. 0x80 – 0xFF are for vendor adapters. Use a value not in the database. Remember that all vendors use this range.
dsp_name	This is the value of type field of the PdDv object class. It provides a unique name for your adapter. Check the database to make sure your name is unique.
dsp_desc	Description of the adapter.

scrn_height	Screen height.
scrn_width	Screen width.
color1–16	The 16 default colors to be used by the VDD.
belongs_to	Indicates that the LFT subsystem can use this device. The string graphics is the default value.

Refer to *AIX Version 4.1 Technical Reference, Volume 5: Kernel and Subsystems* for more information on the LFT interface.

CuDep

The **CuDep** object class tells the LFT subsystem that your device is available for use. The LFT subsystem reads this class to find out which devices it can open and use as a graphics output device. You create an object in this object class based on the *belongs_to* attribute in the **PdAt** object class for that adapter.

GAI

The **GAI** object class is in the **/usr/lib/objrepos** directory. This object class requires two entries to tell the location of two modules that need to be loaded during X server startup. (Although these entries are not needed directly by the display driver to run the LFT, they are discussed here because all entries to the ODM should be handled by the define methods of the display driver.) The required entries are:

Adapter_Id	This field is derived from the <i>display_id</i> attribute in the PdAt object class. This is the decimal equivalent of the <i>display_id</i> attribute value and is used as a link to the related objects in other object classes.
Module_Key	This is the name of the module to be loaded when the X server starts.
Module_Path	This is the path of the location of the module that is to be loaded. The X server prepends the string /usr/lpp/gai/ to this to form the full path name of the load module.

Configuration Summary

IBM configuration methods for a VDD (that is, the graphics adapter driver) do not create a **/dev** entry for their device. The device is exclusively used by the LFT subsystem. This simplifies supporting device multiplexing.

Instead of a user level **/dev** entry, the LFT subsystem calls kernel services such as **fp_opendev** and **fp_ioctl**. The LFT subsystem gets to the hardware specific functions of the VDD through the *d_dsdptr* pointer of the devsw table that it obtains through the **devswqry** kernel service.

VDD configuration methods should make sure that they save the address of the **phys_displays** structure in the *d_dsdptr* pointer of the devsw table. Initialize the **phys_displays** structure with pointers to various hardware specific functions in the VDD. If you decide to create a **/dev** entry for the adapter device, you need to incorporate privilege checks and protection mechanisms in your **open** and **close** routines.

The system finds out about your adapter through the **PdDv** object class and adds it to the system through your configuration methods. You allow your device to be used by the LFT through the *belongs_to* **PdAt** attribute.

Your configuration method puts an object into the **CuDep** object class that tells the LFT subsystem to use your graphics adapter. The system management commands **chdisp** and **lsdisp** use the name and description fields to show information about your adapter. Finally,

the X server uses the *display_id* attribute to find out which ddx to load for your adapter. The information for which one to load is contained in the **GAI** object class.

Display Device Driver Subroutines

The standard display device driver subroutines are as follows:

- **vddconfig** routine
- **vddopen** routine
- **vddclose** routine
- **vddioctl** routine
- **interrupt handler**

Note: All subroutines with names beginning **vdd** (virtual device driver) deal with the low level display driver. Do not confuse these with routines with names beginning **vtt**, which are LFT Interface routines.

Configure the Device (vddconfig)

Purpose

Configures the display device driver.

Syntax

```
int vddconfig (devno, cmd, uiop)
dev_t devno;
int cmd;
struct uio *uiop;
```

Description

The **vddconfig** routine initializes the device driver into the device switch table. It also can terminate the device driver by removing itself from the device switch table.

The **vddconfig** routine is called by the display configure, unconfigure, or reconfigure method. This routine can also provide additional device specific functions relating to configuration such as returning device Vital Product Data. This routine is invoked through the **sysconfig** subroutine by the display configure method.

Parameters passed with this routine are:

devno Device major and minor number.

cmd What function this routine performs. The commands are:

INIT Specifies that the **vddconfig** routine is to perform an initialization function. This involves checking the minor number in *devno* for validity, and installing the device driver's entry points in the device switch table. This is accomplished by using the **devswadd** kernel service along with a **devsw** structure to add the device driver's entry points to the device switch table for the major device number supplied in the *devno* parameter.

This routine also copies the device dependant information from the Device Dependant Structure (DDS) provided by the caller into the device specific data area to be used when the device driver routines are invoked. The address and length of the DDS is described in the **uio** structure pointed to by the *uiop* parameter. The **uiomove** kernel service can copy the DDS into the data area of the device driver.

TERM Terminate the device driver and return any system resources. An unconfigure or reconfigure method, through the configuration entry point of the device driver, uses this command to remove resources and system access to this driver. This routine determines if any opens are outstanding on the specified *devno*. If not, it marks the device as terminated and does not allow any subsequent opens to the device.

All dynamically allocated data areas associated with the specified device number should be freed. If this termination removes the last minor number supported by the device driver from use, the **devswdel** kernel service should be called to remove the device driver's entry points from the device switch table for this *devno*

Devices that can act as console should return an error without terminating the device driver.

QVPD Query VPD is an optional function used by the device's configure method to request the return of device specific vital product data. This information is usually used for diagnostic purposes. For this function, the **uio** structure pointed to by *uio* must be set up by the caller to define an area in the caller's storage in which this routine is to write the vital product data. Use the **uio** kernel service to provide the data copy operation.

uio A pointer to a **uio** structure specifying the location and length of the caller's data area in which to transfer information to or from.

This pointer is pointing to a caller provided **uio** structure that describes the location and length of the device-dependent data structure (for INIT) in which to read the information or to the vital product data area (for QVPD) in which to write the requested information. The **uio** kernel service can facilitate the copying of information out of or into the area described by the **uio** structure. The format of the **uio** structure is defined in the **sys/device.h** header file.

Return Values

The **vddconfig** routine should set the return code to zero if no errors were detected for the operation specified. If an error is returned, the return code is one of the values defined in the **errno.h** header file.

Open a Device (**vddopen**)

Purpose

Initializes the display device driver into the system.

Syntax

```
int vddopen (devno, devflag, chan, ext)
dev_t devno;
int devflag;
int chan;
int ext;
```

Description

The **vddopen** routine prepares a device for operation. The kernel calls **vddopen** when a program uses an **open** or **devopen** subroutine.

The display device driver can only be opened by one process at a time (LFT). The **open** subroutine can enforce this by maintaining a static flag variable, which is set to 1 if the

device is open and zero if not. Each time it is called, **vddopen** checks the value of the flag and, if it is not zero, returns with a return code of EIO to indicate that the device is already open. Otherwise, **vddopen** sets the flag and returns normally. The **vddclose** routine later clears the flag when the device is closed.

The **vddopen** routine should initialize the device. It should allocate the required system resources to the device (such as DMA channels, interrupt levels and priorities) and register its **vddintr** device interrupt handler for the interrupt level required to support the target device (if required).

Parameters passed with this routine are:

devno	Specifies both the major and minor device numbers.
devflag	One of the following values:
DKERNEL	The device was called by a kernel routine using DEVOPEN.
DREAD	The device is being opened for reading only.
DWRITE	The device is being opened for writing.
DAPPEND	The device is being opened for appending.
DNDelay	The device is being opened in nonblock mode.
chan	Channel number (ignored by this device driver).
ext	The extended system call parameter (ignored by this device driver).

Return Values

The **vddopen** routine indicates an error condition to the application program by returning a nonzero return code. The return code should be one of the values defined in the **errno.h** header file.

Close a Device (vddclose)

Purpose

Resets the display device driver.

Syntax

```
int vdd_close (devno, chan, ext)
dev_t devno;
int chan;
int ext;
```

Description

The kernel calls the **vddclose** routine when a program uses a **close** or **devclose** subroutine.

The **vddclose** routine resets the display to prevent generating any more interrupts or DMA requests until it is opened again. It should free DMA channels and interrupt levels allocated for this device. The intent is to free system resources used by this device until they are needed again. The flag that was set by the **vddopen** routine should be reset.

Parameters passed with this routine are:

devno	Specifies both the major and minor device numbers.
chan	Channel number (ignored by this device driver).
ext	The extended system call parameter (ignored by this device driver).

Return Values

The **vddclose** routine indicates an error condition by returning a nonzero return code. The return code should be one of the values defined in the **errno.h** header file.

Device Control (vddioctl)

Purpose

The **vddioctl** routine provides control commands and parameters to the device.

Syntax

```
long vdd_ioctl (devno, cmd, arg, devflag, chan, ext)
dev_t devno;
long cmd;
long arg;
ulong devflag;
long chan;
long ext;
```

Description

The **vddioctl** routine performs special I/O operations. These operations are normally device specific. Within the LFT subsystem it is not used, providing only access to adapter diagnostic functions. However, it is a valid porting strategy for upper level graphics libraries, such as the adapter ddx, to use **vddioctl**. Possible operations include returning adapter bus addresses for I/O, and DMA control.

Parameters passed with this routine are:

devno	Both the major and minor device numbers.
cmd	Command parameter indicating what function this routine should perform.
arg	Parameter from the ioctl subroutine call that specifies an additional argument for the <i>cmd</i> operation.
devflag	Device open or other control flags.
chan	Channel number (typically not used by this device driver).
ext	The extended system call parameter (typically not used by this device driver).

Return Values

The **vddioctl** routine indicates an error condition to the application program by returning a nonzero return code. The return code should be one of the values defined in the **errno.h** header file. Data may be returned to the user application by use of the **copyout** kernel service.

Programming Notes

There are *no* required ioctl commands. This entry point should always return success if there are no vendor-specific commands.

Interrupt Handling

Graphics adapters which generate interrupts need an interrupt handler. The interrupt handler is made known to the system during the device driver initialization. The specific interrupts to be handled are device dependent.

LFT Interface Routines

All the following routines take a pointer to the **vtmstruc** structure associated with the virtual terminal. The structure contains terminal specific data and pointers. The VDD uses this structure to obtain the pointer to its local terminal dependent data area, as well as to determine the position to move the cursor to when applicable. The following is a list of all the routines necessary to define the LFT interface:

- Activate (**vtact**)
- Copy full line (**vttcfl**)
- Clear rectangle (**vttclr**)
- Copy line (**vttcpl**)
- Deactivate (**vttdact**)
- Define cursor (**vttdefc**)
- Initialize (**vttnit**)
- Move cursor (**vttmovc**)
- Scroll display (**vttscr**)
- Terminator (**vttterm**)
- Draw text (**vtttext**)

Note: Structures are described in “Display Driver Structure Descriptions”, on page 11-31.

Activate (**vtact**)

Synopsis

The **vtact** routine switches the Virtual Display Driver into the active state and gives the virtual display driver exclusive access to the display hardware.

Description

This routine copies the presentation space in the frame buffer and establishes the correct position of the hardware cursor. The presentation space is first cleared and then the cursor is placed in the upper left corner.

Also, the color palette is loaded.

Syntax

Use this routine with the following call:

```
rc = (*vp->display->vtact)(vp);
struct vtmstruc *vp;
```

The parameter passed with this routine is:

vp Pointer to the **vtmstruc** structure associated with this terminal.

Return Values

The return value 0 indicates successful completion.

Programming Notes

Only one instance of the device driver can be in the activated state at one time.

Calling this routine when the device driver is already activated causes unpredictable effects on subsequent operations.

You *must* call the initialize (**vttnit**) routine (see page 11-25) prior to the first call to this routine. Failure to do so causes unpredictable effects on the operation of this and subsequent routines.

Copy Full Lines (vttcfl)

Synopsis

The **vttcfl** routine copies the entire or partial content of the presentation space by one or more full lines positions, either up or down. The source and destination points are within the presentation space and the resulting information which lies beyond the absolute lower right hand or upper left hand corner of the presentation space is lost.

Description

The command copies a sequence of one or more consecutive full lines of character and attribute pairs in either direction within a presentation space and truncates the modified content at the presentation space boundaries.

Use this command with the Clear Rectangle (**vttclr**) routine (see page 11-22) to insert new lines into the presentation space.

In addition, the cursor can optionally be removed or left shown on the screen.

Syntax

Use this routine with the following call:

```
rc = (*vp->display->vttcfl) (vp, source_row, dest_row,
length, update_cursor);
struct vtmstruc *vp;
long source_row;
long dest_row;
long length;
ulong update_cursor;
```

Parameters passed with this routine are:

- vp** Pointer to the **vtmstruc** associated with this terminal.
- source_row** Row number at which to begin the line copy operation.
- dest_row** Row number to which the line copy operation will copy the first row of the source.
- length** Number of lines which are to be copied.
- update_cursor** This is a Boolean value. If a device does not have a hardware cursor and this parameter is set to zero, the cursor is not visible when this command returns to the caller. Otherwise, the cursor is visible.

Note: If a display has a hardware sprite cursor, this routine does not affect the visibility of the cursor. Rather, UPDATE_CURSOR specifies whether the cursor is moved to the position specified in *vp->mparms.cursor*.

Return Values

The return value 0 indicates successful completion.

Programming Notes

The actual amount of data copied varies, depending on the physical display adapter.

The cursor position used for positioning is contained in the **vtmstruc** structure pointed to by *vp*.

Clear Rectangle (**vttclr**)

Synopsis

The **vttclr** routine clears any specified rectangular area of the screen.

Description

This routine stores a space in each character position of the rectangular area with any combination of the following attributes:

- Foreground color (one of 16 different colors)
- Background color (one of 16 different colors)
- Font (one of 8 different fonts)
- Underscore
- Reverse image
- Blink
- Bright
- Non-display

If the real display does not have a hardware cursor, then the cursor may optionally be displayed or not displayed when this routine returns to the caller. If the real display has a hardware cursor, the cursor is always displayed when the routine returns to the caller.

Syntax

Use this routine with the following call:

```
rc = (*vp->display->vttclr)(vp, screen_pos, attr, update_cursor);
struct vtmstruc *vp;
struct vtt_box_rc_parms *screen_pos;
ulong attr;
ulong update_cursor;
```

Parameters passed with this routine are:

- vp** Pointer to the **vtmstruc** structure associated with this terminal.
- screen_pos** Pointer to a **vtt_box_rc_parms** structure with the coordinates of the upper-left and lower-right corners of the rectangular area that is cleared. The **vtt_box_rc_parms** structure is defined in "Display Driver Structure Descriptions" on page 11-31.
- attr** Attributes of each character in the specified rectangular area.
- update_cursor** This is a Boolean value. If a device does not have a hardware cursor and this parameter is set to zero, the cursor is not visible when this command returns to the caller. Otherwise, the cursor is visible.

Note: If a display has a hardware sprite cursor, this command does not affect the visibility of the cursor. Rather, the *update_cursor* parameter specifies whether the cursor is moved to the position specified in the **vtmstruc** structure.

Return Values

The return value 0 indicates successful completion.

Programming Notes

If a rectangular area is not valid (such as coordinates outside the presentation space, upper-left and lower-right coordinates switch), the result of running this routine is unpredictable.

The cursor position used for positioning is contained in the **vtmstruc** structure pointed to by *vp*.

Copy Line Segment (**vttcpl**)

Synopsis

The **vttcpl** routine copies a specified segment in one or more consecutive lines either left or right.

Description

This routine copies a sequence of one or more consecutive character/attribute pairs in either direction within a line and truncates the modified content at the line boundaries.

The contents of the preceding lines are not affected. Succeeding consecutive lines, however, can be similarly copied if the number of lines requested to operate on is greater than one.

Use this routine with the Clear Rectangle (**vttclr**) routine (see page 11-22) to obtain the repetitive inline “move” function.

If the starting point is the first position in the first line and the number of lines requested is equal to the number of rows in the presentation space, then the entire screen is scrolled right. Otherwise, there is a partial screen scroll.

Syntax

Use this routine with the following call:

```
rc = (*vp->display->vttcpl) (vp, rc, update_cursor);
struct vtmstruc *vp;
struct vtt_rc_parms *rc;
ulong update_cursor;
```

Parameters passed with this routine are:

- vp** Pointer to the **vtmstruc** structure associated with this terminal.
- rc** Pointer to a **vtt_rc_parms** structure containing the “copy to” and “copy from” information. This structure is defined in “Display Driver Structure Descriptions” on page 11-31.
- string_length** Number of character/attribute pairs to be copied in each line.
- string_index** Index of the first character to display.
- start_row** Beginning line on which to operate.
- start_column** Starting position, within each line, of the source to be copied from.
- dest_row** Last line on which to operate.
- dest_column** Starting position, within each line, of the destination to be copied to. If the destination column is bigger than the starting column, the specified line segment is copied to the right. Otherwise, it is copied to the left.
- update_cursor** This is a Boolean value. If a device does not have a hardware cursor and this parameter is set to zero, the cursor is not visible when this routine returns to the caller. Otherwise, the cursor is visible.

Note: If a display has a hardware sprite cursor, this routine does not affect the visibility of the cursor. Rather, *update_cursor* specifies whether the cursor is moved to the position specified in the **vtt_cursor** structure.

Return Values

The return value 0 indicates successful completion.

Programming Notes

The actual amount of data copied varies, dependent on the physical display adapter.

The cursor position used for positioning is contained in the **vtmstruc** structure pointed to by *vp*.

Deactivate (vttddact)

Synopsis

The vttddact routine switches the virtual display driver into the inactive state.

Description

Recall that the device driver model maintains a presentation space model. When the display hardware has a character buffer, that buffer will, in general, be used to contain and maintain the presentation space data while a given instance of the device driver is active (the so called integral presentation space case). If no hardware character buffer is used to maintain the presentation space the device driver must maintain a presentation space while either Active or Inactive.

Syntax

Use this routine with the following call:

```
rc = (*vp->display->vttddact) (vp);
struct vtmstruc *vp;
```

The parameter passed with this routine is:

vp Pointer to the **vtmstruc** structure associated with this terminal.

Return Values

The return value 0 indicates successful completion.

Programming Notes

Calling this routine when the device driver is already deactivated causes unpredictable effects on subsequent operations.

Define Cursor (vttdefc)

Synopsis

The **vttdefc** routine changes the shape of the cursor to one of six predefined shapes.

Description

Select the shape of the cursor from a set of shapes which is dependent on the display device. Selectors 6 through 255 are reserved. The default (selector 2, the double underscore) is used if you specify a reserved selector.

This routine can also reposition and turn off the cursor.

Syntax

Use this routine with the following call:

```
rc = (*vp->display->vttdefc) (vp, selector, update_cursor);
struct vtmstruc *vp;
uchar selector;
ulong update_cursor;
```

Parameters passed with this routine are:

vp Pointer to the **vtmstruc** structure associated with this terminal.

selector Cursor shape to display. The available shapes are:

0	No cursor
1	Single Underscorer
2	Double Underscore
3	Half Blob
4	Double Line
5	Full Blob
Other values	Default to double underscore.

update_cursor

This is a Boolean value. Whether the cursor should be made visible after its shape is changed.

Note: If a display has a hardware sprite cursor, this command does not affect the visibility of the cursor. Rather, `Show_Cursor` specifies whether the cursor is moved to the position specified in `Cursor_Pos`.

Return Values

The return value 0 indicates successful completion.

Programming Notes

The cursor is invisible on all displays if the selector value is zero (a null shape is chosen). The `update_cursor` parameter does not affect the visibility of the cursor on displays with hardware cursors.

If you specify an invalid cursor position, the results are unpredictable.

On display adapters that do not support a hardware cursor, the cursor shape is generated with an exclusive-or operation.

The cursor position used for positioning is contained in the **vtmstruc** structure pointed to by `vp`.

Initialize (vttinit)

Synopsis

The **vttinit** routine initializes the internal state of the virtual display driver. It also allows the caller to select up to eight different fonts for subsequent use by the virtual terminal (see the Draw Text (**vtttext**) subroutine on page 11-29). All characters in all the selected fonts must be the same size.

Description

This routine initializes the internal state of the virtual display driver. A list of eight font IDs may optionally be passed to this procedure.

All characters in all fonts in the list *must* be the same size. Any font whose character box size differs from that of the first font in the list is replaced by the first font in the list.

To allow custom font selection, `vttinit` is passed a `font_id`. Using the **chfont** command, the user can create a `font_id` which is stored in the ODM to be used at LFT configuration time. If `font_id` is not equal to `-1`, then the index is used to load the font table. If the `font_id` is equal to `-1`, then there has been no custom font selection, and the LFT uses the first available font in the font table.

If the font list is validated or a default font is found:

- The presentation space (PS) is initialized with space characters and the canonical attribute of each character in the PS is set to zero.
- The cursor is moved to the upper left corner of the PS and displayed if the virtual display driver is active.
- The default cursor shape (double underscore) is selected.

Syntax

Use this routine with the following call:

```
rc = (*vp->display->vttinit)(vp, font_ids, ps_size);
struct vtmstruc *vp;
struct fontpal *font_ids;
struct ps_s *ps_size;
```

Parameters passed with this routine are:

vp	Pointer to the virtual-terminal-specific data area.
font_ids	Pointer to a fontpal structure that contains the font IDs to use for the various fonts selected by the attribute encoding. The font order is important as it matches the font selector enumerated scalar in the canonical representation of attributes (see “vtt_cp_parms” on page 11-31). This parameter is optional. If it is not specified, one default font will be selected.
ps_size	Pointer to a ps_s structure that contains the width and height (in characters) of the presentation space.

Return Values

This routine sets the height and width of the presentation space in the `ps_s` parameter. If it is unable to allocate the local data, `vttinit` returns `ENOMEM`. The return value `0` indicates successful completion.

Programming Notes

All fonts specified in this routine should be the *same size* because of the default actions.

Call this routine *before* the first call to the activate (**vttact**) routine (see page 11-20). If you do not, there are unpredictable effects on subsequent operations.

The PS size is calculated using the following formulas (all division is integer division):

$$\text{Presentation Space Height} = \frac{\text{Height of the real screen (in picture elements)}}{\text{Number of rows in a character box}}$$

$$\text{Presentation Space Width} = \frac{\text{Width of the real screen (in picture elements)}}{\text{Number of columns in a character box}}$$

Move Cursor (vttmovc)

Synopsis

The **vttmovc** routine moves the cursor to the indicated position.

Description

The current cursor shape is repositioned to the specified character row and column. The cursor is always visible after this routine if the cursor has been defined as a visible shape, and this virtual terminal is active.

Syntax

Use this routine with the following call:

```
rc = (*vp->display->vttmovc) (vp);  
struct vtmstruc *vp;
```

The parameter passed with this routine is:

vp Pointer to the **vtmstruc** structure associated with this terminal.

Return Values

The return value 0 indicates successful completion.

Programming Notes

If an invalid position is specified, the results are unpredictable.

Scroll (vttscr)

Synopsis

The **vttscr** routine scrolls the entire contents of the display screen up or down.

Description

The current display screen contents are moved the indicated number of lines up or down. When scrolling up, the lines moved off the screen at the top are discarded and the indicated number of lines at the bottom are cleared to blanks with the attributes provided. When scrolling down, the lines moved off the screen at the bottom are discarded and the indicated number of lines at the top are cleared to blanks with the attributes provided. The cursor may also be repositioned and made invisible.

Syntax

Use this routine with the following call:

```
rc = (*vp->display->vttscr)(vp, lines, attributes,
update_cursor);
struct vtmstruc *vp;
long lines;
ulong attributes;
ulong update_cursor;
```

Parameters passed with this routine are:

- vp** Pointer to the **vtmstruc** structure associated with this terminal.
- lines** Number of lines to scroll. A positive value moves the screen contents toward the top. A negative value moves the screen contents toward the bottom of the screen.
- attributes** Attributes of the blanks to be inserted in the new lines that appear at either the top or bottom of the screen.
- update_cursor** This is a Boolean value. Whether the cursor should be made visible after the scroll.

Note: If a display has a hardware sprite cursor, this command does not affect the visibility of the cursor. Rather, `Show_Cursor` specifies whether the cursor is moved to the position specified in `Cursor_Pos`.

Return Values

There are no return values.

Programming Notes

If a position is not valid, the results are unpredictable.

The cursor position used for positioning is contained in the **vtmstruc** structure pointed to by *vp*.

Terminate (vttterm)

Synopsis

The **vttterm** routine prevents interrupts to a virtual terminal and must be issued when the virtual terminal closes.

Description

A virtual terminal in the process of closing no longer requires interrupts from its I/O devices. In fact, after the virtual terminal process terminates, virtual interrupts directed to that process cause a system failure. To prevent such problems, the virtual terminal must call the terminate routine before terminating itself.

Syntax

Use this routine with the following call:

```
rc = (*vp->display->vttterm)(vp);
struct vtmstruc *vp;
```

The parameter passed with this routine is:

vp Pointer to the **vtmstruc** structure associated with this terminal.

Return Values

The return value 0 indicates successful completion.

Programming Notes

Not issuing this routine when closing causes a system failure.

For display adapters that generate interrupts serviced by the virtual terminal, the routine ensures that the adapter is not reinitialized when the virtual terminal process terminates.

For device drivers which allocate storage for their presentation space, this routine will free that storage and all the local data areas.

Draw Text (**vtttext**)

Synopsis

The **vtttext** routine draws a string of qualified ASCII characters into the refresh buffer and presentation space buffer of the display device.

Description

Each supplied ASCII character is drawn into the refresh buffer and/or presentation space buffer beginning at the specified starting row and column. At the end of this character drawing operation, the cursor is redrawn at the specified new location.

Each character drawn is mapped to the appropriate font range by two mechanisms. First, each character will be logically ANDed by the value supplied in the code point mask parameter. It is expected that this parameter will contain either 0xFF or 0x7F which will have the effect of wrapping the code point in either a seven-bit range or an eight-bit range. Next, the code point base parameter is added to each supplied character to find the correct symbol in the font.

Each character will be drawn with a device specific attribute derived from the canonical attribute specification in the *vtt_attr* parameter.

The cursor can be moved or redrawn by this routine if the *update_cursor* parameter is set to true (1). In this case, a cursor of the type specified by the Define Cursor (**vttdefc**) routine (see page 11-24) will be moved or redrawn, as appropriate for the given hardware at the location specified in the **vtmstruc** parameters.

Note: The no show state (zero) of *update_cursor* does not guarantee that the cursor will be invisible after this operation. Whether or not it is invisible is a device-dependent characteristic.

Recall that the device driver model maintains a presentation space model. Usually, where the display hardware has a character buffer, that buffer is used to contain and maintain the presentation space data while a given instance of the device driver is active. This is the *integral presentation space case*

If no hardware character buffer is used to maintain the presentation space, this is the *disjoint refresh buffer case* in which the driver must maintain an off adapter presentation space while active or inactive.

Syntax

Use this routine with the following call:

```
rc = (*vp->display->vtttext)(vp, ascii_string, rc, cp,
update_cursor);
struct vtmstruc *vp;
char *ascii_string;
struct vtt_rc_parms *rc;
struct vtt_cp_parms *cp;
ulong update_cursor;
```

Parameters passed with this routine are:

- vp** Pointer to the **vtmstruc** structure associated with this terminal.
- ascii_string** Adjustable extent array of characters which are to be drawn on the display. The length of this string must be greater than or equal to the *string_length* parameter in the **vtt_rc_parms** structure.
- rc** Pointer to a **vtt_rc_parms** structure containing the ASCII string, row, and column information.
- string_length** How many of the characters in the ASCII string are to be drawn. Note that this drawing operation performs no end-of-line check so that unpredictable results will occur if the combination of *string_length* and *start_column* causes an attempt to write off the end of the line.
- string_index** Index of the first character in *ascii_string* that is to be displayed.
- start_row** Specifies, in the presentation buffer and display refresh buffer, the row number of the first character drawn. This parameter has a unity origin.
- start_column** Specifies, in the presentation buffer and display refresh buffer, the column number of the first character drawn. This parameter has a unity origin.
- dest_row**
dest_column Not used.
- cp** Pointer to a **vtt_cp_parms** structure containing the ASCII string, and row and column information.
- cp_mask** 8-bit value which is logically ANDed with each ASCII character before font translation. Legal values for this parameter are either 0xFF or 0x7F.
- cp_base** Value which allows each ASCII character to be translated to a final font set, larger than 256 code points. The ASCII character is masked by *cp_mask*. The masked ASCII character is logically added to the value of *cp_base* to determine the actual character, in a 10-bit selection, to be displayed.
- attributes** Canonical representation of the desired attributes which should be used when drawing the character string. See the discussion of canonical attribute representation under “vtt_cp_parms” on page 11-31.
- cursor** A **vtt_cursor** structure that defines the x and y position of the cursor.
- update_cursor** This parameter is a Boolean value and controls whether the cursor should be redrawn at the end of this draw routine. The purpose is to allow the LFT subsystem to control how much APA display processing is wasted in undrawing and redrawing cursors.

Note: If a display has a hardware sprite cursor, this routine does not affect the visibility of the cursor. Rather, *update_cursor* specifies whether or not the cursor is moved to the position specified in the **vtt_cp_parms** structure.

Return Values

There are no return values.

Programming Notes

Use this routine *only* to draw in one *single* line of the display at a time. Length specifications that would imply a wrap to next line in the middle of the call will cause unpredictable results in displays with invisible refresh buffer space at the end of the visible line.

The cursor position used for positioning is contained in the **vtmstruc** structure pointed to by *vp*.

Display Driver Structure Descriptions

vtt_rc_parms

The **vtt_rc_parms** structure is supplied as a parameter to several of the virtual display driver routines, notably the Draw Text routine, the two copy routines and the Read Screen Segment routine.

```
/* *****  
/* row column length structure interfaces to VDD Routine      */  
/* *****  
struct vtt_rc_parms {  
    ulong string_length; /* length of character string that */  
                        /* must be displayed                */  
    ulong string_index; /* index of the 1st char to display */  
    long  start_row;    /* starting row for draw/move      */  
                        /* operations, unity based.        */  
    long  start_column; /* starting column for draw/move   */  
                        /* operations, unity based         */  
    long  dest_row;     /* destination row number for move */  
                        /* operations, zero based          */  
    long  dest_column; /* destination column number for move */  
                        /* operations, zero based          */  
};
```

vtt_box_rc_parms

The **vtt_box_rc_parms** structure is supplied as a parameter to the Clear Rectangle routine entry points.

```
/* *****  
/* the vtt_box_rc_parms structure is supplied as a parameter to the */  
/* VDD clear rectangle routine.                                     */  
/* *****  
struct vtt_box_rc_parms {  
    long row_ul; /* row number of upper-left corner */  
                        /* of the upright rectangle        */  
    long column_ul; /* col number of upper-left corner */  
                        /* of the upright rectangle        */  
    long row_lr; /* row number of lower-right corner */  
                        /* of the upright rectangle        */  
    long column_lr; /* col number of lower-right corner */  
                        /* of the upright rectangle        */  
};
```

vtt_cp_parms

The **vtt_cp_parms** structure is supplied as a parameter to the Draw Text routine. The **vtt_cursor** substructure is passed as a parameter to several procedures such as Define Cursor, Move Cursor, and Scroll Up.


```

/*****
/* code point base/mask and cursor positioning parameter      */
/* structure for use in vtttext replace text                  */
/*****
struct vtt_cp_parms
{
    ulong cp_mask;      /* code point mask for implementing */
                        /* 7 or 8 bit ascii                */
    long cp_base;      /* code point base, added to code */
                        /* point if base >= 0             */
    ushort attributes; /* attribute bits                  */
    struct vtt_cursor cursor; /* cursor x and y position      */
};

```

The definitions and default states of the canonical attributes are:

- FG COLOR** Foreground color specification in writing into the frame buffer. The default value is device dependent.
- BG COLOR** Background color specification to be used in writing into the frame buffer. The default value is device dependent.
- FONT_SELECT** Select for the font to be used in writing into the frame buffer. A value of zero selects the first font ID supplied in the VTTINIT font ID array as the active font for drawing. A value of seven selects the eighth font ID array value as the active font for drawing. The default value is zero.
- NO_DISP** Select non-displayed mode for characters. One equals non-displayed, zero equals displayed. Default value is zero.
- BRIGHT** Select bright (intensified) display mode. One equals intensified. Zero equals normal intensity. Default value is zero.
- BLINK** Select blinking representation. One equals blinking characters. Zero equals non-blink. Default value is zero.
- REV_VIDEO** Select reverse video representation. One equals reversed image. Zero equals normal image. Default value is zero. To get the effect of reverse video on all adapters, this bit must be set to one. Swapping the foreground and background colors will not yield the desired result on all adapters because these fields may in fact be ignored by some virtual device drivers.
- UNDERSCORE** Select underscoring of characters. One equals underscore each character. Zero equals no underscore. Default value is zero. A default color table is provided for any display type. The size and organization of color tables is device dependent.

No single device supports all the canonical attributes in any configuration. For example, when a canonical attribute is supplied to the Draw Text routine, the canonical form is interpreted to something usable by the given device. When a subsequent call to the Read Screen Segment routine is made, the internally stored attributes will be mapped back to canonical form. This back transform has some loss of information in that canonical attributes that are unsupported by a given device will be set to a default state instead of the originally supplied state.

font_data

The **font_data** structure is initialized by the LFT and is used by the VDD to get the necessary font information for displaying text.

```

struct font_data
{
    ulong font_id;           /* font_data index (1 based) */
    ulong font_attr;        /* 0=plain, 1=bold, 2=italic */
    ulong font_style;       /* 0=apa,*/
    long  font_width;       /* charbox width  in pels*/
    long  font_height;      /* charbox height in pels*/
    long  *font_ptr;        /* pointer to font file*/
    ulong font_size;        /* size of font structure*/
};

```

phys_displays

The **phys_displays** structure is defined in `/usr/include/sys/display.h`.

```

/*****
/* presentation space structure */
*****/

struct ps_s {
    long ps_w;           /* presentation space width */
    long ps_h;           /* presentation space height */
};

/*****
/* cursor positioning structure used as parameter to VDD routine */
*****/

struct vtt_cursor {
    long x;
    long y;
};

struct vtmstruc {
    struct phys_displays *displays; /* display this VT is using */
    struct vtt_cp_parms  mparms;    /* attribute+cursor position*/
    char                 *vttld;    /* pointer to VTT local data*/
    off_t                vtid;      /*virtual terminal id = 0 */
    uchar                vtm_mode;  /* mode = KSR */
    int                  font_index; /* -1 means use 'best' font */
    int                  number_of_fonts; /* number of fonts found */
    struct font_data     *fonts;    /* font information */
    int                  (*fsp_eng)(); /* font request enqueue func*/
};

```

Device Dependent Structure (DDS)

The Device Dependent Structure is specific to the display adapter. The contents must contain the addresses, interrupt levels, DMA areas, and other configuration information set by the configuration method.

The following is an example Device Dependent Structure:

```
struct display_dds
{
    int adapter_busaddr;
    int int_level;
    int int_class;
    int slot_number;
    int color_tablet[16];
    int dms_channel;
    int dma_address;
    int screen_height_mm;
    int screen_weight_mm;
    int display_id;
};
```

Graphics Adapter Interface (GAI) 2D Adapter Load Modules

The device-dependent portion (ddx) of the X server code (the portion of the server code that actually manipulates the adapter) is located in dynamically loadable modules under the current GAI architecture. These modules are loaded at server invocation time when a user requests that the server run on a specific adapter via command-line flags, or the default adapter(s) in the absence of a specific request. Adapter modules not requested in this instantiation of the server are not loaded.

The interface between the device-independent (dix) portion of the server and the device dependent (ddx) subsystem is documented in *Strategies for Porting the X1 Sample Server* by Angebrannt and others. The porting layer specified in the porting guide is a well-defined interface which allows third party vendors to supply those device-dependent portions of the X server required to support a new device. In the current IBM GAI architecture, the ddx resides in dynamically loadable, device-dependent code modules (loadddx's) and are required to support execution of the X server on display adapters. By utilizing a dynamic loading scheme, the GAI architecture facilitates the independent addition of new driver modules by third-party vendors, as access to object modules for static binding is not necessary.

The loadable module technique is utilized to support the addition of extensions to the X Window System as well. Loading only the required extensions for a particular invocation of the server limits the system resources from being consumed. An additional benefit is that third-party extension writers can add their X server extensions without requiring server integration and rebuilding.

This information highlights the process of initializing and destroying 2D minimal RMS adapter instances within the GAI architecture. Details regarding the implementation of an X device dependent subsystem (ddx) can be found in standard documents, such as *The X Window System Services* by Israel and Fortune.

Loadable DDX Interface

The process by which the X server and one or more adapters is initialized is divided into three steps:

- Selecting the adapter or adapters through command-line flags or default actions
- Loading the appropriate load module and dynamic binding with X server executable
- Initializing the hardware and the X server screen structure

Selection of Adapters

Command-line flags allow the user to specify a number of characteristics of the X server display. Characteristics relating to adapter selection include the number of screens for the display, configuration of the screens for the display, and default visual and default depth of the display. Command-line flags for adapters include the logical name returned from the AIX **lsdisp** command as they are stored in the Object Data Manager database (the ODM database contains the configuration information for the entire system). As an example, the **lsdisp** command returns the following information if the system is configured with a Color Graphics Display Adapter and a POWER_Gt1 Graphics Adapter:

DEV_NAME	SLOT	BUS	ADPT_NAME	DESCRIPTION
sga0	0J	sys	POWER_Gt1	Graphics Adapter
gda0	03	mca	colorgda	Graphics Adapter

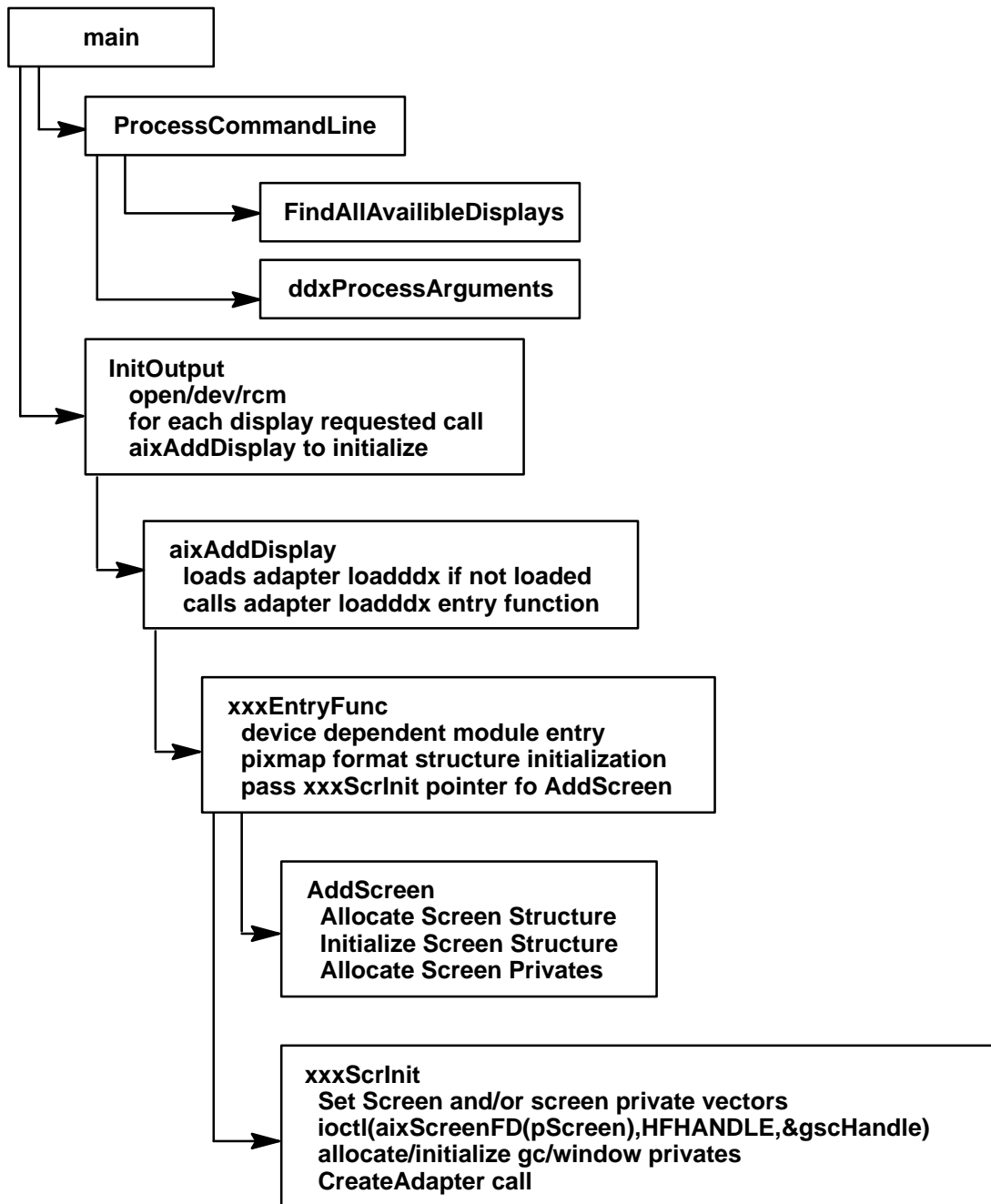
If no display characteristic command-line flags are specified, the default actions are used. These default actions include:

- Use of all configured adapters. If there is only one graphics adapter installed, the X server will be initialized on that graphics adapter. If more than one adapter is configured in the system, the X server will be initialized in multihead mode.
- Use of PseudoColor visual for the default visual (or GrayScale visual for grayscale adapters).
- Use of default depth (usually 8).

If the user chooses to bypass the default server configuration options and initialize on less than all available display adapters, the `-P Row Column DeviceName` command-line option must be specified. If the user chose to initialize the X server on Color Graphics Display Adapter (assuming output from **ldisp** above), the appropriate command line options would include the flag to initialize on Color GDA only, `-P11 gda0`.

As part of the adapter initialization process, the X server main function calls an operating system dependent subroutine, **ProcessCommandLine**. **ProcessCommandLine** calls **FindAllAvailableDisplays** function, which queries the ODM database for all available graphics adapters in the system and sets up a **DisplayRec** structure for every adapter found. The **DisplayRec** structures are actually stored in a global dynamic array, **AvailableDisplays**. **FindAllAvailableDisplays** returns to **ProcessCommandLine** which calls a device-dependent command line parsing function, **ddxProcessArgument**. **ddxProcessArgument** handles some GAI-specific global variable initializations based on the command line flags specified for wrapping pointer in x and y directions, backing-store configuration, extension loading, and adapter initialization. It is in **ddxProcessArgument** that the remaining structure members in the adapter **DisplayRecs** are initialized to specify users interest in X initialization on a given adapter. When **ddxProcessArgument** encounters a `-P` flag, it uses the logical name associated with the flag as a key in the ODM **CuDv** database and retrieves the associated entry. If the entry does not exist, an error message is displayed. When an unexpected option or invalid value is specified for a flag, the X server displays the usage message and ignores the option or terminates initialization, depending on the severity of usage error. Both **ProcessCommandLine** and **ddxProcessArgument** validate command line options and display error messages.

Assuming that there are no errors detected by **ddxProcessArgument**, **InitOutput** is called from the **main** subroutine, passing a pointer to the **screenInfo** structure. The pixmap and bitmap format information is initialized in the **screenInfo** structure and for each requested adapter **aixAddDisplay** is called to load the **loaddx** if it has not been bound in (on server reset, the module would have already been loaded) and evoke the drivers entry function, for example **xxxEntryFunc** which sets device specific **pixmapFormat** information and calls the MIT dix function **AddScreen** which "increases" the size of the **screenRecs** (internal server structure) array and calls the device-specific initialization routine specifying the adapters screen number, a pointer to its **ScreenRec**, and the command line argument vector and argument count where the adapter will be initialized, screen vectors set, and access to the adapter given to X and rms configured via **CreateAdapter** call.



Server Initialization Sequence for the Selection and Initialization of Displays

X Server Initialization Subroutines

Note: The following subroutines are for informational purposes only.

ddxProcessArgument

Purpose

Processes the command line arguments for the X server.

Syntax

```
int ddxProcessArgument (argc, argv, i)
int argc;
char *argv[];
int i;
```

Description

The **ddxProcessArgument** subroutine is an AIXwindows supplied subroutine. There is one for the X server rather than one for each ddx. It is called for each flag one flag at a time. Its parameters include the *argc* and *argv* parameters which were used for this invocation for the X server and the current index into *argv*. **ddxProcessArgument** looks at the current argument and returns zero if the argument is not a device-dependent one, and otherwise returns a count of the number of elements of *argv* that are part of this one argument.

Parameters

<i>argc</i>	Specifies the number of options in the <i>argv</i> list.
<i>argv</i>	Specifies an array of options that were passed to the X Server at its invocation.
<i>i</i>	Specifies the current index into <i>argv</i> list.

Flags

ddxProcessArgument processes the following arguments (not a complete list):

Mouse Processing Flags

Flag Name	Variable Set
-wrapx	aixXWrapScreen Specifies mouse behavior when the mouse's hot spot reaches the left or right borders of any root window.
-wrapy	aixYWrapScreen Specifies of mouse behavior when the mouse's hot spot reaches the top or bottom borders of any root window.
-wrap	aixXWrapScreen, aixYWrapScreen Sets both the -wrapx and -wrapy flags.

Other Flags

Flag Name	Variable Set
-bs	aixAllowBackingStore Enables backing store support on all screens.
-nobs	aixAllowBackingStore Disables backing store support on all screens.
-x <extention name>	how_many_extensions, extension_list Tracks number of and list of requested extensions.
-T	aixDontZap Disables the Ctrl-Alt-Backspace key sequence used to kill the server.

- wp** **aixWhitePixelText**
Allows you to specify a WhitePixel color.
- bp** **aixBlackPixelText**
Allows you to specify a BlackPixel color.
- n** **dpy**
Specifies the connection number of the server.
- layer <int> layer**
Specifies the requested layer for the default visual if the ddx supports overlays and underlays.

Processing of Displays Specified on the Command Line

The following flags can be used to specify the display on the command line. If a display is not specified, by default, the X server is invoked on all available displays in a multihead configuration by slot order. The first active display found becomes the equivalent of `-P11` (Screen 0). The next active screen found becomes the equivalent of `-P12` (Screen1).

- force** Specifies that the X server may be invoked from tty.
- P Row Col DisplayName**
Specifies that the X server is to be invoked on the display specified by the `display_name` parameter. Valid values for this field are those found in the `logical_name` field of the **CuDv**.

In all of the flags, the `display_name` parameter is specified by the logical name of the device. You can find out what this logical name is by executing the **lsdisp** command.

Note: Only the **ProcessCommandLine** and **ddxProcessArgument** subroutines should process command-line flags.

FindAllAvailableDisplays

Purpose

Initializes the **DisplayRec** structure for the X server.

Syntax

```
void FindAllAvailableDisplays (argc, argv)
int argc;
char **argv;
```

Description

The **FindAllAvailableDisplays** subroutine queries the ODM to find out information about all possible graphics displays that are available. It returns an initialized **DisplayRec** structure for each display and sets the **NumAvailableDisplays** global variable to the number of displays attached to the system.

The information needed to obtain a list of all available displays is contained in two ODM classes, the **PdAt** (predefined attributes) and the **CuDv** (customized devices). To find all the active displays attached to a system, two ODM queries are needed—one to find all the graphics devices, and one to subset the available devices from all the graphics devices.

For each display attached to the hardware, the **FindAllDisplays** subroutine initializes the following information:

- dev_id** Specifies the GAI adapter ID. This is used as a key to the ODM query of the GAI load modules. This value is stored in the predefined attribute **display_id**.

name	Set to the name of the display as returned by the lsdisp command (i.e. ppr0). This value is the logical name of the device as retrieved from the PdAt query of all the graphics devices.
loadModule	Set to the loadddx load module for the specified adapter. It is a concatenation of the Module_path as stored in the GAI class and the GAI_PATH . The value stored is the full path not the partial path as retrieved from the GAI.
EntryFunc	Set to NULL . This gets initialized when the loadddx load module is actually loaded.
x_pos	Set to 0. This is initialized to the appropriate value later in the initialization cycle.
y_pos	Set to 0. This is initialized to the appropriate value later in the initialization cycle.
requested	Set to False . Set to True later in the initialization cycle when the adapter is actually requested.
display_number	Set to the current value of the variable used to control the for loop.

Parameters

The *argc* and *argv* (count and vector) command-line arguments are passes in for processing of the **-force** flag. The **-force** flag allows the X server to be started from a tty.

InitOutput Subroutine

Purpose

Initializes the **screenInfo** *imageByteOrder* bitmap information, and pixmap format information.

Syntax

```
void InitOutput (screenInfoPtr, argc, argv)
ScreenInfo *screenInfoPtr;
int argc;
char *argv[];
```

Description

The **InitOutput** subroutine initializes the **screenInfo** *imageByteOrder* bitmap information, and pixmap format information. The **InitOutput** subroutine obtains a private index for the screen and one for the AIX extensions. These indexes are global variables named **aixExtScreenPrivateIndex** and **aixScreenPrivateIndex**. They are used as private indexes for the structure containing the vectors for the AIX private extensions (cursor and colormap) and the AIX screen private structure and vectors.

If there were no screens specified on the command line, the subroutine processes the default screens. The default screens are all currently available screens, the first screen found in the ODM.

The **InitOutput** subroutine opens the **/dev/rcm** file once for each invocation of the X server and not done once per head in a multihead environment. The subroutine stores the file descriptor from the open in the **screenInfoPriv** structure.

For each screen, the **InitOutput** subroutine calls subroutines to load and initialize or reinitialize the adapter. The **serverGeneration** global variable can be used to determine if the X server is being invoked or reset.

Parameters

<i>screenInfoPtr</i>	Specifies a pointer to the screenInfo structure. The screenInfo structure contains a section of per server information as well as an array of screen information.
<i>argc</i>	Specifies the number of options in the <i>argv</i> list.
<i>argv</i>	Specifies an array of options that were passed to the X server at its invocation.

Device-Dependent Initialization Subroutines

Note: These routines must be provided by the **loadddx** module.

xxxentryFunc Subroutine

Purpose

Returns information from the loading of the **loadddx** module.

Syntax

```
int xxxentryFunc (index, pScreen, argc, argv)
register int index;
register ScreenPtr pScreen;
int argc;
char **argv;
```

Description

The address of the entry point of the loadable driver is returned when the **load** system call is performed on the **loadddx** module. The function is invoked in the **InitOutput** subroutine and performs some or all of several actions. This function must be provided by the load module driver.

The **xxxentryFunc** function initializes any appropriate pixmap formats and then calls **AddScreen** subroutine to allocate the per screen structure and provide addition dix initialization of the **screenInfo** structure. The **AddScreen** subroutine allocates the screen structure, links it appropriately into dix structures and then calls a ddx supplied routine to finish the initialization.

A pointer to **xxxentryFunc** is stored in the *EntryFunc* field of the **DisplayRec** structure for this adapter. In practice, this function can be called whatever the driver implementor chooses. By convention it is referred to as **xxxentryFunc**.

Parameters

<i>index</i>	Specifies the index of this screen in the array of screens for the X server.
<i>pScreen</i>	Specifies a pointer to the screen structure.
<i>argc</i>	Specifies the number of options in the <i>argv</i> list.
<i>argv</i>	Specifies an array of options that were passed to the X server at its invocation.

Return Values

Success
Failure

xxxScrInit

Purpose

Initializes the screen structure for each screen and performs device dependent initialization as required.

Syntax

```
Bool xxxScrInit (index, pScreen, argc, argv)  
register int index;  
register ScreenPtr pScreen;  
int argc;  
char **argv;
```

Description

There is one **xxxScrInit** subroutine for each adapter load module. Its address is passed to the **AddScreen** dix subroutine and is called from the **AddScreen** subroutine through the function pointer to do device dependent initialization of the per-screen structure as well as window and gc devPrivates. The **xxx** in the subroutine name is replaced with some meaningful prefix, such as the adapter name. This subroutine performs the following:

- Gets a **gscHandle** from the RCM by making an ioctl call with a GSC_HANDLE parameter.
- Calls the **CreateAdapter** function for the 2D model.
- Registers this process as a graphics process by making the **aixgsc** system call with the MAKE_GP parameter.
- Initializes **pScreen** structure. This includes any visual structures for the screen.
- Allocates and initializes **pScreenPrivate** vectors (calls **aixAllocatePrivates**) and initializes AIX extension vectors—these are the vectors for the cursor and colormap extensions.
- Allocates and initializes any internally used private structures.
- Performs any other device-specific initialization.

Parameters

<i>index</i>	Specifies the index of this screen in the array of screens for the X Server.
<i>pScreen</i>	Specifies a pointer to the screen structure.
<i>argc</i>	Specifies the number of options in the <i>argv</i> list.
<i>argv</i>	Specifies an array of options that were passed to the X Server at its invocation.

Return Values

True	xxxScrInit succeeded
False	xxxScrInit failed.

xxxCloseScreen

Purpose

Closes the screen during X Server reset.

Syntax

```
Bool pScreen→CloseScreen (scrnNum, pScreen, argc, argv)  
int scrnNum;  
ScreenPtr pScreen;  
int argc;  
char **argv;
```

Description

The **xxxCloseScreen** subroutine is a device-dependent routine that must be provided by the driver implementor. This subroutine is called for each initialized screen when the server is reset or closed down.

The **xxxCloseScreen** subroutine should free any privately allocated structures. If dynamically allocated space is not freed, the driver will cause the X server to leak memory across server resets. In addition, a process similar to that outlined in the **xxxScrInit** subroutine must be followed to relinquish access to the adapter.

Whereas in the **xxxScrInit** routine the MAKE_GP subfunction of the aixgsc system call was called to gain access to the adapter, the UNMAKE_GP subfunction of **aixgsc** must be called in the **xxxCloseScreen** routine to relinquish access to the adapter. Failure to do so will prevent the success of subsequent **aixgsc** MAKE_GP calls in **xxxScrInit** if the server is resetting and not actually terminating.

Parameters

<i>scrnNum</i>	Number of screen
<i>pScreen</i>	Specifies the screen to be closed.
<i>argc</i>	Number of command line parameters
<i>argv</i>	Command line argument vector

Server Termination

Each implementation of the X server must provide two routines that terminate the X server. The **ddxGiveUp** subroutine is called when the X server terminates normally. The **abortDDX** subroutine is called when the X server encounters an unrecoverable error. Both of these subroutines close the screen for each active screen of the display. There is no change in either subroutine.

Adapter Access and the aixgsc System Call

Because the X server is a user process, special considerations have to be taken to ensure that the DDX is able to access the adapter. This is done through the **aixgsc** system call with the MAKE_GP parameter.

The MAKE_GP parameter marks the calling process as a graphics process for the specified adapter and returns a segment base address of the adapter and a device-specific structure that should contain device-specific offset addresses. The device-specific structure may also contain additional information.

The device-specific portion of the information returned by the **aixgsc** system call is actually set in the low-level device driver **make_gp** routine. This is a hardware-specific function that is a member of the **phys_display** structure. (See “Display Device Driver” on page 11-14.) This **make_gp** routine fills in addresses and other information required by the DDX.

Once the **make_gp** routine of the display device driver has been properly implemented, the **aixgsc** system call will return the correct addresses within the loadable DDX module. This enables the X server to access the graphics adapter.

Implementation Details

The **aixgsc** system call takes a pointer to the following structure (defined in **/usr/sys/include/aixgsc.h**) as an argument:

```
typedef struct _make_gp {
    int          error;          /* error report */
    caddr_t      segment;       /* segment base address */
    genericPtr   pData;        /* device specific adapter addresses */
    int          length;        /* length of device specific data */
    int          access;        /* access authority */
    # define SHARE_ACCESS 0     /* process shares access to adapter */
    # define EXCLUSIVE_ACCESS 1 /* process cannot share access */
} make_gp;
```

The following is a sample **make_gp** routine for a display device driver:

```
static long xxx_make_gp(pdev, pproc, map, len, trace)
gscDev *pdev;
rcmProcPtr pproc;
struct xxx_map *map;
int len;
int *trace;
{
    struct ddsent * dds;

    /* This routine fills in the register mappings for this instance */
    /* returns to the gsc call. */

    dds = (struct ddsent *) pdev->devHead.display->odmdds;
    map->io_addr = dds->io_bus_addr_start;
    map->cp_addr = dds->io_bus_mem_start;
    map->vr_addr = dds->vram_start;
    map->dma_addr[0] = dds->dma_range_start;
    map->screen_height_mm = dds->screen_height_mm;
    map->screen_width_mm = dds->screen_width_mm;
    return(0);
}
```

With the above **make_gp** routine, the following is a typical fragment of code to bind the graphics adapter load module into the GAI architecture and allow access to the adapter within the **xxx_ScrInit**:

```

#include <gai/gai.h>
#include <gai/misc.h>
#include <sys/rcm_win.h>
#include <sys/aixgsc.h>
#include <sys/rcmioctl.h>
#include <aixPrivates.h>
extern gAdapterPtr gaiDevAdpIndex[MAXSCREENS];
gsc_handle gscHandle;
Bool xxxScrInit(int index, ScreenPtr pScreen, int argc, char **argv)
{
    gAdapterPtr pAdapter = NULL;
    make_gp makegp_info;
    struct xxx_map xxx_makegp_info;      /* see RCM */
                                        /* documentation */

    int count = 0;
    char *cmdlist[1] = {(char *) NULL};
    strcpy (gscHandle.devname, aixScreenName(pScreen));
    if (ioctl (aixScreenFD(pScreen), GSC_HANDLE, &gscHandle) < 0) {
        return (FALSE);
    }
    pAdapter = CreateAdapter (aixDeviceID(pScreen),
                             gscHandle.handle, count, cmdlist);
    makegp_info.error = 0;
    makegp_info.pData = (char *) &xxx_makegp_info;
    makegp_info.length = sizeof(struct xxx_map);
    makegp_info.access = EXCLUSIVE_ACCESS;
                                        /* not a direct access device */
    if (aixgsc(gscHandle.handle, MAKE_GP, &makegp_info)) {
        return (FALSE);
    }
    vram_bits = makegp_info.segment +
                ((struct xxx_map *) makegp_info.pData)->vr_addr;
    /* continue with remainder of adapter initialization ,
       return TRUE on success, FALSE otherwise. */
}

```

Minimum Resource Management Subsystem (RMS) for 2D Adapters

The Resource Management Subsystem (RMS) is responsible for managing resources common to different rendering models. These resources include cursors, color palettes, buffers, client clip information, fonts, window geometries, and so on. It is not necessary to provide full-function RMS support for an adapter that supports only the X rendering model, as X handles its own resources.

In order to provide the minimal RMS support required by the Graphics Adapter Interface (GAI) architecture, the 2D graphic adapter load module must call the RMS **CreateAdapter** function in the adapter initialization subroutine and the RMS **DestroyAdapter** function in the adapter uninitialize or destroy subroutine.

Implementation

Adapters that are not direct access adapters and will never be direct access adapters can write an abbreviated version of the Resource Management Support. The following calls must be made by a 2D adapter that is intended to display any of the SOFT 3D APIs, such as OpenGL or SoftPHIGS. Otherwise, these calls are recommended, but not required.

- CreateAdapter** This routine is the basic initialization routine of the GAI RMS library. It provides the necessary initialization to create a graphics process within the RCM (Rendering Context Manager) and permits the graphics process to begin configuration and operation. The routine is typically called in the ddx specific adapter initialization routine.
- DestroyAdapter** This routine destroys all data structures allocated for the adapter and conveys information to the RCM that invalidates the graphics process' adapter access privileges. This routine is typically called in the ddx specific routine vectored through when the X server terminates or resets.

Adapters that use the minimal RMS will not physically have an RMS load module and will store the value **NORMS** in the rms ODM entry.

The calling sequence for the **CreateAdapter** and **DestroyAdapter** functions are as follows:

```
gAdapterPtr CreateAdapter(int aid, GSC_HANDLE gsc_handle, int
argc, char** argv);
```

The **CreateAdapter** function is called from the **xxxScrInit** subroutine.

```
pAdapter->pProc->DestroyAdapter(pAdapter);    /* pAdapter is */
/* the adapter pointer returned from CreateAdapter */
```

The **DestroyAdapter** function is typically called from the **xxxCloseScreen** routine.

Include Files

In addition to the standard X server include files, the compilation units (c source files) which contain the implementation of the device dependent initialization and uninitialization routines must include the following header files:

- <sys/rcm_win.h>
- <gai/gai.h>
- <gai/misc.h>
- <sys/aixgsc.h>
- <sys/rcmioctl.h>
- <X11/ext/aixPrivates.h>

These header files resolve type definitions for GAI specific data types and external function declarations for GAI-specific functions used by 2D graphics adapter load modules and input adapter load modules. In addition, they make available the use of the AIX Graphics System Call (aixgsc). At initialization time, an aixgsc call must be made to the device independent portion of the RCM to register the process (in the case of a 2D graphics adapter ddx module bound with the X server, the X process) as a graphics process.

Configuring the 2D Adapter into the ODM Database

Information about where loadable modules exist is stored in the ODM database. These records are located in the GAI class of the database in **/usr/lib/objrepos**. At load time, this database is queried based on a unique adapter ID for the correct modules to load. There are two entries that are required for 2D graphics adapters: one for the RMS module and one for the DDX module.

Since it is not necessary to provide full-function RMS support for 2D adapters (see previous section), a special keyword "NORMS" has been developed to provide the minimal level of RMS support; "NORMS" should be placed in the Module_Path field of the ODM entry.

The following is the structure definition for the GAI record in the ODM database:

ODM Record

```
class GAI {
    long Adapter_Id;
    char Module_Key[11];
    vchar Module_Path[256];
    vchar Processor_Id[16];
};
```

Field definitions for the GAI record in the ODM database.

Adapter_id Specifies the GAI adapter ID for the adapter. New adapters are assigned this value during development. This is obtained from the `display_id` attribute of the ODM PdDv Object Class and stored in the `dev_id` field of the appropriate **display_rec** structure when the server initializes.

The ID is given a value of the form 0x04xxmmnn, where:

04	Fixed.
xx	An adapter-specific ID.
0x00 – 0x7f	Reserved for IBM use.
0x80 – 0xff	Provided for vendor adapters.

Note: To avoid duplication of adapter IDs within this range, it is advised that you contact IBM for a list of available IDs and not rely on using values not currently installed in the database.

mm	Specifies the adapter state for vendor adapters. 0x00 indicates that the adapter is functional. No other value is allowed at this time.
-----------	---

nn	Differentiates between multiple instances of the same adapter type (max 4).
-----------	---

Module_Key Specifies the key for the module. Typically, this key is “rms”, “ddx” or one of the 3D or extension keys.

Module_Path Specifies the path name of the load module for this key and adapter. The module path is used along with the **GAI_PATH** environment variable to find the load module.

Processor_Id Identifies the processor. Intended for future use to enable the developer to load modules based on processor type. Currently, the only supported value of this field is 0 (zero), the default, defined as ANY_PROCESSOR.

Usually, ODM records are gathered in a file to be added at configuration time. For example, a `samp2d.add` file for an adapter with an ID of 0x4330000 may contain the following records:

```
GAI:
Adapter_Id      = 0x4330000
Module_Key      = "rms"
Module_Path     = "NORMS"
Processor_Id    = 0
```

```
GAI:
Adapter_Id      = 0x4330000
Module_Key      = "ddx"
Module_Path     = "sampddx/loadddx"
Processor_Id    = 0
```

Porting Input Devices

This section describes how to add a new input device to the X Server using the X11 Input Extensions. It consists of the following areas:

- Input Device Driver Overview (on page 11-48)
- Block and Wakeup Handling (on page 11-51)
- Event Processing (on page 11-53)
- Input Load Module (on page 11-54)
- ODM Database Entry for Input Devices (on page 11-60)
- Sample Input Device Load Module (on page 11-61)

Input Device Driver Overview

The AIX system supports only keyboards and displays as part of the low-function Terminal (LFT) subsystem. Typically, the X server uses the mouse as its core pointer, but it could also be a tablet or other device that has a valuator with at least two axes and one pointer. Devices with valuators of two or more axes are supported on the system as extension input devices. The AIX system provides built-in load modules for supported extension input devices. Unsupported input devices require their own load modules.

The X server accesses extension input devices through the X11 Input Extension. This information describes adding a new input device to the AIXwindows X server. You must understand how the X11 Input Extension works before trying to port a new input device. This information should be used with the *X11 Input Extension Porting* document from the X Consortium.

The three requirements for adding a new extension input device to the AIXwindows X server are:

- A device driver (or standard tty device driver)
- A load module that contains the device-dependent subroutines for accessing the extension input device
- Entries in the Object Data Manager (ODM) database.

Device Driver

The device driver is the interface between the X server and the extension input device hardware. If the extension input device can use the serial port, you may want to use the standard tty device driver rather than write a device driver. The requirements for the device driver vary, depending upon whether you want to use the X server input ring. If you use the standard tty device driver, you may need to maintain your own input ring. For information on using the tty device driver, see Chapter 10 “STREAMS-Based TTY Subsystem Interface.”

To use the X Server input ring to store events, follow the specification for adding events to the input ring. The device driver must send a signal to the X server to notify it that there is input pending. The load module for the extension input device is responsible for sending the input ring information of the X server to the device driver. The **deviceProc** subroutine must initialize the **inputInfoDevPriv** structure for the extension input device.

If you do not want to use the X server input ring, then the load module and the device driver must coordinate how and where input device events are to be stored. The load module must register with the X server the two values it will use to check for input pending. The module must also determine the subroutine to call if input is available. This can be done using

AddInputCheck subroutine. You will need to follow this method if you plan on using the standard tty device driver. You must also register a block handler and a wakeup handler.

InputInfo Structure

The **InputInfo** structure, defined in the `/usr/include/X11/ext/inputstr.h` file, contains information about the input devices. The **inputInfo** global variable must be used to access the elements of the structure.

The **InputInfo** structure is defined as follows:

```
typedef struct {
    int          numDevices; /* total number of devices */
    DeviceIntPtr devices;   /* all devices turned on */
    DeviceIntPtr off_devices; /* all devices turned off */
    DeviceIntPtr keyboard;  /* fast path to keyboard device */
    DeviceIntPtr pointer    /* fast pointer to the pointer */
                          /* device */
    #ifdef AIXV4
        inputInfoDevPriv infoPriv; /* the global input specific */
                                   /* control block */
    #endif
} InputInfo;
```

The fields of the **InputInfo** structure are defined as follows:

numDevices Specifies the number of devices that are known to the X server.

devices Specifies a pointer to a linked list of device structures that represent devices that are turned on.

off_devices Specifies a pointer to a linked list of device structures that represent devices that are turned off.

keyboard Specifies a pointer that is a fast path to the core keyboard structure.

pointer Specifies a pointer that is a fast path to the core pointer structure.

infoPriv Specifies a pointer to the private structure for the input control blocks.

inputInfoDevPriv Structure

The **InputInfo** structure contains a private area that is used for input control blocks. The following structure is used as the **InputInfo** private area. The **MAX_DEVICES** variable is defined in the `inputstr.h` file. Use the **inputInfo** global variable to access this structure. The **devices** and **eventHandlers** elements must be set in the **deviceProc** subroutine if the X server input ring is to be used to store the device events.

```
struct _inputInfoDevPriv {
    struct inputring *inputRing;
    struct DeviceIntRec *devices [MAX_DEVICES] ;
    void (*eventHandlers [MAX_DEVICES]) ();
} inputInfoDevPriv, *inputInfoDevPrivPtr;
```

```
#define INPUT_RING_SIZE 4096
```

The fields of the **inputInfoDevPriv** structure include are defined as follows:

inputRing Specifies a pointer to the input ring where the input devices store events. The **inputring** structure is defined in system include file `/usr/include/sys/inputdd.h`.

devices Specifies an array of pointers to the **DeviceIntRec** structure for each device and serves as a fast path to the structures for each device. The

index into the array is the device ID stored in the **DeviceIntRec** structure for the extension input device.

`eventHandlers`

Specifies an array of subroutines (one for each possible device) that format the raw event as read from the input ring into the proper X server format. The index into the array is the device ID stored in the **DeviceIntRec** structure for the extension input device.

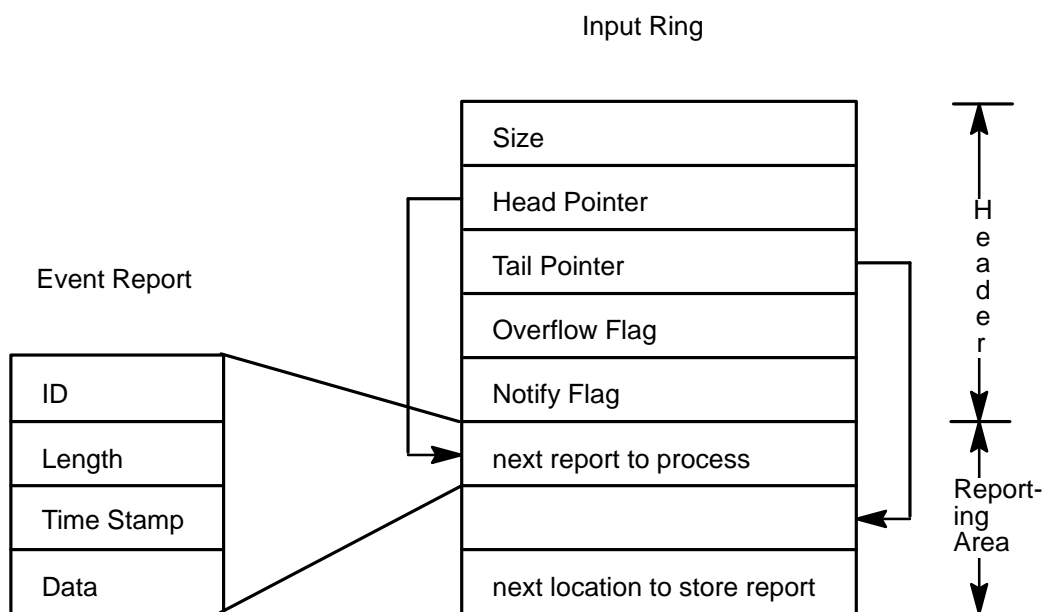
inputring and ir_report Structures

The **inputring** and **ir_report** structures are defined in the **inputdd.h** include file. Refer to the **inputdd.h** file for details of the structure. The first 10 bytes in each **ir_report** structure contain the ID of the device returning the event, the size of the event, and a time stamp. The Input Ring Structure figure depicts the location of this data within the context of the input ring.

X Server Input Ring

Once a device has been opened, extension input device drivers can add events to the X server input queue. The input ring is allocated in X server memory and is stored as part of the device private area of the input control block in the **inputInfo** structure. To use the input queue, the **deviceProc** subroutine for each extension input device must register the input ring data with the device driver. This data is the pointer to the **inputring** structure and the device ID. The device ID for each device is the `id` field of its **DeviceIntRec**.

The following diagram illustrates the X Server input ring:



Input Ring Structure

All input devices that receive input events from device drivers by this mechanism share the input ring. The input queue contains a head and tail pointer. The device driver adds device events to the ring at the tail pointer, and the X server reads events from the ring at the head pointer. The header of each raw device event contains a device ID, the size of the event data, and a timestamp in milliseconds. The offset to the next event is computed from size value. The locking strategy uses priority levels to ensure that only one input driver touches

the input ring at any given time. Thus, all input device drivers must execute with an **INTCLASS3** priority whenever they modify the input ring.

When the queue head pointer is equal to the tail pointer, it is assumed that the queue is empty. When the device driver attempts to add an event that would cause the tail pointer to be equal to or logically greater than the head pointer, it sets the **ir_overflow** flag and flushes the input event that caused the overflow. Additionally, all events are flushed as long as the **ir_overflow** flag in the ring header is set to **IROVERFLOW**, even if there is room on the ring for an event report.

Data stored in the input ring are in the form of raw device events, but they must be reformatted into the type of event that the X server returns to the client. The events may be core device events (those originating from a pointer or keyboard) or extension device events, such as valuator motion, button press, button release, key press, key release or proximity. Since the X server may not have knowledge of every raw device event format stored in the queue, it calls a device-specific subroutine to reformat the event into one of the event types that the X server understands. These subroutines are stored in the input device private structure in the **processRawInputEvents** field.

A fast path to each of the extension input device's **processRawInputEvents** is the **eventHandlers** array stored in the private area of the **inputInfo** structure. When the device is opened, the **deviceProc** subroutine stores a pointer to the device's **processRawInputEvents** subroutine in the **eventHandlers** index that corresponds to the device's ID as stored in the **id** field of the **deviceIntRec** structure. The **processRawDeviceEvents** formats the event into the appropriate X events and calls the subroutine stored in the **processInputEvents** field of the **DeviceIntRec** for the extension input device.

SIGMSG Signal

Input device drivers using the AIX input event queue provide event notification to the X server using the **SIGMSG** signal. The **SIGMSG** signal is defined in the **/usr/include/sys/signal.h** file. The **ir_notifyreq** flag in the input ring header specifies when the **SIGMSG** signal is to be sent. If the flag is set to **IRSIGEMPTY**, the driver sends a signal only when it enqueues an event report into an empty input ring. If the **ir_notifyreq** flag is set to **IRSIGALWAYS**, then the device driver sends a **SIGMSG** signal each time it enqueues an event.

Notification is provided for the event that causes a queue overflow; however, once queue overflow has been posted to the input ring, no other event notification is provided until the overflow indicator has been cleared. Queue overflow is posted to the input ring when the **ir_overflow** flag is set to **IROVERFLOW**.

If the input device to be added connects to the system through a tty port, a block and wakeup handling scheme must be used. See "Block and Wakeup Handling" for more information.

Block and Wakeup Handling

If you use the standard tty device driver, you must notify the X server to enable the extension input device driver for input pending. This can be accomplished by calling the **AddEnabledDevice** subroutine, which adds the file descriptor for the device driver to the select mask. To remove the file descriptor for the device driver from the select mask, use the **RemoveEnabledDevice** subroutine. Both subroutines should be called in the **deviceProc** subroutine for the extension input device when the extension input device is turned on or off. The **xxxBlockHandler** and **xxxWakeupHandler** subroutines must be

registered with the X Server by calling the **RegisterBlockAndWakeupHandlers** subroutine to control what the X Server should do before and after calling the **select** subroutine.

Note: You must replace the **xxx** portion of the subroutine name to make it unique.

The **xxxBlockHandler** and **xxxWakeupHandler** subroutines are removed by calling **RemoveBlockAndWakeupHandlers** when the extension input device is turned off in a **deviceProc** subroutine of the device.

The **AddEnabledDevice**, **RemoveEnabledDevice** and **RemoveBlockAndWakeupHandlers** subroutines are common X Server routines. For more information about these subroutines, refer to *The X-Window Servey* by Elias Israel and Erik Fortune, Digital Press.

Additional information about the **xxxBlockHandler** and **xxxWakeupHandler** subroutines follows.

xxxBlockHandler Subroutine

Purpose

Allows for any device-specific action before the **select** subroutine is processed.

Syntax

```
void xxxBlockHandler (blockData, ppTimeout, pReadmask)
char *blockData;
struct timeval **ppTimeout;
unsigned *pReadmask;
```

Description

The **xxxBlockHandler** subroutine allows for any device-specific action before the X server processes the **select** subroutine. The **xxxBlockHandler** subroutine is called before the X server calls the **select** subroutine. For most extension devices, this subroutine is an empty routine. You must register block and wakeup handlers in pairs, and in most cases, only a wakeup handler is needed for input devices.

Parameters

<i>blockData</i>	Points to the private data used by the block handler subroutine. Its value is the same as the value of the blockData parameter that was passed to the RegisterBlockAndWakeupHandlers subroutine.
<i>ppTimeout</i>	Returns a pointer to the timeout value. This value is used to determine the maximum amount of time the block is allowed to last.
<i>pReadmask</i>	Points to a data structure that defines which descriptors are to be blocked on. The data structure is an array of unsigned long data elements, and the bit in the array must be set so that it corresponds to the value of the file descriptor.

xxxWakeupHandler Subroutine

Purpose

Allows for reading from ready file descriptors after **select** processing.

Syntax

```
#define FD_SETSIZE 128 /* limit of server connections */
#include <sys/select.h>
void xxxWakeupHandler (result, pReadmask)
```

```
int result;
unsigned *pReadmask;
```

Description

The **xxxWakeupHandler** subroutine is called after the X server returns from the **select** subroutine. If the descriptor for the extension input device driver is ready for reading, the two long values that check for input should be set not equal to each other.

The macro **FD_ISSET(*n*, *pReadMask*)** checks the file descriptor that you are using to see if the **select** subroutine selected your ID.

Parameters

<i>result</i>	Indicates the return value from the select subroutine that specifies the number of descriptors that are ready for reading.
<i>pReadMask</i>	Points to a data structure that defines which descriptors are ready for reading.

Event Processing

Processing events involves obtaining events from devices and delivering them to clients. If the X server input ring is used, each device event in the input ring is passed to the **eventHandlers** subroutine, which was initialized in the **inputInfoDevPriv** structure. Otherwise, the input processing subroutine passed to the **AddInputCheck** subroutine will be called to process pending input. Both of these methods should use the **processRawEvents** subroutine, which is part of the load module.

AddInputCheck Subroutine

Purpose

Adds additional input checking for extension input devices that are not using the X server input ring.

Syntax

```
void AddInputCheck (p1, p2, proc)
void *p1, *p2;
void *proc ();
```

Description

The **AddInputCheck** subroutine adds extra input checking for extension input devices that do not use the X server input ring. The **AddInputCheck** subroutine adds a level of indirection to the current scheme of checking for input by comparing of two long data elements. Instead of one set of long data elements to check for pending input, there is an array of longs. The **AddInputCheck** subroutine adds a new set of long data elements as well as an input processing procedure into this array. Two identical values indicates no input whereas two different values indicates input is pending. The long values can be pointers or hard-coded values.

Parameters

<i>p1</i> , <i>p2</i>	Specifies a pointer to a long data element. The values can be a mask, a hard coded value, an integer, or a pointer.
<i>proc</i>	Specifies a procedure to call if the values pointed to by the <i>p1</i> and <i>p2</i> parameters are not equal.

RemoveInputCheck Subroutine

Cancels the **AddInputCheck** subroutine.

Syntax

```
void RemoveInputCheck (proc)  
void *proc ();
```

Description

The **RemoveInputCheck** subroutine removes additional input checking for extension input devices that do not use the X server input ring. This subroutine is called when the **deviceProc** subroutine of the device contains the **DEVICE_OFF** value.

Parameters

proc Specifies the procedure that corresponds to the input check that should be removed.

Input Load Module

The subroutines that manipulate each extension input device are stored in a dynamically loadable module. When the X server receives an **X_ListInputDevices** protocol request, it loads each of the extension load modules and calls the entry point for each load module. Load modules remain loaded until the X server terminates or resets.

InputDevPrivate structure

The **InputDevPrivate** structure must be initialized in the entry point subroutine of the extension input device load module. This structure is part of the **/usr/include/X11/ext/aiXInput.h** file and is defined as follows:

```
typedef struct {  
    pointer    client_list;  
    short     ctrlstructsize;  
    int       (*legalModifier) ();  
    int       (*deviceProc) ();  
    int       (*openDevice) ();  
    void      (*closeDevice) ();  
    int       (*setDeviceMode) ();  
    int       (*setDeviceValuators) ();  
    int       (*getDeviceControl) ();  
    int       (*changeDeviceControl) ();  
    void      (*processRawInputEvents) ();  
    int       fd;  
    int       idx;  
    pointer    vendor_specific;  
} InputDevPrivate;
```

Field definitions for the **InputDevPrivate** structure:

client_list Specifies a linked list of all the clients that have this device open. This element must be initialized to NULL.

ctrlstructsize Specifies the size of the device control structure that this device accepts. This field is used in the **ChangeDeviceControl** protocol. If the extension input device does not support changing controls, this element must be set to 0.

<code>legalModifier</code>	Specifies a pointer to a subroutine that is used for devices that have keys. This field returns True or False if the key passed as an input parameter can be used as a modifier key. If the extension input device does not use a modifier key, this element must be set to False .
<code>deviceProc</code>	Specifies a pointer to the deviceProc subroutine. The deviceProc subroutine initializes the device, turns the device on and off, and closes it.
<code>openDevice</code>	Specifies a pointer to the subroutine to call when the OpenDevice protocol is specified. The initial value of this field is NoopDDA() .
<code>closeDevice</code>	Specifies a pointer to the subroutine to call when the CloseDevice protocol is specified. The initial value of this field is NoopDDA() .
<code>setDeviceMode</code>	Specifies a pointer to the subroutine to call when the SetDeviceMode protocol is specified. If the extension input device does not support mode changing, this element must be set to NULL.
<code>setDeviceValuators</code>	Specifies a pointer to a subroutine to call when the SetDeviceValuators protocol is specified. If the extension input device does not support valuators or initializing of valuators, this element must be set to NULL.
<code>getDeviceControl</code>	Specifies a pointer to a subroutine to obtain the current device control setting for the device. If the extension input device does not support querying of device controls, this element must be set to NULL.
<code>changeDeviceControl</code>	Specifies a pointer to a subroutine to set the device control settings for a device. If the extension input device does not support changing device controls, this element must be set to NULL.
<code>processRawInputEvents</code>	Specifies an input subroutine that formats the raw events extracted from the input ring.
<code>fd</code>	Specifies the file descriptor of the extension input device.
<code>idx</code>	Specifies an index to device-specific information. Do not initialize this field.
<code>vendor_specific</code>	Specifies a pointer to any vendor private data used in the load module of the extension input device.

ExtlnitInput Subroutine

Purpose

Initializes the input private structure for the extension input device.

Syntax

```
int ExtlnitInput (inputDevPrivate)
InputDevPrivate *inputDevPrivate;
```

Parameters

inputDevPrivate Specifies a pointer to an empty input device private structure and returns an initialized structure.

Description

The **ExtlnitInput** subroutine for the extension input device initializes the **devPrivate** structure of the **deviceIntRec** structure. The **ExtlnitInput** subroutine is the entry to the load module and is returned from the **load** system call.

Return Values

0	Indicates successful completion.
1	Indicates an error occurred.

deviceProc Subroutine

Purpose

Initializes, turns on, opens, or closes the extension input device.

Syntax

```
int xxxdeviceProc (pDev, onoff)
DevicePtr pDev;
int onoff;
```

Description

The **deviceProc** subroutine is a device-specific procedure that is called from the **EnableDevice**, **DisableDevice**, **InitAndStartDevices**, **CloseDevice**, and **ExtlnitInput** subroutines. This subroutine performs the function specified by the *onoff* parameter for the device specified by the *pDev* parameter.

Typical actions that occur during device initialization include setting up device private information and calling the **dix** subroutines to initialize device-specific information. These subroutines include **InitializePointerDeviceStruct** and **InitializeKeyboardDeviceStruct** to initialize the pointer and keyboard. The **deviceProc** subroutine must initialize extensions, device class structures, and the global state of the device. The extension class device structures are created by using the following subroutines:

- **InitButtonClassDeviceStruct**
- **InitKeyClassDeviceStruct**
- **InitValuatorClassDeviceStruct**
- **InitValuatorAxisClassDeviceStruct**
- **InitFocusClassDeviceStruct**
- **InitProximityClassDeviceStruct**
- **InitKbdFeedbackClassDeviceStruct**
- **InitPtrFeedbackClassDeviceStruct**

The **MakeAtom** subroutine is called with the name of the device. The **AssignTypeAndName** subroutine is called with the device ID and the atom returned from the **MakeAtom** subroutine. These subroutines set the device name and type in the device structure.

The **deviceProc** subroutine must initialize the default values for global state in the **devPrivate** section of the **DeviceIntRec** structure for the specified device.

Parameters

<i>pDev</i>	Specifies a pointer to the structure that devices the device
<i>onoff</i>	Specifies the action to be taken. Valid values for <i>onoff</i> are: DEVICE_INIT Creates structures.

DEVICE_ON Opens the hardware device driver if one exists. The device is marked as on and calls the **AddEnabledDevice** subroutine. Other OS-specific actions occur to make the device available. Typically, the **deviceProc** subroutine calls the **AddEnabledDevice** subroutine to add the device to the list of “on” devices and the `on` field of the device structure is set to **True**.

If an extension device does not receive its events from input ring, then the **AddInputCheck** and **RegisterBlockAndWakeupHandlers** subroutines are called to set up the input checking and processing mechanism for the extension device.

DEVICE_OFF Marks the device off. Typically, **deviceProc** calls the **RemoveEnabledDevice** subroutine to delete the device from the list of “on” devices and adds the device to the list of initialized devices. The `on` field of the device structure is set to **False**.

If the **RegisterBlockAndWakeupHandlers** subroutine was called during the **DEVICE_ON** state, the **RemoveBlockAndWakeupHandler** subroutine turns off the block and wakeup handlers for the device.

DEVICE_CLOSE

If the **AddInputCheck** subroutine was called during the **DEVICE_ON** stage, the **RemoveInputCheck** subroutine turns off input checking and processing for the extension input device in the server.

Return Values

0	Indicates successful completion.
1	Indicates an error occurred.

setDeviceMode Subroutine

Purpose

Sets the mode of a device.

Syntax

```
int xxxsetDeviceMode (pDev, client, mode)
DeviceIntPtr pDev;
ClientPtr client;
int mode;
```

Description

The **setDeviceMode** subroutine sets the mode of the device to the mode passed as an input parameter. This mode can be **Absolute** or **Relative**.

Parameters

<i>pDev</i>	Specifies the device
<i>client</i>	Specifies the client opening the input device
<i>mode</i>	Specifies the mode. Valid values are Absolute or Relative .

Return Values

BadDevice	Indicates an invalid device ID.
BadMatch	Indicates the device does not support modes.
BadMode	Indicates the specified mode is invalid.
DeviceBusy	Indicates another client is using this device.

setDeviceValuators Subroutine

Purpose

Initializes the valuators on a device.

Syntax

```
int xxxsetDeviceValuators (pDev, client, valuators, first_valuator, num_valuators)  
DeviceIntPtr pDev;  
ClientPtr pClient;  
int *valuators;  
int first_valuator;  
int num_valuators;
```

Description

The **setDeviceValuators** subroutine initializes one or more of the valuators of the device to the values passed in the *valuators* parameter.

Parameters

<i>pClient</i>	Specifies the client setting the valuators the input device.
<i>pDev</i>	Specifies the device.
<i>valuators</i>	Specifies an array of valuator values to be set.
<i>first_valuator</i>	Specifies the index of the first valuator to be set.
<i>num_valuators</i>	Specifies the number of valuators to be set.

Return Values

Success	
BadMatch	Indicates the device does not permit valuators to be set.
BadValue	Indicates an invalid value for the valuator was specified.
BadDevice	Indicates an invalid device ID was specified.

getDeviceControl Subroutine

Purpose

Returns the current device control setting for the specified device.

Syntax

```
int xxxgetDeviceControl (pClient, pDev control)  
ClientPtr client  
DevicePtr pDev;  
xDeviceCtl *control;
```

Description

The **getDeviceControlProc** subroutine for each device obtains the specified controls for the device. The **xDeviceControl** structure is passed directly to the client. The client must link third-party vendor code and initialize the appropriate subroutines to interpret the data.

Parameters

pClient Specifies the client that is requesting the control change.
pDev Specifies a pointer to the device.
control Specifies the control to obtain.

Return Values

Success
BadValue Indicates the X server does not support specified control
BadMatch Indicates the X server supports specified control, but the device does not.
AlreadyGrabbed Indicates the device was grabbed by another client.

changeDeviceControl Subroutine

Purpose

Changes the device control settings for the specified device.

Syntax

```
int xxxchangeDeviceControl (pDev, control)  
DevicePtr pDev;  
xDeviceCtl *control;
```

Description

The control procedure for each device changes the specified controls.

Parameters

pDev Specifies a pointer to the device.
control Specifies the change to the device control.

Return Values

Success
Failure

processRawInputEvents Subroutine

Purpose

Formats a raw device event into an event that can be sent to an X Client.

Syntax

```
void xxxprocessRawInputEvents (pDev, rawEvents)  
DeviceIntRecPtr pDev;  
pointer *rawEvent;
```

Description

The **processRawInputEvents** subroutine takes the raw device event as retrieved from the input queue or input device and formats it into an event that can be sent to an X client.

In order to map raw device event into X event, information about the current screen and cursor position is needed. This type of information is accessible globally as follows:

```
extern int DeviceButtonPress, DeviceButtonRelease;

/* extension device event types */
extern int DeviceMotionNotify, DeviceValuator;

/* extension device event types */
extern int ProximityIn, ProximityOut;

/* extension device event types */
extern int lastEventTime;
extern ScreenInfo screenInfo; /* screen information */
extern int aixCurrentScreen; /* current screen number */
extern InputInfo inputInfo; /* global InputInfo,
                             defined in dix/devices.c */
extern int AIXCurrentX; /* current X position of cursor */
extern int AIXCurrentY; /* current Y position of cursor */
extern long GetTimeInMillis(); /* function to get current time */
```

Parameters

<i>pDev</i>	Specifies a pointer to the DeviceIntRec structure for the device.
<i>rawEvent</i>	Specifies a pointer to the raw device event for the device.

ODM Database Entry for Input Devices

Information about each input device is stored in the ODM (Object Data Manager) database. The following information will be stored for each device.

```
XINPUT {
    char DeviceName[64]; /* VENDOR'S NAME FOR DEVICE */
    vchar ModuleName[MAXPATHLEN] /* PATH NAME TO LOAD MOCULE */
    char GenericName[15]; /* ONE OF 18 GENERIC DEVICES */
    vchar Class[16]; /* PdDv CLASS */
    vchar SubClass;
    vchar connwhere;
};
```

The fields of the **XINPUT** ODM record structure are defined as follows:

DeviceName	Specifies a vendor specific name for the device.
ModuleName	Specifies the relative name of the load module. This name is used in conjunction with the XINPUTPATH environment variable to find the name of the XINPUT load module. The default path is set to: /usr/lpp/X11/lib/load .
GenericName	Specifies the name of the device as used in name field of the deviceIntRec structure. There are 18 predefined names defined in the XI.h include file. Vendors are encouraged to use these names for interoperability.
Class	Same as the class name in the device driver CuDv record.
SubClass	Reserved for future use.
connwhere	Reserved for future use.

ODM Input Device Record Example

The **.add** file for the IBM input devices will contain the information for the input device. For instance, the keyboard device will have the following record:

```
XINPUT:
DeviceName = "Keyboard Device"
GenericName= "KEYBOARD"
ModuleName = "loadkeyboard"
Class      = "Keyboard"
```

Sample Input Device Load Module

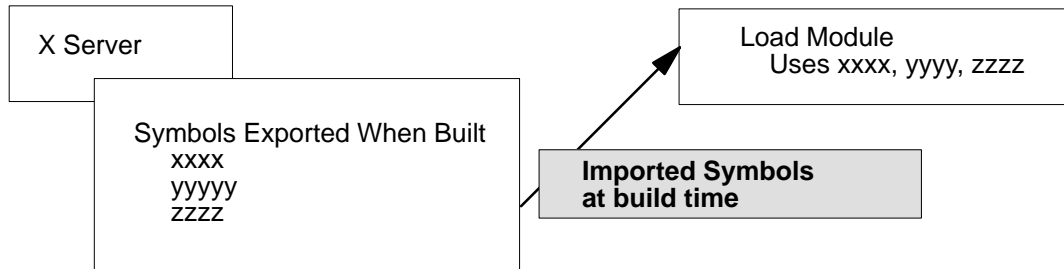
A sample Input Device load module is included in the AIX Version 4.1 distribution and is located in the directory **/usr/lpp/X11/Xamples/extensions/server/load**.

If this directory is not on your system, you may have not loaded the **X11.samples.ext** LPP. To see if this is on your system execute:

```
lslpp -h | grep samples.ext
```

Building a Dynamically Loadable Module

The figure shows how dynamically loadable modules are built.



Build of Dynamically Loadable Modules

To build a dynamically loadable module, the loading module and the loaded module must cooperate. The loading module must export a list of symbols that will be needed by the loadable module. The loadable module must import the symbols needed from the loading module. These symbols are resolved at the time of the **load** system call. In addition, each dynamically loadable module must provide the name of a subroutine to be used as the entry point to this module.

AIX provides the following exports files that can be used by a dynamically loadable module as import files:

- X.exp** Contains symbols from the device independent (dix) portion of the X server.
- Xi.exp** Contains symbols from the device independent portion of the X11 Input Extension.

The following example is a sample makefile used to build a dynamically loadable module. This module imports symbols from **X.exp** and **Xi.exp**. The **ExtlnitInput** subroutine is used as the entry point to this module:

```
# Makefile for creating loadable object module

# specify entry point
EPNT = SampExtInit

# import files
IMPS = -bI:/usr/lpp/X11/bin/X.exp \
       -bI:/usr/lpp/X11/bin/Xi.exp

# include file directories
INCS = -I/urs/include/X11          \
       -I/usr/include/X11/ext      \
       -I/usr/include/X11/extensions

sample:
       cc SampExt.c -o LoadSampExt -e${EPNT} ${INCS} ${IMPS}
-bM:SRE
```

Debugging Load Modules

Remember that symbols defined in the load module are not available to the debugger until the load module is loaded into memory. So, in order to set a breakpoint inside a function defined in the load module, you have to wait until the module is loaded. To do this:

- Get into the debugger (`dbx /usr/lpp/X11/bin/X`) and issue a **set** subcommand. This will list all the variables defined in the debug environment.

- Use the **unset** subcommand to unset the variable **\$ignoreload**:

```
unset $ignoreload
```

- Issue the **run** subcommand at the debug prompt:

```
run -P11 dev_name
```

where *dev_name* is the index of the device you want to debug from the **lsdisp** command (for example, `ppr0`). The process of loading and unloading generates signals which the debugger catches and pauses. At this time you set the breakpoint.

As an example, say that there are two displays on the system:

```
>lsdisp
DEV      SLOT    BUS      APT_NAME    DESCRIPTION
ppr0     00-01    mca      POWER_G4    Midrange Graphics Adapter
gda0     00-03    mca      colordga    Color Graphics Adapter
```

Get into the debugger:

```
dbx /usr/lpp/X11/bin/X          (user action)
dbx>unset $ignoreload          (user action)
dbx>run -P11 gda0              (user action)
```

At this point the LFT subsystem takes over the cursor and you will see a screen flash on the vendor display. Type Shift-Action to get the cursor back to the `ibmdisp` display.

```
dbx> stopped due to load/unload (system response)
dbx> stop in xyz                (user action, xyz() is
                                the function you want to
                                break in.)
```

List of X Server Porting Subroutines

Load modules for the AIXwindows X server must contain the following subroutines.

X Server Initialization

The following X Consortium subroutines must be included to meet X server requirements. One instance is required for each server.

ddxProcessArgument

Processes the command-line arguments for the X server.

FindAllAvailableDisplays

Initializes the **DisplayRec** structure for the X server.

InitOutput

Initializes the **screeninfo** *ImageByteOrder*, bitmap information, and pixmap format information and obtains private indexes for the screen and AIX extensions.

Device-Dependent Initialization

The following subroutines must be supplied to meet ddx requirements:

xxxentryFunc Returns information from the loading of the loadddx module.

xxxScrInit Initializes the screen structure for each screen and performs device-dependent initialization as required.

xxxCloseScreen

Closes the screen during X server reset.

Block and Wakeup Handling (Input Devices)

Third-party vendors need to write these subroutines for their input extension load modules.

xxxBlockHandler

Allows for any device-specific action before the select subroutine is processed.

xxxWakeupHandler

Allows for reading from ready file descriptors after **select** processing.

Event Processing (Input Devices)

The following are provided callable routines for input extension device event processing:

AddInputCheck

Adds additional input checking for extension input devices that are not using the X server input ring.

RemoveInputCheck

Cancels the **AddInputCheck** subroutine.

Input Load Module (Input Devices)

The following subroutines must be supplied for input extension load modules:

ExtInitInput Initializes the input private structure for the extension input device.

deviceProc Initializes, turns on, opens, or closes the extension input device.

setDeviceMode

Sets the mode of a device.

setDeviceValuators

Initializes the valuators on a device.

getDeviceControl

Returns the current device control setting for the specified device.

changeDeviceControl

Changes the device control settings for the specified device.

processRawInputEvents

Formats a raw device event into an event that can be sent to an X client.

Related Information

Angebrannt, Susan, Drewry, Raymond, Karlton, Philip, Newman, Todd, Packard, Keith and Scheifler, Robert W. *Strategies for Porting the X v1 Sample Server*. Massachusetts Institute of Technology. 1991.

Fortune, Erik, and Israel, Elias. *The X-Window Server*. Digital Press.

Gettys, James, Newman, Ron and Scheifler, Robert W. *Xlib—C Language X Interface, MIT X Consortium Standard, X Version 11, Release 5*. MIT X Consortium 1991.

Patrick, Mark, and Sachs, George. *X11 Input Extension Library Specification. MIT X Consortium Standard. X Version 11, Release 5*. Hewlett-Packard Company, Ardent Computer, and the Massachusetts Institute of Technology. 1989, 1990, 1991.

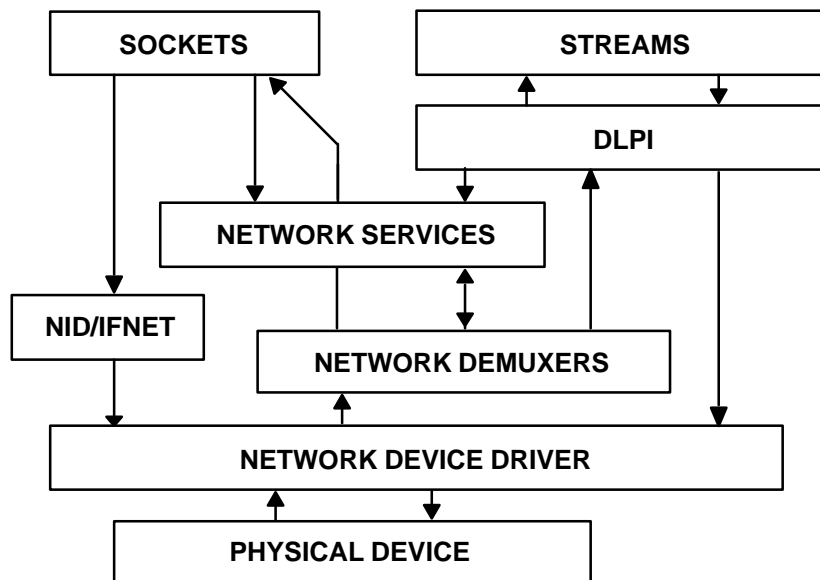
Patrick, Mark, and Sachs, George. *X11 Input Extension Protocol Specification. MIT X Consortium Standard. X Version 11, Release 5*. Hewlett-Packard Company, Ardent Computer, and the Massachusetts Institute of Technology. 1989, 1990, 1991.

Sachs, George. *X11 Input Extension Porting Document. MIT X Consortium Standard. X Version 11, Release 5*. Hewlett-Packard Company and the Massachusetts Institute of Technology. 1989, 1990, 1991.

Scheifler, Robert W. *X Window System Protocol, MIT X Consortium Standard, X Version 11, Release 5*. MIT X Consortium 1991.

Chapter 12. Implementing a Network Device Driver

A CDLI device driver has several components, as shown in the CDLI Device Driver Structure figure.



CDLI Device Driver Structure

Several of those components are system-supplied, but a device driver writer typically writes the following components:

Network Device Driver (NDD)

Defines a simple interface to network based devices that can be used by both the sockets network interface layer (IFNET) and the STREAMS DLPI data link layers. (See “Writing a Network Device Driver” on page 12-2.)

Network Demuxer (ND)

Provides common data-link receive functionality. The demuxer specifies receive filters that are used to distribute network packets. (See “Writing a Network Demuxer” on page 12-13.)

Network Interface Driver(NID)

The AIX Network Interface Driver (NID) is a layer of software between a network device driver (NDD) and an AIX network layer. This layer is required for all network device drivers that have to be made available to a network layer. (See “Writing a Network Interface Driver” on page 12-22.)

Writing a Network Device Driver

Network device drivers, including the system provided network device drivers (for example, Ethernet and Token-Ring), are implemented as loadable kernel extensions in AIX. The following general guidelines apply:

- By convention, device driver kernel extensions are installed in the **/usr/lib/drivers** directory. Some system utilities may assume that this is the case.
- When building the kernel extension the relevant base system exports must be imported. The system provided device drivers import the following: **kernexp.exp**, **syscalls.exp**, **streams.exp** (used by some drivers). The system export files are located in the **/usr/lib** directory.
- Device driver writers must decide how much of the kernel extension to pin. The major consideration is that code executed on the interrupt level should be pinned. In general, the system provided device drivers are pinned in their entirety.
- Network device drivers must be configured and loaded into the system.

Overview of Network Device Driver Changes in AIX Version 4.1

To accommodate the changes for the CDLI interface that the communications subsystem now uses, the network device driver's interfaces in AIX Version 4.1 are different from the interfaces in AIX Version 3.2. The parameters passed to the open, close, ioctl and output functions have been changed. The ioctl function has a new list of commands to service. The receive interrupt routine is now required to pass up frames via a required interface. There are also some minor changes to the way a network device driver is loaded and terminated.

Network Device Driver Initialization and Termination

Device driver initialization involves execution of the kernel extension's configuration entry point. This function is designated when the network device driver is built and is called by the system when the device driver is loaded. Usually, the **ifconfig** command is responsible for configuring the network device driver. The format of the entry should be as follows:

```
(*networkdd_config) (int cmd, struct uio *uio)
```

The entry has the following parameters:

- `cmd` indicates what type of configuration operation should be performed. The network device driver should recognize the following values for `cmd` (see **/usr/include/sys/device.h**) and can define additional device specific commands:

CFG_INIT	Initialize the device.
CFG_TERM	Terminate the device.
CFG_QVPD	Return vital product data.
CFG_UCODE	Download microcode.

- `uio` is a pointer to a **uio** structure whose data area contains a **ndd_config_t** structure (refer to **/usr/include/sys/ndd.h**). This **ndd_config_t** structure contains the configuration information for the network device driver

When called with the CFG_INIT command the device driver's configuration function should perform the following tasks:

- Do any pinning of modules or data structures required by the device driver.
- Do any lock initialization required by the driver.

- Initialize the device control structures, including allocating memory for the **ndd** structure.
- Call the **ns_attach** kernel service to add the device driver to the list of available network devices.

When called with the CFG_TERM command the device driver's configuration function should perform the following tasks:

- Do any unpinning of modules or data structures required by the device driver.
- Free any lock resources.
- Free any device control structures.
- Call the **ns_detach** kernel service to remove the device driver from the list of available network devices.

Warning: DLPI or the socket network interface layer may still have references to the device driver's **ndd** structure. This structure should only be deleted from the system when the device driver is certain that these references have been removed.

When called with the CFG_QVPD command the device driver's configuration function should return the devices vital product data in the **ndd_config_t** structure.

CFG_UCODE is device specific.

The configuration entry point can be called from the process environment only.

Sample network device driver code for configuration and unconfiguration follows:

```
xx_config(
    int          cmd,          /* command being processed */
    struct uio   *p_uio)      /* pointer to uio structure */
{
    xx_dev_ctl_t *p_dev_ctl = NULL; /* point to dev_ctl area */
    int rc = 0;                /* return code */
    int i;                     /* loop index */
    ndd_config_t ndd_config;    /* config information */

    /*
     * Use lockl operation to serialize the execution of the config commands.
     */
    lockl(&CFG_LOCK, LOCK_SHORT);
    if (!xx_initd) {
        /* perform first time initialization
         *
         * set all locks
         * init device driver structures
         */
    }
    pincode(xx_open);          /* pin the entire driver */

    uiomove((caddr_t) &ndd_config, sizeof(ndd_config_t), UIO_WRITE, p_uio);

    /*
     * find the device in the dev_list if it is there
     */
    p_dev_ctl = xx_dd_ctl.p_dev_list;
    while (p_dev_ctl) {
        if (p_dev_ctl->seq_number == ndd_config.seq_number)
            break;
        p_dev_ctl = p_dev_ctl->next;
    }
    switch(cmd) {
        case CFG_INIT:
```

```

if (p_dev_ctl) {
    rc = EBUSY;
    break;
}
/*
 * Allocate memory for the dev_ctl structure
 */
p_dev_ctl = xmalloc(sizeof(xx_dev_ctl_t), MEM_ALIGN,
    pinned_heap);

bzero(p_dev_ctl, sizeof(en3com_dev_ctl_t));

/*
 * Initialize the locks in the dev_ctl area
 */

.....
/*
 * Add the new dev_ctl into the dev_list
 */
p_dev_ctl->next = en3com_dd_ctl.p_dev_list;
xx_dd_ctl.p_dev_list = p_dev_ctl;
xx_dd_ctl.num_devs++;
/*
 * Copy in the dds for config manager
 */
if (copyin(ndd_config.dds, &p_dev_ctl->dds,
    sizeof(en3com_dds_t))) {
    rc = EIO;
    break;
}
p_dev_ctl->seq_number = ndd_config.seq_number;
/* save the dev_ctl address in the NDD correlator field */
p_dev_ctl->ndd.ndd_correlator = (caddr_t)p_dev_ctl;
p_dev_ctl->ndd.ndd_addrflen = XX_NADR_LENGTH;
p_dev_ctl->ndd.ndd_hdrflen = XX_HDR_LEN;
p_dev_ctl->ndd.ndd_physaddr = WRK.net_addr;
p_dev_ctl->ndd.ndd_mtu = XX_MAX_MTU;
p_dev_ctl->ndd.ndd_mintu = XX_MIN_MTU;
p_dev_ctl->ndd.ndd_type = NDD_ISO???.;
p_dev_ctl->ndd.ndd_flags = NDD_BROADCAST | NDD_SIMPLEX;
p_dev_ctl->ndd.ndd_open = xx_open;
p_dev_ctl->ndd.ndd_output = xx_output;
p_dev_ctl->ndd.ndd_ctl = xx_ctl;
p_dev_ctl->ndd.ndd_close = xx_close;
p_dev_ctl->ndd.ndd_specstats = (caddr_t)&(XXSTATS);
p_dev_ctl->ndd.ndd_specflen = sizeof(XXSTATS);

/* perform device-specific initialization */
.....
/* add the device to the NDD chain */
if (rc = ns_attach(&p_dev_ctl->ndd)) {
    return(rc);
}
break;
case CFG_TERM:
/* Does the device exist? */
if (!p_dev_ctl) {
    rc = ENODEV;
    break;
}

/*
 * Make sure the device is in CLOSED state.
 * Call ns_detach and make sure that it is done
 * without error.

```

```

        */
        if (p_dev_ctl->device_state != CLOSED || ns_detach(&(p_dev_ctl->ndd))) {
            rc = EBUSY;
            break;
        }
        /*
        * Remove the dev_ctl area from the dev_ctl list
        * and free the resources.
        */
        break;
    case CFG_QVPD:
        /* Does the device exist? */
        if (!p_dev_ctl) {
            rc = ENODEV;
            break;
        }
        if (copyout((caddr_t)&p_dev_ctl->vpd, ndd_config.p_vpd,
                    (int)ndd_config.l_vpd)) {
            rc = EIO;
        }
        break;
    default:
        rc = EINVAL;
    }
}
/* if we are about to be unloaded, free locks */
if (!xx_dd_ctl.num_devs) {
    /* free locks here */
    xx_initied = FALSE;
}
unpincode(xx_open);          /* unpin the entire driver */
unlockl(&CFG_LOCK);
return (rc);
}

```

CDLI – Device Driver Interface

In AIX Version 4.1, network device drivers should be entered only through the kernel CDLI users. STREAMS and sockets are implemented above CDLI, so that users gain access to the network device drivers through the standard socket and STREAMS application interface. The kernel establishes a **ndd** structure (refer to **/usr/include/sys/ndd.h**) for all network devices. This structure defines the entry points that the device driver must support. The following is a list of these entry points:

- **ndd_open**
- **ndd_close**
- **ndd_output**

ndd_open Entry Point

This entry point has the following format:

```
ndd_open(struct ndd *ndd)
```

The parameter **ndd** is a pointer to the system **ndd** structure for this network device. The file **/usr/include/sys/ndd.h** contains the definition of this structure.

Device driver users (DLPI or the socket network interface layer) gain access to the device through a call to the **ns_alloc** kernel service. This kernel service opens the network device, if required, by a call to the **ndd_open** entry point.

When the **ndd_open** function is called, the device driver should allocate the necessary system resources (such as DMA channel, interrupt level and priority). It should register its interrupt handler with the system using the **i_init** kernel service and initialize the device.

After the open is successful, the device driver is responsible for ORing the `ndd_flag` field with `NDD_UP`. The device driver should have ORed this field with `NDD_RUNNING` upon successful initialization.

ndd_open can be called from the process environment only.

A code fragment for a sample network device driver open routine follows:

```
xx_open(
    ndd_t          *p_ndd)          /* pointer to the ndd in the dev_ctl area */
{
    xx_dev_ctl_t  *p_dev_ctl = (xx_dev_ctl_t *) (p_ndd->ndd_correlator);
    int rc;
    /*
     * Set the device state and NDD flags
     */
    p_dev_ctl->device_state = OPEN_PENDING;
    p_ndd->ndd_flags = NDD_BROADCAST | NDD_SIMPLEX;
    /* set up locks, register interrupt handler */
    .....
    .....
    p_ndd->ndd_flags |= (NDD_RUNNING | NDD_UP);
    return(0);
}
```

ndd_close Entry Point

This entry point has the following format:

```
ndd_close(struct ndd *nnd)
```

The parameter `nnd` is a pointer to the system **nnd** structure for this network device.

DLPI or the socket network interface layer relinquishes access to a network device by calling the **ns_free** kernel service. This kernel service will close the device, when the user reference count reaches 0, by a call to the **ndd_close** entry point.

When the **ndd_close** function is called, the device driver should free its system resources (including all mbufs) and deregister its interrupt handler through a call to the **i_clear** kernel service. On entry to the close routine the device driver should remove the `NDD_UP` and `NDD_RUNNING` flags from the `ndd_flags` field.

ndd_close can be called from the process environment only.

A code fragment for an example Network Device Driver close routine follows:

```

xx_close(
    ndd_t      *p_ndd)    /* pointer to the ndd in the dev_ctl area */
{
    xx_dev_ctl_t  *p_dev_ctl = (xx_dev_ctl_t *) (p_ndd->ndd_correlator);
    int ipri;
    if (p_dev_ctl->device_state == OPENED) {
        p_dev_ctl->device_state = CLOSE_PENDING;
        /* wait for the transmit queue to drain */
        while (p_dev_ctl->device_state == CLOSE_PENDING &&
            (p_dev_ctl->tx_pending || p_dev_ctl->txq_len)) {
            DELAYMS(1000);          /* delay 1 second */
        }
    }

    /*
     * deactivate the adapter
     */
    p_dev_ctl->device_state = CLOSED;
    p_ndd->ndd_flags &= ~(NDD_RUNNING | NDD_UP | NDD_LIMBO | NDD_DEAD);

    unlock_enable(ipri, &SLIH_LOCK);

    /* cleanup all the resources allocated for open */
    .....
    unpincode(xx_open);
    return(0);
}

```

ndd_output Entry Point

This entry point has the following format:

```
ndd_output(struct ndd *ndd, struct mbuf *data)
```

The entry point has the following parameters:

- ndd is a pointer to the system **ndd** structure for this network device.
- data is a pointer to the mbuf chain to be transmitted.

DLPI or the socket network interface layer calls **ndd_output** *directly* to output data on the network device.

The first mbuf in each packet chain will be of the M_PKTHDR format (see **/usr/include/sys/mbuf.h**). Multiple mbufs may hold the packet and will be linked to data via the `m_next` field. Multiple packets on the transmit are supported and these will be linked to data via the `m_nextpkt` field. `m_pkthdr.len` is set equal to the total length of the individual packet.

On successful transmit, the device driver is responsible for freeing all mbufs. On failure, this is the callers responsibility. In the case of a multiple packet transmission, if any of the packets are transmitted the device driver should return success and free the mbufs.

On output, the device driver should check the value of `ndd->ndd_trace`. If this entry point is not null then the device driver should call the trace point for each packet which is transmitted (chained packets must be unchained before this call). The format of the call is:

```
(*ndd_trace) (struct ndd *ndd, struct mbuf *data, caddr_t *hp,
    ndd->ndd_trace_arg)
```

Parameters for **ndd_trace** have the following meanings:

- ndd is a pointer to the system ndd structure for this network device.
- data is a pointer to the mbuf chain to be transmitted.
- hp is a pointer to the start of the data in the mbuf being transmitted
- ndd_trace_arg is a cookie set by the trace routine.

The trace routine does not free the mbuf chain.

If transmission fails because of queue overruns, the device driver should return EAGAIN.

ndd_output may be called from the process or interrupt environment.

A code fragment for an example Network Device Driver output routine follows:

```

xx_output(
    ndd_t          *p_ndd,          /* pointer to the ndd in the dev_ctl area */
    struct mbuf    *p_mbuf)        /* pointer to a mbuf (chain) */
{
    xx_dev_ctl_t  *p_dev_ctl = (xx_dev_ctl_t *) (p_ndd->ndd_correlator);
    struct mbuf   *p_cur_mbuf;
    struct mbuf   *buf_tofree;
    int bus;
    int first;
    if (p_dev_ctl->device_state != OPENED) {
        return(ENETDOWN);
    }
    /*
     * if there is a transmit queue, put the packet onto the queue.
     */
    if (p_dev_ctl->txq_first) {
        /*
         * if the txq is full, return EAGAIN. Otherwise, queue as
         * many packets onto the transmit queue and free the
         * rest of the packets, return no error.
         */
        .....
        .....
    }
    while (p_mbuf) {
        p_cur_mbuf = p_mbuf;
        /*
         * If there is txd available, try to transmit the packet.
         */
        if (!(WRK.txd_avail->flags & XX_IN_USE)) {

            if (!xx_xmit(p_dev_ctl, p_cur_mbuf, bus)) {
                /*
                 * Transmit OK, free the packet
                 */
                p_mbuf = p_mbuf->m_nextpkt;
                p_cur_mbuf->m_nextpkt = NULL;;
                m_freem(p_cur_mbuf);
                first = FALSE;
            }
            else {
                /*
                 * Transmit error. Call hardware error recovery
                 * function. If this is the first packet,
                 * return error. Otherwise, free the reset packets
                 * and return error.
                 */
                p_ndd->ndd_genstats.ndd_oerrors++;
                if (first) {

                    return(ENETDOWN);
                }
            }
        }
    }
}

```

```

    }
    else {
        /* increment the error counter */
        while (p_cur_mbuf = p_mbuf) {
            p_mbuf = p_mbuf->m_nextpkt;
            p_cur_mbuf->m_nextpkt = NULL;
            m_freem(p_cur_mbuf);
        }
        return(0);
    }
}
/* if there is txd available */
else {
    while (p_cur_mbuf = buf_tofree) {
        p_ndd->ndd_genstats.ndd_xmitque_ovf++;
        p_ndd->ndd_genstats.ndd_opackets_drop++;
        buf_tofree = buf_tofree->m_nextpkt;
        p_cur_mbuf->m_nextpkt = NULL;
        m_freem(p_cur_mbuf);
    }
    return(0);
}
} /* while */
return(0);
}

xx_xmit(
    xx_dev_ctl_t *p_dev_ctl, /* pointer to the device control area */
    struct mbuf *p_mbuf,     /* pointer to the packet in mbuf */
    int bus)                /* handle for I/O bus accessing */
{
    ndd_t *p_ndd = &(p_dev_ctl->ndd);
    int count;
    int offset;
    int rc;
    int pio_rc = 0;
    /*
     * Call ndd_trace if it is enabled
     */
    if (p_ndd->ndd_trace) {
        (*(p_ndd->ndd_trace))(p_ndd, p_mbuf,
            p_mbuf->m_data, p_ndd->ndd_trace_arg);
    }
    /* increment the tx_pending count */
    p_dev_ctl->tx_pending++;
    /*
     * copy data into transmit buffer and do processor cache flush
     */
    m_copydata(p_mbuf, 0, count, p_txd->buf);
    /*
     * Pad short packet with garbage
     */
    if (count < XX_MIN_MTU)
        count = XX_MIN_MTU;
    p_txd->tx_len = count;
    .....
    .....
    /*
     * tell the adapter how many bytes to send
     * clear the status and set the EOP and EL bit.
     */
    .....
    /* start watchdog timer */
    w_start(&(TXWDT));
    return(0);
}

```

ndd_ctl Entry Point

This entry point for the IOCTL system call has the following format:

```
ndd_ctl(struct ndd *nnd, int cmd, caddr_t arg, int length)
```

This entry point has the following parameters:

- `nnd` is a pointer to the system `nnd` structure for this network device.
- `cmd` is the IOCTL.
- `arg` is the address of the ioctl arguments. This address will be in kernel memory.
- `length` is the length of `arg`.

The ioctls for a network device driver are shown below:

NDD_ADD_FILTER

Add receive data filter. `arg` points to the address of the filter function to add.

NDD_DEL_FILTER

Remove receive data filter. `arg` points to the address of the filter function to remove.

NDD_ADD_STATUS

Add status filter. `arg` points to the address of the filter function to add.

NDD_DEL_STATUS

Delete status filter. `arg` points to the address of the filter function to delete.

NDD_CLEAR_STATS

Clear all statistics maintained by the network device driver.

NDD_DISABLE_ADDRESS

Disable a multicast address. Remove the `NDD_ALTADDRS` flag in the `nnd` structure. `arg` contains the multicast address.

NDD_ENABLE_ADDRESS

Enable a multicast address. Set the `NDD_ALTADDRS` flag in the `nnd` structure. `arg` contains the multicast address.

NDD_GET_STATS

Get general statistics from the network device. General statistics are maintained by the device driver in the `nnd_genstats` field of the `nnd`. `arg` points to a user buffer where the `nnd_genstats` information should be placed.

NDD_GET_ALL_STATS

Get all statistics from the network device. All statistics means general statistics maintained by the device driver in the `nnd_genstats` field of the `nnd` and additional specific statistics maintained by the device driver in a structure pointed to by the `nnd_specstats` field. Some device drivers may only update specific statistics when IOCTL commands are processed. `arg` points to a device-specific structure that contains the `nnd_genstats` structure followed by device-specific information.

NDD_MIB_QUERY

Return the MIB structure identifying which options the device

supports. `arg` is a pointer to a structure of type `generic_mib_t` in the kernel address space.

NDD_MIB_GET

Get device specific MIBs. `arg` is a pointer to a structure of type `generic_mib_t` in the kernel address space.

NDD_DUMP_ADDR

Return the address of remote dump routine in `arg`.

NDD_MIB_ADDR

Get all receive addresses for the device. `arg` is a pointer to a structure of type `ndd_mib_addr_elem_t` in the kernel address space.

NDD_ENABLE_MULTICAST

Enable receipt of all multicast packets. Set `NDD_MULTICAST` flag in `ndd` structure.

NDD_DISABLE_MULTICAST

Disable receipt of all multicast packets. Remove `NDD_MULTICAST` flag in `ndd` structure.

NDD_PROMISCUOUS_ON

Set promiscuous mode on if the adapter allows this. Promiscuous mode allows the adapter to receive all the packets transmitted on the network. Set `NDD_PROMISC` flag in `ndd` structure.

NDD_PROMISCUOUS_OFF

Set promiscuous mode off. Remove `NDD_PROMISC` flag in `ndd` structure.

Network device drivers are not required to support all of the preceding IOCTLs. For IOCTLs not supported, the driver should return `EOPNOTSUPP`.

Additional device specific IOCTLs may also be supported.

DLPI or the socket network interface layer calls `ndd_ctl` *directly* if one of the preceding IOCTLs is issued on a stream or a raw (`AF_NDD`) socket.

`ndd_ctl` can be called from the process or interrupt environment.

Device Driver – CDLI Interface

This consists of the following kernel services and functions:

- `ns_attach` and `ns_detach`
- `nd_receive`
- `nd_status`

`ns_attach` and `ns_detach` Kernel Services

These have the following format:

```
ns_attach(struct ndd *ndd)
ns_detach(struct ndd *ndd)
```

The parameter `ndd` is a pointer to the system `ndd` structure for this network device.

The preceding two kernel services are called by network device driver configuration functions to add or remove their network device from the system's list of available devices.

These services are discussed, in detail, in the initialization and CDLI–Device Driver Interface sections.

ns_attach and **ns_detach** can be called from the process environment only.

nd_receive Function

Network device drivers pass receive data to the system through calls to the **nd_receive** function that is specified in the **ndd** structure for the network device. This function may be null. The format of the call is:

The format of the call is:

```
(*nd_receive(struct ndd *ndd, struct mbuf *data))
```

The parameters are as follows:

- **ndd** is a pointer to the system **ndd** structure for this network device.
- **data** is a pointer to the **mbuf** chain to be received.

AIX supports multiple protocols, concurrently, on the same network device for both sockets and STREAMS users. To accomplish this receive packets are passed to network demuxers. Network demuxers are loadable kernel extensions whose function is to route packets to the appropriate user. Generally, a unique demuxer will be required for each type of network device. This is because knowledge of the physical layer data headers is required for packet routing. Network device driver writers will need to provide a network demuxer with their device driver if one of the system-provided demuxers is not satisfactory for their network device. Because both device drivers and network demuxers are loadable kernel extensions it is possible to add a completely new network device to AIX Version 4.1 without requiring any modifications to the base operating system. The device driver can either bind in a network demuxer or have the system add a demuxer appropriate for the type of network device being added. When DLPI or the socket network interface layer calls **ns_alloc** to register use of a network device driver, the system checks if the **ndd_demuxer** field in the **ndd** structure is null. If the field is null (the device driver has not bound in a demuxer) then the system searches the known demuxers in the system trying to find a demuxer for this type of network device. The interface type is set in the **ndd_type** field of the **ndd** structure and describes classes of network devices which have the same physical layer characteristics such as 802.3 Ethernet and Token Ring. All devices of the same type can use a common demuxer. If a match is found, **ns_alloc** sets the **ndd** receive and status function pointers to the demuxers. It also sets the **ndd_demuxsource** field to 0 (system provided). AIX Version 4.1 provides demuxers for the following interface types: **NDD_ISO88023**, **NDD_ISO88025** and **NDD_FDDI** (see **/usr/include/sys/ndd.h**). Device driver writers for network devices of these types may elect to use the system demuxers instead of providing their own.

Typically, the **nd_receive** function is called by the device driver's receive interrupt routine. It is only called for receive frames that are not in error for any reason. The **mbuf** pointer passed to the **nd_receive** function may contain more than one **mbuf** chain.

nd_receive may be called with interrupts disabled.

nd_status Function

Network device drivers pass status event information to the system through calls to the **nd_status** function that is specified in the **ndd** structure. This function may be null.

The format of the call is:

```
(*nd_status(struct ndd *ndd, struct ndd_statblk *status))
```

The call has the following parameters:

- `ndd` is a pointer to the system **ndd** structure for this network device.
- `status` is a pointer to the status block. This structure is defined in **/usr/include/sys/ndd.h**

The **nd_status** function is typically called by the device driver's receive interrupt routine upon receipt of a bad frame. If the adapter has support for status interrupts, it should be called by the status interrupt routine. The status structure of type `ndd_statblk` should have the code field set to `NDD_BAD_PKTS` and values set in the option fields before **nd_status** is called. After calling the **nd_status** function, the interrupt routine should free the mbufs associated with the bad frame, or frames.

nd_status can be called from the process or interrupt environment.

Writing a Network Demuxer

Network demuxers, including the system provided ISO88023, ISO88025 and FDDI demuxers, are implemented as loadable kernel extensions in AIX Version 4.1. The following general guidelines apply:

- By convention, network demuxer kernel extensions are installed in the **/usr/lib/drivers** directory. Some system utilities may assume that this is the case.
- When building the kernel extension the relevant base system exports must be imported. The system provided demuxers import the following: **kernex.exp**, **syscalls.exp**. The system export files are located in the **/usr/lib** directory.
- Network demuxer writers must decide how much of the kernel extension to pin. The major consideration is that the demuxer's receive and status functions will be called with interrupts disabled. In general, the system provided demuxers are pinned in their entirety.
- Network demuxers must be configured and loaded into the system.

Network demuxers are loadable kernel extensions whose function is to route incoming packets to the appropriate user. Generally, a unique demuxer will be required for each type of network device. This is because knowledge of the physical layer data headers is required for packet routing. DLPI and the socket network interface layer make direct calls to network demuxers to register particular packet types to be received. Device drivers call network demuxers to route packets and status events to the correct users.

Demuxer Initialization

Network demuxer initialization involves execution of the kernel extension's configuration entry point. This configuration entry point is called by the configuration method of the network device driver during the system boot up. This function is invoked when the network demuxer is built and is called by the system when the kernel extension is loaded. The format of the entry should be as follows:

```
(*network_demuxer_config) (int cmd, struct uio *uio)
```

The entry has the following parameters:

<code>cmd</code>	Designates a particular demuxer command.
<code>uio</code>	Pointer to a uio structure. This is generally ignored by demuxers.

Command values that should be recognized by the network demuxer are:

<code>CFG_INIT</code>	Initialize the demuxer.
<code>CFG_TERM</code>	Terminate the demuxer.

Other demuxer specific commands can be defined.

When called with the `CFG_INIT` command the demuxer's configuration function should perform the following tasks:

- Do any pinning of modules or data structures required by the demuxer.
- Do any lock initialization required by the demuxer.
- Initialize the demuxer control structures.
- Call the `ns_add_demux` kernel service to add the demuxer to the system's list of available demuxers. This step is not required if the demuxer is to be bound to a device driver or is not to be made available to other network devices.

When called with the `CFG_TERM` command the demuxer's configuration function should perform the following tasks:

- Verify that there are no active users of the demuxer.
- Do any unpinning of modules or data structures required by the demuxer.
- Free any lock resources.
- Free any demuxer control structures.
- If required, call the `ns_del_demux` kernel service to remove the demuxer from the system's list of available demuxers.

nd_add_filter Function

The `nd_add_filter` function adds a receive filter for the routing of received data. It is invoked by the CDLI routines and has the following entry point:

```
(*nd_add_filter)(struct ndd *nddp, caddr_t filter, int len, ns_user_t ns_user)
```

The entry point has the following parameters:

<code>nddp</code>	Pointer to the system ndd structure for this demuxer.
<code>filter</code>	Address of a user defined structure which contains information that the demuxer needs such as filter type. See <code>/usr/include/sys/cdli.h</code> for the ns_8022 structure as an example of what is used for this parameter with 802.2 networks.
<code>len</code>	Length in bytes of the filter parameter.
<code>ns_user</code>	Pointer to the ns_user structure that describes the user of the filter. The structure is defined in <code>/usr/include/sys/cdli.h</code> .

The `nd_add_filter` field of the **ns_demuxer** structure filter is set to the address of this function during configuration of the demuxer. The function should perform the following tasks:

- Set any necessary data structure locks
- Perform sanity checks on the fields of the **ns_user** structure
- Handle the appropriate filter type specified in the filter parameter, if applicable.
- Set the `nddp->ndd_specdemux` field to a pointer to a private demuxer control structure that maintains a list of which filters have been added. Be sure to save **ns_user** information here; it will later be retrieved by the **nd_receive** function.
- Call the `nddp->ndd_ctl` entry point of the Network Device Driver to register the filter with the NDD. The calling format is:

```
(( *nddp->ndd_ctl )) (nddp, NDD_ADD_FILTER, filter, len);
```

- If any errors are returned from the preceding call, invoke the appropriate demuxer delete filter function.
- Release all locks set upon entry into the function.
- Return 0 for success, **errno** otherwise.

nd_del_filter Function

The **nd_del_filter** function for the demuxer deletes a previously specified receive filter. It is invoked by the CDLI routines and has the following entry point:

```
(*nd_del_filter)(struct ndd *nddp, caddr_t filter, int len)
```

The entry point has the following parameters:

<code>nddp</code>	Pointer to the system ndd structure for this demuxer.
<code>filter</code>	Address of a user defined structure which contains information about the filter that is to be deleted.
<code>len</code>	Length in bytes of the filter parameter.

The `nd_del_filter` field of the **ns_demuxer** structure filter is set to the address of this function during configuration of the demuxer. The function should perform the following tasks:

- Set any necessary data structure locks.
- Check the `len` parameter for correctness and verify that a filter has previously been added.
- If applicable, handle the appropriate filter type specified in the filter parameter.
- Retrieve the filter list from `nddp->ndd_specdemux`, search for the filter to be deleted and free any storage associated with that filter.
- Call the **nddp->ndd_ctl** entry point of the Network Device Driver to inform it that the filter is being deleted. On return from this entry point, check for nonzero error return codes. The calling format is:

```
(( *nddp->ndd_ctl )) (nddp, NDD_DEL_FILTER, filter, len);
```

- Release all locks set upon entry into the function.
- Return 0 for success, **errno** otherwise.

nd_add_status Function

The **nd_add_status** function adds a filter for routing asynchronous status. It is invoked by the CDLI routines and has the following entry point:

```
(*nd_add_status)(struct ndd *nddp, caddr_t filter, int len, ns_stater *ns_stater)
```

The entry point has the following parameters:

<code>nddp</code>	Pointer to the system ndd structure for this demuxer
<code>filter</code>	Address of a user defined structure which contains information that the demuxer needs such as filter type. See <code>/usr/include/sys/cdli.h</code> for the ns_com_status structure as an example of what is typically used for this parameter.
<code>len</code>	Length in bytes of the filter parameter.

`ns_statuser` Pointer to structure describing the status user. The structure is defined in `/usr/include/sys/cdli.h`.

The `nd_add_status` field of the `ns_demuxer` structure filter is set to the address of this function during configuration of the demuxer. The function should perform the following tasks:

- Set any necessary data structure locks.
- Check the `len` parameter for correctness and verify that a filter has previously been added.
- Call the `dmx_add_status` kernel service to add the status filter. The calling format is:

```
dmx_add_status(nddp, filter, ns_statuser)
```

- If no errors were encountered, call the `nddp->ndd_ctl` entry point of the Network Device Driver to register the status filter. On return from this entry point, check for nonzero error return codes. The calling format is:

```
(( *nddp->ndd_ctl )) (nddp, NDD_ADD_STATUS, filter, len);
```

- Release all locks set upon entry into the function.
- Return 0 for success, `errno` otherwise.

nd_del_status Function

The `nd_del_status` function deletes a previously added status filter. It is invoked by the CDLI routines and has the following entry point:

```
(*nd_del_status)(struct ndd *nddp, caddr_t filter, int len)
```

The parameters are the same as the parameters of the companion `nd_add_status` function.

The `nd_del_status` field of the `ns_demuxer` structure filter is set to the address of this function during configuration of the demuxer. The function should perform the following tasks:

- Set any necessary data structure locks.
- Check the `len` parameter for correctness and verify that a filter has previously been added.

Call the `dmx_del_status` kernel service to delete the status filter. The calling format is:

```
dmx_del_status(nddp, filter)
```

- If no errors were encountered, call the `nddp->ndd_ctl` entry point of the Network Device Driver to unregister the status filter. On return from this entry point, check for nonzero error return codes. The calling format is:

```
(( *nddp->ndd_ctl )) (nddp, NDD_DEL_STATUS, filter, len);
```

- Release all locks set upon entry into the function.
- Return 0 for success, `errno` otherwise.

nd_receive Function

The `nd_receive` function is invoked by the Network Device Driver for receive packets.

The `nd_receive` field of the `ns_demuxer` structure filter is set to the address of this function during configuration of the demuxer. The function should perform the following tasks:

- Set any necessary data structure locks.
- If there are no receive filters, free all the mbufs and return.
- Grab each mbuf in the chain and parse the header data to determine which type of additional processing needs to be done. For example, an Ethernet driver will send the **nd_receive** function 802.3 packets, or standard Ethernet type packets, or both of these types of packets. Use the macros `DELIVER_PACKET` or `IFSTUFF_AND_DELIVER` found in `/usr/include/net/nd_lan.h` to send the data to the next layer. For 802.2 style networks, the `dmx_8022_receive` kernel service has been provided to send the data to the next layer. The entry point has the following format:

```
dmx_8022_receive(struct ndd *nddp, struct mbuf *m, int len)
```

For this entry point, the `len` parameter is the size of the MAC header in bytes.

- After exhausting all the mbufs in the chain, release all locks and return. The return value for this function is void.

For more information, see the description of **nd_receive**, on page 12-12, in “Writing a Network Device Driver”.

nd_status Function

The **nd_status** function is invoked by the Network Device Driver for all status conditions. It distributes asynchronous status to the appropriate network services users.

The `nd_status` field of the **ns_demuxer** structure filter is set to the address of this function during configuration of the demuxer. The function should perform the following tasks:

- Set any necessary data structure locks.
- Pass status information to the next layer by calling the **dmx_status** kernel service. The calling format is:

```
dmx_status(struct ndd *nddp, struct ndd_statblk *status)
```

- Release all locks set upon entry into the function. The return value of this function is void.

For more information, see the description of **nd_status**, on page 12-12, in “Writing a Network Device Driver”.

nd_response Function

The **nd_response** function is invoked by the higher level protocols if they choose to perform 802.2 LLC Exchange Station ID or TEST Link Frame processing. It has the following entry point:

```
(*nd_response)(struct ndd *nddp, struct mbuf *m, int llcoffset)
```

The entry point has the following parameters:

<code>ndd</code>	Pointer to the system ndd structure for this demuxer.
<code>m</code>	Pointer to a mbuf structure which may not contain more than one packet which contains a XID or TEST packet.
<code>llcoffset</code>	Byte offset to the start of the <code>llc</code> header in the mbuf.

The `nd_response` field of the **ns_demuxer** structure filter is set to the address of this function during configuration of the demuxer. The function should perform the following tasks:

- Set any necessary data structure locks.

- Copy the MAC source address for the MAC destination address and copy the `nndp->nnd_physaddr` to the MAC source address.
- If the MAC destination address indicates routing information is present, turn off the routing control bits in the MAC header.
- Call the device driver output routine (`*nndp->nnd_output`)(`nndp,m`) to send the packet to the driver. If a nonzero error is returned, free the mbuf.
- Release all locks set upon entry into the function. The return value of this function is void.

DLPI/Socket – Network Demuxer Interface

When a network demuxer is bound to a network device driver either by the device driver or through a `ns_alloc` call (`ns_alloc` is described in the appendix of this book), a linkage is created between the `nnd` structure and the demuxer's `ns_demuxer` structure (defined in `/usr/include/sys/cdli.h`). This is done by setting the `nnd_demuxer` field in the `nnd` structure equal to the address of the demuxer's structure. When DLPI or the socket network interface issues a request to the system (using `ns_add_filter` or `ns_del_filter`) to start receiving or stop receiving packets or status on this network interface, the system invokes the demuxer's entry points for adding or deleting filters or status as defined in the `ns_demuxer` structure. A sample fragment of code for `ns_add_filter` is:

```

/*****
 *
 *   ns_add_filter() - Pass "add filter" request on to demuxer
 *
 *****/
ns_add_filter(nndp, filter, len, ns_user)
    struct nnd      *nndp;      /* specific interface */
    caddr_t         filter;
    int             len;
    struct ns_user  *ns_user; /* the details */
{
    return((*nndp->nnd_demuxer->nd_add_filter)
           (nndp, filter, len, ns_user));
}

```

Thus these kernel service are really entry points into the network demuxer. The syntax of the calls are as follows:

```

ns_add_filter(struct nnd *nnd, caddr_t filter, int len,
             struct ns_user *ns_user);

```

The following parameters apply to `ns_add_filter`:

<code>nnd</code>	Pointer to the system <code>nnd</code> structure for this network device. Typically, a demuxer user obtains this pointer through a prior call to <code>ns_alloc</code> .
<code>filter</code>	Address of the filter function to be added. This is demuxer dependent.
<code>len</code>	Length of the filter.
<code>ns_user</code>	Pointer to a <code>ns_user</code> structure. This structure is defined in <code>/usr/include/sys/cdli.h</code> .

“Sample Code – DLPI Call to `ns_add_filter`”, on page 12-21, shows how DLPI calls `ns_add_filter`. First, a `ns_alloc` call is made to obtain the `nnd` pointer. DLPI builds the interface name from the basename of the device being opened (for example, “`tr`” for `/dev/dlpi/tr`) and the physical point of attachment (for example, 0) passed by the user in the `DL_ATTACH_REQ`. Next the sample illustrates how DLPI handles a `DL_SUBS_BIND_REQ` request to receive packets from this interface that match a specific 802.2 logical link control

(LLC) header. Note that DLPI simply passes the users request intact to the demuxer. Also note how the **ns_user** structure is built. The `isr` field points to the DLPI interrupt handler and `isr_data` is set to the internal address of this Stream. When the demuxer calls DLPI on a receive packet it also returns `isr_data` and DLPI does not need to do **any** packet demultiplexing to locate the target user. The `net_isr` field is NULL indicating that the DLPI interrupt routine is to be called on the interrupt level. `pkt_format` tells the demuxer how much of the LLC header to remove on input (and how much to expect on output). DLPI simply passes the format which the user has set through an IOCTL.

ns_add_status has the following syntax:

```
ns_add_status(struct ndd *ndd, caddr_t filter, int len, struct
ns_stater *ns_stater);
```

The following parameters apply to **ns_add_status**:

`ndd` Pointer to the system ndd structure for this network device. Typically, a demuxer user obtains this pointer through a prior call to **ns_alloc**.

`filter` Address of the status function to be added. This is demuxer dependent.

`len` Length of the filter.

`ns_stater` Pointer to a **ns_stater** structure. This structure is defined in **/usr/include/sys/cdli.h**.

When **ns_add_status** is called, the demuxer should issue a **ndd_ctl** with the NDD_ADD_STATUS operation.

The call to reverse **ns_add_filter** is **ns_del_filter**. The call to reverse **ns_add_status** is **ns_del_status**.

ns_del_filter has the following syntax:

```
ns_del_filter(struct ndd *ndd, caddr_t filter, int len)
```

The function arguments for **ns_del_filter** are the same as those for **ns_add_filter**.

ns_del_status has the following syntax:

```
ns_del_status(struct ndd *ndd, caddr_t filter, int len)
```

The function arguments for **ns_del_status** are the same as those for **ns_del_filter**.

The demuxer should issue a **ndd_ctl** with the NDD_DEL_STATUS operation when the **ns_del_status** entry point is called.

Many of the demuxer functions for 802.2 LANs can be handled by common routines. A demuxer can register to use these functions by setting `nd_use_nsdmx` in the demuxers **ns_demuxer** structure to the true value. Setting `nd_use_nsdmx` true causes the **dmx_init** function to be called when **ns_alloc** binds a network demuxer to a system **ndd** structure.

Device Driver – Network Demuxer Interface

The Network Device Driver determines which demuxer to use by setting the `ndd_type` field during configuration in its config entry point routine. (For more information see “Writing a Network Device Driver”, on page 12-2.) A list of defined `ndd_types` can be found in `/usr/include/sys/ndd.h`; a typical name is `NDD_ISO88023`. The command `ifconfig` calls the configuration method of the NDD, which calls the configuration method for the demuxer. In addition, `ifconfig` calls the configuration method for the associated NID for the interface. (For more information, see “Loading and Initialization” in “Network Interface Driver Functions”, on page 12-22.)

The demuxer’s configuration entry point passes to the `ns_add_demux` function a structure of type `ns_demuxer` that contains fields that are set to the address of functions that the CDLI and Network Device Driver routines eventually call. For the NDD, this means that the `nd_receive` and `nd_status` fields of the `ndd` structure are used. These fields contain function pointers to the appropriate demuxer functions. (For more information, see the discussion of the `nd_receive` function, on page 12-16, and the discussion of the `nd_status` function, on page 12-17.)

Sample Code – DLPI Call to ns_add_filter

```
/*
 * dlb_attach - attach interface to this module
 *
 */
staticf MBLKP
dlb_attach(dlb, mp)
    DLBP dlb;
    MBLKP mp;
{
    int len;
    int ppa;
    char *name;
    NDDP ndd;
    ...
    ppa = ((dl_attach_req_t *) (mp->b_rptr))->dl_ppa;
    name = mknddname(dlb->dlb_nddname, ppa);
    if (ns_alloc(name, &ndd)) {
    ...
    }
    ...
}
/*
 * dlb_subs_bind - bind a SNAP
 *
 * if dlb_isap != 0xAA, then either
 * - we bound to some non-snap sap
 * - we have already done a subs_bind
 */staticf MBLKP
dlb_subs_bind(dlb, mp)
    DLBP dlb;
    MBLKP mp;
{
    dl_subs_bind_req_t *dlsbr = (dl_subs_bind_req_t *)mp->b_rptr;
    snap_t *snap = (snap_t *) (mp->b_rptr + dlsbr->dl_subs_sap_offset);
    int err = 0;
    struct      ns_8022      dl;
    struct      ns_user      ns_user;
    ...
    else if (dlsbr->dl_subs_bind_class != DL_HIERARCHICAL_BIND)
        err = DL_NOTSUPPORTED;
    else if (dlsbr->dl_subs_sap_length != sizeof(snap_t))
        err = DL_BADADDR;
    else if ((char *)snap + sizeof(snap_t) > mp->b_wptr)
        err = DL_BADADDR;
    ...
    ns_user.isr = (int)dlb_intr;
    ns_user.isr_data = (caddr_t)dlb;
    ns_user.protoq = nilp(struct ifqueue);
    ns_user.netisr = NULL;
    ns_user.ifp = nilp(struct ifnet);
    ns_user.pkt_format = dlb->dlb_pkt_format;
    dl.filtertype = NS_8022_LLC_DSAP_SNAP;
    dl.dsap = dlb->dlb_isap;
    bcopy(snap, dl.orgcode, sizeof(snap_t));
    if (err = ns_add_filter(dlb->dlb_ndd, &dl, sizeof(dl), &ns_user))
        return dlb_error(mp, DL_SUBS_BIND_REQ, err, dlb);
    ...
}
```

Writing a Network Interface Driver

The AIX Network Interface Driver (NID) is a layer of software between a network device driver (NDD) and an AIX network layer. This layer is required for all network device drivers that have to be made available to a network layer. This discussion concentrates on NIDs for an Internet Protocol (IP) network layer, though you can easily modify an NID to support other network layers.

An AIX Network Interface Driver provides a uniform interface to the IP layer. The NID passes output packets from the IP layer to the Network Device Driver (NDD). The device driver insulates the NID from the hardware but does not hide the special anomaly from the type of underlying physical network.

Each type of physical network has unique access requirement. For Ethernet, this is the Ethernet header. For a Token Ring, it is the MAC/LLC header. These different anomalies of the network are handled inside the NID, providing a uniform interface to the network protocols.

Basic Functions of a Network Interface Driver

The NID is a dynamically loadable kernel extension similar to a device driver. You must load or add the NID to the kernel through a configuration method. A typical NID performs the following basic functions:

- Provides a uniform interface from the network layer to the network device driver.
- Translates an IP address to a hardware address for the underlying device driver.
- Builds the communication device driver specific protocol header (see the Data Packet for Ethernet figure, on page 12-27).
- Communicates with the network device driver.

A specific NID may perform more functions than those previously listed.

Summary of NID Changes in AIX Version 4.1

Several functions which were in AIX Version 3.2 NIDs have been removed from the NIDs and placed into the Network Demuxer. For example, the Receive Data and Status Interrupt function has been moved to the Network Demuxer. (For more information on this, see "Writing a Network Demuxer", on page 12-13.) The **add_arp_iftype**, **del_arp_iftype**, and **find_arp_iftype** kernel services are no longer available, making it the responsibility of NID writers to either supply their own ARP and address resolution routines or to use the one supplied by the system. There is an all new set of NID specific IOCTL calls. Finally, there are some small changes regarding initialization and termination of the NID.

Network Interface Driver Functions

A network interface driver (NID) performs the following functions:

- Loading and initializing
- Communicating with the NDD
- Translating network addresses to hardware addresses
- Handling NID specific IOCTL calls
- Terminating and unloading.

Loading and Initialization

The configuration method, **ifconfig**, loads the AIX NID kernel extension.

The **ifconfig** command configures the correct NID with the correct Network Device Driver (NDD). The non-numeric portion of the interface parameter from the command line is used as the NDD name. The NID name is the same as the NDD name with `if_` prepended to the NDD name. For example, the following command tells **ifconfig** to configure NID `if_en` and NDD `en` in `/usr/lib/drivers`:

```
ifconfig en0 1.1.1.1 up
```

The **ifconfig** command is run automatically when the system is started, and configures all NIDs that have been defined in the ODM database. For your NID to be automatically configured, your interface must be defined in the CuDv ODM object class with the parent attribute set to `inet0` and the name attribute set to the name of your interface. In addition, the CuAt ODM object class must have a corresponding `netaddr` attribute defined for the interface.

If you choose not to have **ifconfig** configure your NID automatically, you can manually issue the command as shown above, or place the command in one of the shell scripts that are run when the system is started.

After loading the NID, the configuration method calls the NID kernel extension entry point with the **CFG_INIT** command. This initializes the AIX NID.

Steps during initialization include:

- Pin the code and data for the NID kernel extension if not already pinned. The code might be pinned already if the configuration method for this NID has been invoked earlier for another adapter of the same type. If there are multiple adapters of the same type, the same NID code services all these adapters. It is not necessary to load multiple copies of an NID for multiple adapters of the same type.
- Open and initialize the underlying NDD. This is done by the **ns_alloc** subroutine in the following sample code.
- Initialize an **ifnet** structure and call the **if_attach** kernel service. This is done by **xx_attach** in the following sample code. The **if_attach** kernel service adds a NID to the network interface list, which is a linked list of **ifnet** structures.

Initialize the **ifnet** structure with the address of the **output** routine (`ifp->if_output`), **ioctl** routine (`ifp->if_ioctl`), and **reset** routine (`ifp->if_reset`).

The following sample code illustrates loading and initializing:

```
struct xx_softc {
    struct arpcom xx_ac; /* common part */
    struct ndd *nndp;
};

config_xx(cmd, uio)
int cmd;
struct uio *uio;
{
    struct device_req device;
    int error = 0;
    int unit;
    struct xx_softc *xxp;
    char *cp;
    int type;
    if ( (cmd != CFG_INIT) || (uio == NULL) )
        return(EINVAL);
    if (uiomove((caddr_t) &device, (int)sizeof(device), UIO_WRITE, uio))
        return(EFAULT);
```

```

lockl(&if_xx_lock, LOCK_SHORT);
if (init == 0) {
    if (ifsize <= 0)
        ifsize = IF_SIZE;
    if (error = pincodex(config_xx))
        goto out;
    xx_softc = (struct xx_softc *)
        xmalloc(sizeof(struct xx_softc)*ifsize, 2, pinned_heap);
    if ( (xx_softc == (struct xx_softc *)NULL) ||
        (xx_softc == (struct xx_softc *)NULL) ) {
        unpincodex(config_xx);
        unlockl(&if_xx_lock);
        return(ENOMEM);
    }
    bzero(xx_softc, sizeof(struct xx_softc) * ifsize);
    init++;
}
cp = device.dev_name;
while(*cp < '0' || *cp > '9') cp++;
unit = atoi(cp);
if (unit >= ifsize) {
    error = ENXIO;
    goto out;
}
if (!strcmp(device.dev_name, "xx", 2)) {
    type = IFT_XX;
    xxp = &xx_softc[unit];
} else {
    error = EINVAL;
    goto out;
}
error = ns_alloc(device.dev_name, &xxp->nddp);
if (error == 0)
    xx_attach(unit, type);
else
    bsdlog(LOG_ERR,
        "if_xx: ns_alloc(%s) failed with errno = %d\n",
        device.dev_name, error);
out:
unlockl(&if_xx_lock);
return(error);
}
xx_attach(unit, type)
unsigned    unit;
unsigned    type;
{
    register struct ifnet    *ifp;
    register struct xx_softc    *xxp;
    extern int    xx_output();
    extern int    xx_ioctl();
    xxp    = &xx_softc[unit];
    ifp    = &xxp->xx_if;
    ifp->if_name    = "xx";
    ifp->if_mtu    = ??;
    ifp->if_type    = type;
    ifp->if_unit    = unit;
    bcopy(xxp->nddp->ndd_physaddr, xxp->xx_addr, 6);
    ifp->if_flags    = IFF_BROADCAST | IFF_NOTRAILERS;
    /* Check if the adapter supports local echo */

```

```

if (xyp->nndp->nnd_flags & NDD_SIMPLEX)
    ifp->if_flags |= IFF_SIMPLEX;
ifp->if_output = xx_output;
ifp->if_ioctl = xx_ioctl;
ifp->if_addrln = ??;
ifp->if_hdrln = ??;
ifp->if_mtu = ??;
ifp->if_unit = unit;
ifp->if_name = "xx";
ifp->if_init = xx_init;
ifp->if_output = xx_output;
ifp->if_ioctl = xx_ioctl;
ifp->if_type = IFT_XX;
ifp->if_addrln = ??;
ifp->if_hdrln = ??;
/* packet filter support */
ifp->if_flags |= IFF_BPF; /* Enable bpf support */
ifp->if_tap = NULL; /* Inactive tap filter */
ifp->if_arpres = arpresolve;
ifp->if_arpnev = revarpinput;
ifp->if_arpinput = arpinput;
if_attach(ifp);
}

```

Communicating with the IP

The TCP/IP and NIDs are different kernel extensions that are loaded separately. In order to communicate with each other, these kernel extensions use the functions and data structures of the base kernel extension. These kernel services include:

- Address Family Domain kernel services
- Network Interface Device kernel services
- Routing and Interface Address kernel services.

Some of these services are:

add_input_type

Adds an interface type to the Network Input Table.

del_input_type Deletes an input type from the Network Input Table.

find_input_type

Finds an input type from the Network Input Table.

if_attach Adds a network interface to the network interface list.

if_detach Deletes a network interface from the network interface list.

ifunit Returns the **ifnet** structure for the requested interface.

ifa_ifwithaddr Locates an interface based on a complete interface address.

ifa_ifwithstaddr

Locates the point-to-point interface with a given destination address.

ifa_ifwithnet Locates an interface on a specific network.

if_down Marks an interface as down.

if_nostat Changes statistical elements of the interface array to zero in preparation for an attach operation.

Outgoing Packets

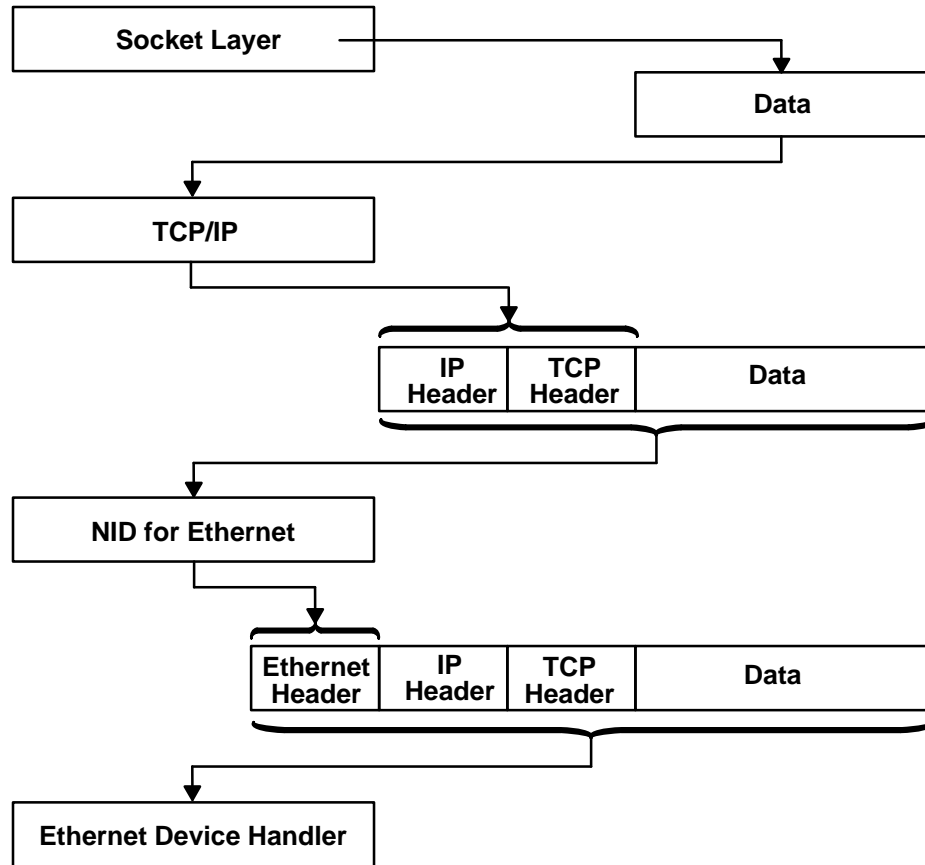
For the outgoing packets, the routing code in the IP layer determines the correct NID to use by scanning the linked list of **ifnet** structures. The routing code uses Interface Address kernel services like **ifa_ifwithaddr**, **ifa_ifwithstaddr**, and **ifa_ifwithnet** to locate the appropriate **ifnet** structure. For more detailed information on kernel services, see the *AIX Version 4.1 Technical Reference, Volume 5: Kernel and Subsystems*

The **ifnet** structure is initialized by each NID during its initialization phase. Once the appropriate NID is located, the IP layer calls the **output** routine of the NID:

```
(*ifp->if_output)(ifp, m, dst, rt)
    struct ifnet *ifp;
    struct mbuf *m;
    struct sockaddr *dst;
    struct rtable *rt;
```

- ifp** A pointer to the **ifnet** structure. This is required as there may be multiple adapters of the same type that are serviced by the same NID. The **if_unit** field in the **ifnet** structure is used to determine the appropriate adapter.
- m** The mbuf chain containing the data packet. This mbuf chain is freed by the NID or the device handler.
- dst** A pointer to the socket address of the destination of the packet.
- rt** A pointer to a routing table entry for this packet. A null value indicates there is no routing table entry.

The following Data Packet for Ethernet figure illustrates the modifications to an outgoing packet from the socket layer to the Ethernet device handler.



Data Packet for Ethernet

The **output** routine provides no guarantee for the transmission of packets. There is no acknowledgment of a successful delivery. The errors returned are those that can be detected immediately, such as the interface is down, no memory buffers, and address family not supported. If the error is detected after the call is returned, the protocol is *not* notified.

The following sample code shows typical code for an **output** routine:

```

/*****
 *      xx_output() - output packet to network
 *****/
xx_output(ifp, m, dst, rt)
register struct ifnet *ifp;
register struct mbuf *m;
struct sockaddr_xx *dst;
struct rtentry *rt;
{
    register struct xx_softc      *xxp;
    register struct xx_hdr        *hdrp;
    struct xx_hdr                 hdr;
    register int                   hdr_len;
    register int                   error = 0;
    struct mbuf                    *mcopy = 0;

    if ((ifp->if_flags & (IFF_UP|IFF_RUNNING)) != (IFF_UP|IFF_RUNNING)) {
        error = ENETDOWN;
        ++xxp->if_snd.ifq_drops;
    }
}

```

```

        goto out;
    }

    if ((xyp->if_flags & IFF_SIMPLEX) && (m->m_flags & M_BCAST))
        mcopy = m_copy(m, 0, (int)M_COPYALL, M_DONTWAIT);

    xyp = &xx_softc[xyp->if_unit];

    switch (dst->sa_family) {
        /*
         * call necessary ARP resolve routine or any other resolve
         * here
         */
        hdrp = (struct xx_hdr *)&hdr;
        break;
    }

    /*
     * Add local net header.  If no space in first mbuf,
     * allocate another.
     */
    hdr_len = xx_len;
    M_PREPEND(m, hdr_len, M_DONTWAIT);
    if (!m) {
        error = ENOBUFS;
        ++xyp->if_snd.ifq_drops;
        goto out;
    }

    /* copy in the local net headers */
    bcopy((caddr_t)hdrp, mtod(m, caddr_t), hdr_len);

    m = m_collapse(m, xx_MAX_GATHERS);
    if (!m) {
        error = ENOBUFS;
        ++xyp->if_snd.ifq_drops;
        goto out;
    }

    if (m->m_flags & M_BCAST|M_MCAST)
        ++xyp->if_omcasts;

    xyp->if_obytes += m->m_pkthdr.len;
    if (!xyp->nndp) {
        m_freem(m);
        ++xyp->if_oerrors;
    }
    /*
     * call the Network Device Driver's output function
     */
    } else if ((*xyp->nndp->nnd_output)(xyp->nndp, m)) {
        m_freem(m);
        ++xyp->if_oerrors;
    }
    if (mcopy)
        (void) looutput(xyp, mcopy, dst, rt);
    ++xyp->if_opackets;
    m = 0;
out:
    if (m)

```

```

        m_freem(m);
    return (error);
}

```

Communicating with the Device Handler

The communications device handler interface kernel services provide a standard interface between NIDs and AIX communication device handlers.

The **ns_alloc** and **ns_free** services allocate and relinquish use of a Network Device Driver, respectively. Once the NDD has been allocated, outgoing packets to the NDD can be sent by calling the driver's output function found in the `ndd_output` field of the **ndd** structure. This structure is returned after a successful call to **ns_alloc**.

Note: **ns_alloc** and **ns_free** are described in the appendix of this book.

Output Data

The **output** routine for NID is called by the network layer (IP). The NID transmits the data by calling the function found in `ndd_output` of the **ndd** structure. This is shown in the **xx_output** routine in the sample code starting on page 12-27.

Before calling the **ndd_output** function, the output routine might have to build the corresponding link-level header. If the **output** routine is called by IP, it is supplied with a destination address in a **sockaddr** structure. If the `sockaddr` address family is supported by the NID, the NID has to map this `sockaddr` into a link-level address.

This mapping may involve a lookup, or it may require more involved techniques like an Address Resolution Protocol (ARP). For more information on ARPs, see the "Translating Network Addresses to Hardware Addresses" section on page 12-29.

It may not be always necessary to build the link-level header. For example, a point-to-point link may not need a link-level header.

The **ndd_output** function pointer requires two parameters. The first is a pointer to the **ndd** structure of which `ndd_output` is a field. The second is a pointer to the `mbuf` with contains the data to transmit. The first `mbuf` in each packet chain will be of the `M_PKTHDR` format (see `/usr/include/sys/mbuf.h`). Multiple `mbufs` may hold the packet and will be linked to data the `m_next` field. **ndd_output** points to a device driver entry point similar to the following:

```

xx_output(p_ndd, p_mbuf)
nnd_t      *p_ndd; /* pointer to the ndd in the dev_ctl area */
struct mbuf *p_mbuf; /* pointer to a mbuf (chain) */

```

This function returns a zero after successful transmission. If a nonzero return is encountered, it is the NIDs responsibility to free the `mbufs`.

Translating Network Addresses to Hardware Addresses

The network layer provides the NID with the destination network address. If it is not a point-to-point network the NID must translate this network address into a destination hardware address to perform a successful transmission of the packet.

For the existing AIX NIDs, different mechanisms are used for different types of NIDs. For Ethernet, Token Ring, 802.3, and Fiber Distributed Data Interface (FDDI), the Address Resolution Protocol (ARP) is used.

The ARP is defined in the RFC826. The ARP provides a dynamic address-translation mechanism for networks that support broadcast or multicast communication.

The idea of ARP is simple. Whenever a packet needs to be sent out, the NID calls the ARP resolver routine to get the hardware address of the destination. If the address is already

known (in the cache translation table) then that value is returned. If not, the packet is queued and an ARP request is broadcast on the network. The request has the network address of the required destination host.

When the correct destination host receives the ARP request, it sends back an ARP reply providing the requester with the hardware destination address. The original sender host can then update its cache translation table and can transmit the queued-up output packet.

The resolver routines for the previously mentioned networks (Ethernet, Token Ring, 802.3, and FDDI) for the Internet address family are provided by the corresponding NID kernel extension.

The resolver routine is specified by the NID during configuration. The `if_arpres` field of the `ifnet` structure is assigned a pointer to the NID specific ARP resolution routine before calling `if_attach` kernel service. The resolver routine will then be called by the IP layer at the appropriate time by looking at the `ifnet` structure for that particular interface.

The **resolve** routine resolves the IP address into the corresponding hardware address. If successful, `addr_hw` points to the hardware address, and a value of zero is returned to the caller.

If there is no entry in the ARP table (`arptab`), one needs to be created. An entry is created with the network address of `sock_addr_net` and a broadcast ARP request is sent out (using `xx_arpwhohas`). The mbuf structure containing the data is held until the address is resolved. A value of 1 is returned to the caller. The response for the ARP request is eventually received by the `xx_arpinput` routine. This routine would update the ARP table and send out the held packet.

```
(*resolve)(ac, m, sock_addr_net, addr_hw)
struct arpcom *ac;
struct mbuf *m;
struct in_addr *sock_addr_net;
caddr_t *addr_hw;
```

Based on the interface type of the resolver, the `addr_hw` parameter has different semantics:

Ethernet `addr_hw` is used as type `struct ether_header *eh`. It is used to return the value of the Ethernet address.

Token-ring (802.5)

`addr_hw` is used as type `struct ie5_mac_hdr *macp`. It completes the MAC and LCC headers.

802.3 `addr_hw` is used as type `struct ie3_mac_hdr *macp`. It completes the MAC and LCC headers.

FDDI `addr_hw` is used as type `struct fddi_mac_hdr *macp`. It completes the MAC and LCC headers.

The ARP input routine is specified by the NID during configuration. The `if_arpinput` field of the `ifnet` structure is assigned a pointer to the NID specific ARP resolution routine before calling the `if_attach` kernel service. The ARP input routine will then be called by the receive function of the interface demuxer at the interrupt level.

For Local Area Networks, several resolver and arp input routines are present with the NIDs provided with the system. These NIDs, in `/usr/lib/drivers`, are for the following interfaces:

if_en	Standard Ethernet and IEEE 802.3 Ethernet NID
if_fd	FDDI NID
if_tr	IEEE 802.5 token-ring NID

If you choose to use the arp routines provided with the system instead of writing your own, use the following values for the *if_arpres* and *if_arpinput* fields of the **ifnet** structure:

Interface Type	if_arpres Value	if_arpinput Value
Standard Ethernet	arpresolve	arpinput
FDDI	fddi_arpresolve	fddi_arpinput
802.5	ie5_arpresolve	ie5_arpinput
802.3	arpresolve	arpinput

The **whohas** routine broadcasts an ARP request packet asking who has the sock_addr_net Internet Address.

```
(*whohas)(ac, sock_addr_net)
struct arpcom *ac;
struct in_addr *sock_addr_net;
```

The **arptfree** routine frees any ARP table entry specified by *ac. Any mbuf chain of data associated with this arptab entry and held for future transmission, is also freed. (See the previous discussion of the **resolve** routine, about one page earlier, in this same section.)

```
(*arptfree)(at)
struct arptab *at;
```

Handling NID Specific ioctl Calls

The NID supports at least the following **ioctl** commands:

- SIOCSIFADDR** Sets the Network Interface address.
- The SIOCIFADDR command does not update the interface address list. The address list is updated by the network (for example, TCP/IP) layer. See the **ifaddr** structure on page 12-37 for more information.
- The SIOCIFADDR only adds the IP address to the **arpcom** structure for the address in the AF_INET domain.
- SIOCSIFFLAGS** Sets interface flags.
- The SIOCSIFFLAGS command modifies the *if_flags* field in the **ifnet** structure for this NID. The different flags that are stored in the *if_flags* field are as follows (also see the **net/if.h** header file):
- SIOCIFDETACH** Calls the **ns_free** function to deallocate the given device driver.
- FIIOCTL_ADD_FILTER** Calls the **ns_add_filter** CDLI service to add the given filter to the device driver
- FIIOCTL_DEL_FILTER** Calls the **ns_del_filter** CDLI service to remove the given filter
- SIOCADDMULTI** Adds a multicast address. The IOCTL must support subfunctions for adding a multicast address and for enabling a multicast address.
- SIOCDELMULTI** Deletes a multicast address. The IOCTL must support subfunctions for deleting a multicast address and for disabling all multicasts.

The following sample code shows an IOCTL routine:

```
xx_ioctl(ifp, cmd, data)
register struct ifnet          *ifp;
int                            cmd;
caddr_t                        data;
```

```

{
    register struct ifaddr          *ifa = (struct ifaddr *)data;
    register struct xx_softc        *xyp = &xx_softc[ifp->if_unit];
    int                             error = 0;
    struct timestruc_t              ct;

    if (!xyp->nndp) {
        return(ENODEV);
    }

    switch (cmd) {
        case SIOCIFDETACH:
            ns_free(xyp->nndp);
            xyp->nndp = 0;
            break;

        case IFIOCTL_ADD_FILTER:
            {
                ns_8022_t          *filter;
                ns_user_t          *user;

                filter = &((struct if_filter *)data)->filter;
                user = &((struct if_filter *)data)->user;
                error = ns_add_filter(xyp->nndp, filter, sizeof(*filter),
                                    user);
                if (error)
                    bsdlog(LOG_ERR,
                            "if_fd: ns_add_filter() failed with %d return code.\n",
                            error);

                break;
            }

        case IFIOCTL_DEL_FILTER:
            ns_del_filter(xyp->nndp, (ns_8022_t *)data, sizeof(ns_8022_t));
            break;

        case SIOCSIFADDR:
            switch (ifa->ifa_addr->sa_family) {
                case AF_INET:
                    ((struct arpcom *)ifp)->ac_ipaddr =
                        IA_SIN(ifa)->sin_addr;
                    break;

                default:
                    break;
            }
            fd_init();
            ifp->if_flags |= IFF_UP;
            curtime(&ct);
            ifp->if_lastchange.tv_sec = (int)ct.tv_sec;
            ifp->if_lastchange.tv_usec = (int)ct.tv_nsec / 1000;
            break;

        case SIOCSIFFLAGS:
            fd_init();
            break;

        case SIOCADMULTI:
            {
                register struct ifreq *ifr = (struct ifreq *)data;
                void xx_map_ip_multicast();
                char      addr[6];

                /*
                 * Update our multicast list.
                 */
            }
    }
}

```

```

switch(driver_addmulti(ifr, &(xyp->fd_ac),
    xx_map_ip_multicast, &error, addr)) {
case ADD_ADDRESS:
    error = (*(xyp->nndp->ndd_ctl))
            (xyp->nndp, NDD_ENABLE_ADDRESS, addr,
             FDDI_ADDRLEN);

    if(error) {
        int rc;
        driver_delmulti(ifr, &(xyp->fd_ac),
            xx_map_ip_multicast, &rc, addr);
    }
    break;

case ENABLE_ALL_MULTICASTS:
    error = (*(xyp->nndp->ndd_ctl))
            (xyp->nndp, NDD_ENABLE_MULTICAST, addr,
             FDDI_ADDRLEN);

    if(error) {
        int rc;
        driver_delmulti(ifr, &(xyp->fd_ac),
            xx_map_ip_multicast, &rc, addr);
    }
    break;

case 0:
    /* address already enabled */
    break;
case -1:
    /* error */
    break;
}
break;
}

case SIOCDELMULTI:
{
register struct ifreq *ifr = (struct ifreq *)data;
void xx_map_ip_multicast();
char addr[6];
/*
 * Update our multicast list.
 */
switch(driver_delmulti(ifr, &(xyp->fd_ac),
    xx_map_ip_multicast, &error, addr))
{
case DEL_ADDRESS:
    error = (*(xyp->nndp->ndd_ctl))
            (xyp->nndp, NDD_DISABLE_ADDRESS, addr,
             FDDI_ADDRLEN);

    break;

case DISABLE_ALL_MULTICASTS:
    error = (*(xyp->nndp->ndd_ctl))
            (xyp->nndp, NDD_DISABLE_MULTICAST, addr,
             FDDI_ADDRLEN);

    break;

case 0:
    /* address still in use */
    break;
case -1:
    /* error */
    break;
}
break;
}
}

```

```

        default:
            error = EINVAL;
    }
}
return (error);
}

```

Two kernel services have been added to help with multicast addressing. They are **driver_addmulti** to add multicast addresses and **driver_delmulti** to delete multicast addresses.

```

driver_addmulti(ifr, *ac, func, error, mac_address)
struct ifreq *ifr;
struct arpcom *ac;
void func();
int *error;
char *mac_address;

```

The func parameter above is a pointer to a user defined routine in the NID that maps the multicast address. The following is an example of such a routine:

```

void xx_map_ip_multicast(struct sockaddr_in *sin, u_char *lo,
u_char *hi)
{
if (sin->sin_addr.s_addr == INADDR_ANY) {
/*
* An IP address of INADDR_ANY means listen to all
* of the xx multicast addresses used for IP.
* (This is for the sake of IP multicast routers.)
*/
bcopy(xx_ipmulticast_min, lo, 6);
bcopy(xx_ipmulticast_max, hi, 6);
} else {
xx_MAP_IP_MULTICAST(&sin->sin_addr, lo);
bcopy(lo, hi, 6);
}
}
}

```

The function **driver_delmulti** has the same parameters.

```

driver_delmulti(ifr, *ac, func, error, mac_address)
struct ifreq *ifr;
struct arpcom *ac;
void func();
int *error;
char *mac_address;

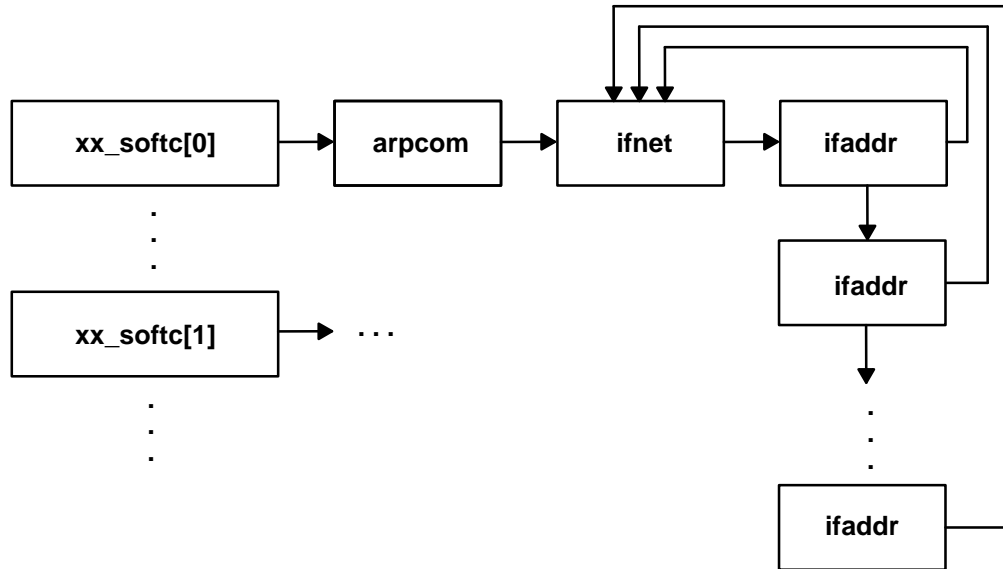
```

Terminating

The IOCTL SIOCIFDETACH discussed in “Handling NID Specific IOCTL Calls”, on page 12-31, supports the termination of the NDD and corresponding NID.

NID and ARP Data Structures

This section lists data structures used for NID and ARP. The following NID Data Structure Relationships figure shows the relationships between various NID and ARP data structures.



NID Data Structure Relationships

xx_softc

The **xx_softc** structure shows *software status* and is an internal structure of the NID. This is a starting point for locating the address of other data structures for the NID.

```

struct xx_softc {
    struct arpcom  xx_ac; /* common part */
    struct ndd    *nddp; /* ptr to device driver struct */
}
  
```

arpcom

The **arpcom** structure is a *network common* structure. It is shared between the NID and the address resolution code. This is the first element in the **softc** structure.

```

struct arpcom {
    struct ifnet  ac_if;
    u_char  ac_hwaddr[MAX_HWADDR];
    struct in_addr ac_ipaddr;
    struct driver_mult *ac_multiaddrs; /* Lock for walking list */
                                        /* of multicast addrs. */
};
  
```

ifnet

The **ifnet** structure is the Network Interface Table. It has the following types of fields:

- Interface identifier (if_name, ...)
- Interface properties (if_mtu, if_flags, ...)
- Interface routines (if_output, if_ioctl, ...)
- Interface statistics (if_ipackets, if_opackets, if_ierrors, ...)

It also maintains a pointer to the linked list interface addresses (if_addrlist) for the interface. The **ifnet** structures for the different interfaces are in a linked list (if_next). The following code sample shows the **ifnet** structure:

```

struct ifnet {
    char    *if_name;           /* name, e.g. ``en'' or ``lo'' */
    short   if_unit;           /* sub-unit for lower level driver */
    u_long  if_mtu;            /* maximum transmission unit */
    u_long  if_flags;          /* up/down, broadcast, etc. */
    short   if_timer;          /* time 'til if_watchdog called */
    int     if_metric;          /* routing metric (external only) */
    struct  ifaddr *if_addrlist; /* linked list of addresses per if */
/* procedure handles */
    int     (*if_init)();       /* init routine */
    int     (*if_output)();     /* output routine (enqueue) */
    int     (*if_start)();      /* initiate output routine */
    int     (*if_done)();       /* output complete routine */
    int     (*if_ioctl)();      /* ioctl routine */
    int     (*if_reset)();      /* bus reset routine */
    int     (*if_watchdog)();   /* timer routine */
/* generic interface statistics */
    int     if_ipackets;        /* packets received on interface */
    int     if_ierrors;         /* input errors on interface */
    int     if_opackets;        /* packets sent on interface */
    int     if_oerrors;         /* output errors on interface */
    int     if_collisions;      /* collisions on csma interfaces */
/* end statistics */
    struct  ifnet *if_next;
    u_char  if_type;            /* ethernet, tokenring, etc */
    u_char  if_addrln;         /* media address length */
    u_char  if_hdrln;          /* media header length */
    u_char  if_index;          /* numeric abbreviation for this if */
/* SNMP statistics */
    struct  timeval if_lastchange; /* last updated */
    int     if_ibytes;          /* total number of octets received */
    int     if_obytes;          /* total number of octets sent */
    int     if_imcasts;         /* packets received via multicast */
    int     if_omcasts;         /* packets sent via multicast */
    int     if_iqdrops;         /* dropped on input, this interface */
    int     if_noproto;         /* destined for unsupported protocol */
    int     if_baudrate;        /* linespeed */

/* stuff for device driver */
    dev_t   devno;              /* device number */
    chan_t  chan;               /* channel of mpx device */
    struct  in_multi *if_multiaddrs; /* list of multicast addresses */
    int     (*if_tap)();        /* packet tap */
    caddr_t if_tapctl;          /* link for tap (ie BPF) */
    int     (*if_arpres)();     /* arp resolver routine */
    int     (*if_arpresv)();    /* Reverse-ARP input routine */
    int     (*if_arpinput)();   /* arp input routine */
    struct  ifqueue {
        struct  mbuf *ifq_head;
        struct  mbuf *ifq_tail;
        int     ifq_len;
        int     ifq_maxlen;
        int     ifq_drops;
    } if_snd;                  /* output queue */
    simple_lock_data_t if_slock; /* statistics lock */
    simple_lock_data_t if_multi_lock;
};

```

ifaddr

The **ifaddr** structure contains information about one interface address. The structures are maintained by the different address families (for example, Internet, osi, and xns). They are allocated and attached by the address families, and *not* by the NID. All the addresses for an interface are linked so that they are easy to locate.

```
struct ifaddr {
    struct sockaddr *ifa_addr;      /* address of interface */
    struct sockaddr *ifa_dstaddr;   /* other end of p-to-p link */
#define ifa_broadaddr ifa_dstaddr /* broadcast address interface */
    struct sockaddr *ifa_netmask;   /* used to determine subnet */
    struct ifnet *ifa_ifp;         /* back-pointer to interface */
    struct ifaddr *ifa_next;       /* next address for interface */
#ifdef _KERNEL
    void (*ifa_rtrequest)(int, struct rtentry *, struct sockaddr *);
#else
    void (*ifa_rtrequest)();       /* check or clean routes (+ or -)'d */
#endif
    struct rtentry *ifa_rt;        /* ??? for ROUTETOIF */
    u_short ifa_flags;            /* mostly rt_flags for cloning */
    u_short ifa_llinfolen;        /* extra to malloc for link info */
};
```

ifreq

The **ifreq** (interface request) structure is used for socket IOCTLs. All interface IOCTLs must have parameter definitions that begin with `ifr_name`. The remainder can be interface-specific.

```
struct ifreq {
#define IFNAMSIZ 16
    char ifr_name[IFNAMSIZ];      /* if name, e.g. "en0" */
    union {
        struct sockaddr ifru_addr;
        struct sockaddr ifru_dstaddr;
        struct sockaddr ifru_broadaddr;
        long ifru_flags;
        int ifru_metric;
        caddr_t ifru_data;
        short ifru_mtu;
    }
};
```

arptab

The **arptab** (ARP table entries) structures contain the network address to link-layer address translation. The table is maintained by the ARP routines.

The entries in this table are hashed for fast retrieval. The table is divided into `ARPTAB_NB` buckets, each of size `ARPTAB_BSIZ`. The network address entry to be stored is hashed based upon its *at_addr* value into the appropriate bucket. To work on this table, `ARPTAB_HASH` and `ARPTAB_LOOK` are defined.


```

struct arptab {
    struct in_addr  at_iaddr; /* internet address */
    u_char          hwaddr[MAX_HWADDR]; /* hardware address */
    u_char          at_timer; /* minutes since last reference */
    u_char          at_flags; /* flags */
    struct mbuf     *at_hold; /* last pkt til resolved/timeout */
    struct ifnet    *at_ifp; /* ifnet assoc with entry */
    union if_dependent if_dependent; /* hdwr dependent info */
};

```

arpreq

Use the **arpreq** (ARP request from user) structure to initiate a socket ioctl request.

```

struct arpreq {
    struct sockaddr  arp_pa; /* protocol address */
    struct sockaddr  arp_ha; /* hardware address */
    int             arp_flags; /* flags */
    u_short         at_length; /* length of hdwr addr */
    union if_dependent ifd; /* hdwr dependent info */
    u_long          ifType; /* interface type */
}

```

Include Files

The include files for the data structures include:

```

net/if.h
net/if_arp.h
sys/mbuf.h
sys/socket.h
sys/devinfo.h

```

Tracing and Debugging for NIDs

The **trace** tool is useful for debugging the NID kernel extension. It is also a useful mechanism for performance analysis and tuning. The tracing code has a small overhead, so the system performance is minimally altered by using tracing code. To add **trace** code, see the chapter on the trace facility in *AIX Version 4.1 General Programming Concepts Volume 2: Debugging Programs* and “Performance Tracing” on page 14-61. The process of tracing is described in the following paragraphs.

Each logical module inside the kernel and kernel extension is allocated a unique hook ID. For example:

- HKWD_SOCKET for the socket module
- HKWD_MBUF for the mbuf module
- HKWD_IFEN for the Ethernet NID.
- HKWD_IFFD for the FDDI NID
- HKWD_IFET for the Ethernet 802.3 NID
- HKWD_IFTR for the Token Ring 802.5 NID
- HKWD_IFSL for the SLIP NID

For each of the hook IDs, there are sub-hook IDs for the procedures contained in the related modules.

For example, for the Ethernet NID output routine, the sub-hook IDs are `hkwd_output_in` and `hkwd_output_out`. A new kernel extension can use tracing by allocating itself a new unique hook ID which does *not* have the same value as an existing hook ID. A combination of hook ID and sub hook ID is used to create a unique trace event in the trace log.

For example, at the beginning of the Ethernet ARP resolve routine, there is the tracing call:

```
TRCHKT(HKWD_IFEN | hkwd_output_in)
```

Start tracing by using the **trace** command. Stop it by calling **trcstop**. Process the trace log by using the **trcrpt** command.

The **net_error** kernel service makes a tracing call to trace the error generated.

```
TRCHKL1(HKWD_NETERR | error_code, ifp)
```

Configuration Method for NID

Since NID is a dynamically loadable kernel extension, it has to be loaded by an application which is called as a configuration method. For the existing NID of AIX this is achieved using the **ifconfig** command.

To add your own NID used by AIX TCP/IP, create a configuration method that loads your own NID. This configuration method will be similar to **ifconfig**.

The configuration method uses the **sysconfig** subroutine to load, unload, and configure the AIX NID.

Chapter 13. Network Interfaces and Protocols

AIX Version 4.1 provides two standard user application interfaces to the network: sockets and STREAMS. The figure CDLI Device Driver Structure, on page 12-1, graphically displays the network interface architecture.

Network device drivers are not required to support the socket and STREAMS interfaces directly. Instead device drivers interface with kernel services which provide the upper layer support. Collectively these kernel services are called the Common Data Link Interface (CDLI). Nor are device drivers required to support direct user access to the driver. Raw socket and STREAMS interfaces are provided for this purpose. Because device drivers only interface with well defined *kernel* services, this architecture greatly simplifies the writing and porting of network device drivers.

AIX Version 4.1 supports user written STREAMS and socket network protocols. For the STREAMS user, both a tli and an xti application library is provided along with the OSF/ 1.2 STREAMS system call framework. The xti application interface is X/Open compliant. STREAMS network protocols written to the TPI (Transport Provider Interface) should be able to directly link into the STREAMS application interface from below. A Data Link Provider Interface (DLPI) STREAMS module is provided to connect STREAMS-based protocols with network device drivers. In general, STREAMS protocols should not have to make direct calls to either CDLI or the network device driver. For socket users, AIX Version 4.1 offers a standard BSD socket framework with a set of kernel services which permit users to dynamically add both communication domains and new protocol switch entries to an existing communications domain, but this feature is currently supported only for the AF_INET address family. In addition, AIX exports a number of low level socket calls such as sbappend, sbdrop and sbreserve which are required by socket based protocols. With these services, user written socket based network protocols should be able to directly link into the socket application interface from below. The standard BSD network interface layer (ifnet) is provided to connect socket based protocols with network device drivers. In general, socket based protocols should not have to make direct calls to either CDLI or the network device driver.

Finally, both network device drivers and network protocols are implemented as loadable kernel extensions. Special commands, (such as **strload** for STREAMS) or special user written commands (similar to **ifconfig** for sockets) accomplish the loading of these kernel extensions.

Detailed Network Interfaces

STREAMS User Interfaces

Both the Transport Layer Interface (TLI) and the X/Open Transport Interface (XTI) define a transport service interface. The TLI and XTI implementations are limited to TCP/IP. TLI and XTI are accessed through a set of library subroutines for the C programming language. These libraries access the kernel through STREAMS messages.

TLI and XTI provide both a connection mode and a connectionless mode of transport services. A transport endpoint specifies a communication path between the transport user and the transport provider. A single transport endpoint may not support both modes of service simultaneously.

TLI and XTI support both synchronous and asynchronous execution modes for handling asynchronous events. In synchronous mode, the user program is blocked until a specific event has occurred. In asynchronous mode, the user process is not blocked and is notified when the specific event has occurred.

Both TLI and XTI support a rich option management facility to provide negotiation, debugging, graceful shutdown, and buffer management functionalities.

The TLI implementation is based on AT&T SVR4 (System V release 4). The transport provider driver, XTISO, and the associated STREAMS module, timod, are based on the Transport Provider Interface (TPI) version 1.5. The XTI implementation complies with the XPG4 standards defined by X/Open.

The user can also directly call the STREAMS **getmsg** and **putmsg** system calls to read from and write to the stream-head message queue.

For more information see “xtiso STREAMS Driver” and “dlpi STREAMS Driver” in *AIX Version 4.1 Technical Reference, Volume 4: Communications*, “Transport Service Library Interface Overview” in *AIX Version 4.1 Communications Programming Concepts*, and the reference information on **getmsg** and **putmsg** in *AIX Version 4.1 Technical Reference, Volume 4: Communications*.

Protocol Interfaces via DLPI

The DLS provider is implemented as a style 2 provider which supports the connectionless mode of communication. The connectionless DLPI service supports local management primitives and data-transfer primitives. Local management primitives allow the DLS user to query and control the DLS provider. Data-transfer primitives enable the DLS user to communicate with a peer DLS user. Each local management primitive and data-transfer primitive is implemented as a STREAMS message. Normal primitives are implemented as M_PROTO messages. High-priority interface acknowledgement primitives are implemented as M_PCPROTO messages.

A style 2 DLS provider can support several physical points of attachment (PPA). The PPA is used to identify one of several of the same type of interface in the system. A DLS user must explicitly identify the PPA using the DL_ATTACH_REQ primitive. For more information, see the DL_ATTACH_REQ primitive manpage.

The DLS provider has been implemented to allow the DLS user the capability of specifying the packet format. Using the M_IOCTL STREAMS message, the DLS user can specify the packet format. If the DLS user does not specify the packet format, the default is NS_PROTO. The packet formats are defined in the **ns_user** structure article.

In order for the DLS provider to generically support all interface types, the DLS provider has been implemented so the DLS user can specify address resolution routines.

For more information, see the reference articles on the DLPI primitives (such as DL_ATTACH_REQ) in *AIX Version 4.1 Technical Reference, Volume 3: Communications* and *AIX Version 4.1 Technical Reference, Volume 4: Communications*.

Writing or Porting STREAMS Network Protocols

This section:

- Lists and explains DLPI interfaces supported by AIX
- Details the AIX interpretations of source and destination address
- Provides sample code for packet format setting
- Points to demuxer docs

DLPI Interfaces Supported by AIX

DLPI supports CDLI-based network interfaces in a generic fashion. This support is enabled by allowing the DLPI user to specify the particular packet format necessary for the transmission media over which the stream is created. This generic interface support allows new CDLI-based device drivers to be added to the operating system without updating the DLPI driver, but with some modifications in the associated demuxer.

The DLPI user is allowed one packet format specification per stream. This packet format must be specified after the attach, and before the bind. Otherwise, an error will be generated. The DLPI user can specify one of the following packet formats per stream: NS_PROTO, NS_PROTO_SNAP, NS_INCLUDE_LLC, and NS_INCLUDE_MAC. These packet formats are defined in the `/usr/include/sys/cdli.h` file. If the user does not specify a packet format, the default packet format is NS_PROTO.

For the DL_UNITDATA_IND primitive, DLPI will provide the header information in the `dl_unitdata_ind_t` structure. If the packet format specified is NS_PROTO or NS_PROTO_SNAP, the MAC and LLC are included in the header, and the data portion of the message will contain only data. If the packet format is NS_PROTO, the DLPI header includes the MAC and LLC without the SNAP. If the packet format is NS_PROTO_SNAP, the DLPI header includes the MAC, LLC, and SNAP. If the packet format specified is or NS_INCLUDE_MAC, the DLPI header will contain only the destination and source addresses. If the packet format is NS_INCLUDE_LLC, only the LLC will be placed in the data portion of the message. If the packet format is NS_INCLUDE_MAC, the MAC and LLC are both placed in the data portion of the message. Thus, the DLPI user must have knowledge of the MAC header and LLC architecture for that interface in order to retrieve the MAC header and LLC from the data portion of the message.

For the DL_UNITDATA_REQ primitive, if the DLPI user had specified either the NS_PROTO, NS_PROTO_SNAP, or NS_INCLUDE_LLC format, the DLPI user must provide the destination address and an optional DSAP in the DLPI header. If the DLPI user does not specify the DSAP, the DSAP specified at bind time will be used. If the DLPI user specifies the NS_INCLUDE_LLC packet format, the user must include only the LLC in the data portion. If the user specifies the NS_INCLUDE_MAC packet format, the DLPI user must provide the full MAC header, including the LLC, in the data portion of the message.

The DLPI user specifies the packet format via the STREAMS `I_STR` IOCTL.

See “Obtaining Copies of the DLPI and TPI Specifications”, on page 13-5, for information on how to get the DLPI specifications.

AIX Interpretations of Source and Destination Addresses

DLPI source and destination addresses are 6 byte hardware addresses.

TLI and XTI source and destination addresses are specified via the **sockaddr_in** structure defined in the **/usr/include/netinet/in.h** file. This structure requires a family, a port number, and an IP address.

Protocol Address Resolution

The DLPI user can provide an address resolution procedure for input and output using the **STREAMS_IOCTL**, or the user can rely on the system default address resolution routines. AIX provides default address resolution routines that are interface specific. The DLPI user must specify the address resolution routine in the **ndd** structure defined in the **/usr/include/sys/ndd.h** file.

AIX STREAMS Loading Convention

The configuration file for DLPI is the **/etc/dlpi.conf** file.

To load: `strload -f /etc/dlpi.conf`

To unload: `strload -uf /etc/dlpi.conf`

For STREAMS modules and drivers that can be accessed by TLI and XTI, the configuration file is the **/etc/xtiso.conf** file.

To load: `strload -f /etc/xtiso.conf`

To unload: `strload -uf /etc/xtiso.conf`

MP Serialization and Locking Options for STREAMS Modules and Drivers

The STREAMS framework has a special set of data structures and synchronizations that enable STREAMS-based devices to operate in a multi-threaded environment.

At configuration time, users supply valid synchronization levels for STREAMS modules and drivers. The synchronization level is specified in the `sc_sqllevel` member of the **strconf_t** structure passed in as an argument into the **str_install** configuration procedure call. In addition, MP-safe and MP-efficient STREAMS drivers and modules are required specify in the `sc_flags` member of the **strconf_t** structure the style of the open routine logically ORed with **STR_MPSAFE**.

Valid synchronization levels are:

SQLVL_QUEUE

Queue Level. This synchronization level enables a separate thread of execution to access either side of the stream simultaneously. This synchronization level provides the finest degree of parallelization. If the user specifies the **SQLVL_QUEUE** level of synchronization, the user may need to provide locks to ensure that only one thread executes within a queue at a time. If the user must provide a locking mechanism, it is recommended that the user employ the `q_lock`.

SQLVL_QUEUEPAIR

Queue Pair Level. This synchronization level guarantees that only one thread of execution can access either queue of the queue pair at a time.

SQLVL_MODULE

Module Level. This level specifies synchronous across all instances of a module, ensuring no more than one thread of execution through all

instances of the module at a time. This is the level of synchronization the user should select in a non MP-safe STREAMS module.

SQLVL_ELSEWHERE

Arbitrary Level. A cooperating group of modules, such as a protocol family, may need to ensure that there is only one thread of execution through the entire group. In this case, the module developer provides a unique name that is used at configuration time. The unique name is specified in the `sc_sqinfo` member of the `strconf_t` structure.

SQLVL_GLOBAL

Global Level. This synchronization level forces a single thread access through all streams. This option is normally used only for debugging. With this level of synchronization, the user requests a single lock for the entire streams system. Only one thread at a time may be executing.

SQLVL_DEFAULT

Default Level. This synchronization level is defined as `SQLVL_MODULE`.

TLI and XTI Interface Protocols

The Transport Provider Interface (TPI) is a STREAMS message interface that specifies the types and allowable sequences of messages passed between the transport user and the transport provider. TPI is defined by a set of primitives which are implemented as STREAMS messages. These STREAMS messages can consist of `M_PROTO`, `M_PCPROTO`, or `M_DATA` message blocks. The TLI and XTI libraries are implemented using TPI primitives.

The TLI module, `timod`, sits between the stream head and the transport provider and helps map TLI messages to TPI primitives. `timod` passes most messages along unchanged.

Obtaining Copies of the DLPI and TPI Specifications

You can obtain copies of the data link provider interface (DLPI) and transport provider interface (TPI) specifications electronically or through direct mail.

A postscript file of the DLPI and TPI specifications may be retrieved electronically by anonymous ftp from the internet host file `ftp.ui.org` residing at IP address 192.203.65.200 under the `pub/osi` directory.

To retrieve the postscript DLPI and TPI specifications through anonymous ftp, use the following example:


```
ftp ftp.ui.org
Connected to ftp.ui.org.
220 uiunix FTP server (Version 1.25 Jul 09 1993) ready.
Name (ftp.ui.org:jhaug):anonymous
ftp> user anonymous
331 Guest login ok, send ident as password.
Password: guest
230 Guest login ok, access restrictions apply.
ftp> cd pub/osi
250 CWD command successful.
ftp> bin
200 Type set to I.
ftp> get dlpi.ps      for TPI: ftp> get tpi.ps
200 PORT command succesful.
150 Opening BINARY mode data connection for dlpi.ps (1476915 bytes).
226 Transfer complete.
1476915 bytes received in 440.4 seconds (3.275 Kbyte/s)
ftp> quit
221 Goodbye.
```

If you do not have access to the internet, you can obtain a copy of the DLPI and TPI specifications directly from UNIX International through the U.S. Postal Service. The mailing address of UNIX International is:

```
UNIX International
Waterview Corporate Center
20 Waterview Boulevard
Parsippany, New Jersey 07054
```

Writing or Porting Socket Network Protocols

Socket protocols, including the system provided TCP/IP and XNS protocols, are implemented as loadable kernel extensions in AIX Version 4.1. The following general guidelines apply:

- By convention, protocol kernel extensions are installed in the **/usr/lib/drivers** directory. Some system utilities may assume that this is the case.
- When building the kernel extension the relevant base system exports must be imported. The system provided protocols import the following: **kernex.exp**, **syscalls.exp**, **sockets.exp** and **statcmd.exp**. The system export files are located in the **/usr/lib** directory. These exports should provide all of the services and data structures required to port standard BSD socket-based protocols.
- If other kernel extensions are using services provided by the new protocol then users must create an export file and export these services. See **/usr/lib/netinet.exp** for an example of a protocol export file.
- Protocol writers must decide how much of the kernel extension to pin. The major consideration is that the system interrupt handlers will call the protocol's handler with interrupts disabled. (By the way, this is not true for the protocol's fast and slow timers.) So at minimum, the protocol's interrupt handler must be pinned. The system provided socket protocols are pinned in their entirety.

Initialization

There are two phases to socket protocol initialization in AIX. The first phase involves execution of the kernel extension's configuration entry point. This function is designated when the kernel extension is built and is called by the system when the kernel extension is loaded. This function should perform the following tasks:

- Do any lock initialization required by the protocol in an MP environment.
- Do any pinning of modules or data structures required by the protocol.
- Add the protocol's communication domain to the system list using the **domain_add** kernel service. **domain_add** will call the domain's initialization function plus initialization function of all of the protocol's listed in the domain's protocol switch table. The later step is phase two of protocol initialization. Definitions related to this are in **/usr/include/sys/domain.h** and **/usr/include/sys/protosw.h**
- Register address resolution functions, loopback handlers and address resolution IOCTLS with the system using the **nd_config_proto** kernel service.

This is illustrated in sample code for a socket protocol's configuration entry point function, on page 13-15.

The second phase of protocol initialization occurs when the protocol initialization function is called by **domain_add**. This is exactly analogous to what happens on BSD systems. Generally, these initialization procedures should contain all of the protocol initialization procedures which are not AIX specific. Protocols being ported from Berkeley systems will require no changes to their initialization procedures with two exceptions:

- If protocol interrupts are to be scheduled using the AIX software interrupt facilities, then this should be initialized at this point. This is done via a call to the **netisr_add** kernel service.
- Perform any lock initialization required by the protocol in an MP environment.

This is illustrated in sample code for a socket protocol's initialization function, on page 13-16.

After this initialization process, the socket protocol is loaded into the kernel, configured into the system's socket framework and initialized. All that remains is for the protocol to notify the system of the types of network packets it wishes to receive, from which network interfaces and the format in which packets are to be exchanged with the system. This is accomplished through a call to the **ns_add_filter** kernel service. Because AIX supports multiple protocols concurrently on the same network adapter and because a protocol may not want to receive packets from all of the network interfaces configured into the system, socket protocols should perform this registration when an address is bound to a network interface.

This is illustrated in sample code for a socket protocol's packet registration function, on page 13-17.

Loading

The system provided socket protocols are loaded and configured by the **ifconfig** command. Because **ifconfig** requires prior knowledge of all the communication address families in the system, it cannot be used to load and configure user written socket protocols. The user must provide a configuration command. The basic logic of this command should be as follows:

1. Check if the protocol is loaded using the **sysconfig** kernel service. If not, load the protocol using the **sysconfig** kernel service.
2. Open a socket.
3. Issue the appropriate socket IOCTL. For additional information, see **/usr/include/sys/ioctl.h** and the reference articles for the socket IOCTLs.

Socket-Protocol Interface

The interfaces that a socket protocol must support are described in the protocol switch structure, defined in **/usr/include/sys/protosw.h**. These interfaces, along with their call semantics, are:

```
void pr_input(defined per communications domain),
int pr_output(defined per communications domain),
void pr_ctlinput(int cmd, struct sockaddr *sa, caddr_t arg)
    /* cmd is one of the PRS commands listed in protosw.h,
       sa is a sockaddr, and
       arg is an optional argument used within
       the protocol family
    */

int pr_ctloutput(int req, struct socket *so,int level,
                int optname,mbuf **optval)
    /* req is a PRCO action listed in protosw.h,
       so is a socket,
       level is an indication of which protocol layer,
       optname is a protocol dependent request value,
       optval is for return results
    */
```

```

int pr_usrreq(struct socket *so, int req, struct mbuf *m,
             struct mbuf *nam, struct mbuf *control)
/* so is the socket,
   req is a PRU request listed in protosw.h,
   m is an optional message chain,
   nam is an optional mbuf containing an address,
   control is an optional mbuf containing control information
   void pr_init(void),
   void pr_fastimo(void),
   void pr_slowtimo(void),
   void pr_drain(void).
*/
*/

```

The socket system calls interact with the protocol solely through the protocol switch structure.

The **pr_init** function is called by the system when the protocol is loaded. After this is accomplished, the system will call the **pr_fastimo** function on a 200 millisecond timer and the **pr_slowtimo** function on a 500 millisecond timer. Unlike other BSD-based systems, AIX calls the protocol's **pr_drain** functions when the system detects a shortage of network memory buffers (mbufs).

Protocols pass data among themselves (for example ip to tcp) using the **pr_input** and **pr_output** functions. **pr_output** moves data towards the network interface and **pr_input** moves data towards the socket system call interface. Control information is passed using the **pr_ctlinput** and **pr_ctloutput** functions. Unlike with **pr_output**, the socket system routines will pass control data down to the protocols using **pr_ctloutput**. The **getsockopt** and **setsockopt** socket system calls are implemented in this fashion.

With the two exceptions noted above, all of the socket-to-protocol interfaces in the system are implemented using **pr_usrreq**. The specific call semantics for each socket system call are as follows:

This is the semantics of the **socket** system call:

```

(pr_usrreq)((struct socket *) so, PRU_ATTACH,
           (struct mbuf *)0, (struct mbuf *)proto,
           (struct mbuf *)0);

```

This is the semantics of the **bind** system call:

```

(pr_usrreq)((struct socket *) so, PRU_BIND,
           (struct mbuf *)0, nam, (struct mbuf *)0);

```

This is the semantics of the **listen** system call:

```

(pr_usrreq)((struct socket *) so, PRU_LISTEN,
           (struct mbuf *)0, (struct mbuf *)0,
           (struct mbuf *)0);

```

This is the semantics of the **close**, **disconnect** system call:

```

(pr_usrreq)((struct socket *) so, PRU_DISCONNECT,
           (struct mbuf *)0, (struct mbuf *)0,
           (struct mbuf *)0));

```

This is the semantics of the **close** system call:

```

(pr_usrreq)((struct socket *) so, PRU_DETACH,
           (struct mbuf *)0, (struct mbuf *)0,
           (struct mbuf *)0);

```

This is the semantics of the **soabort** system call:

```
(pr_usrreq)((struct socket *) so, PRU_ABORT,
            (struct mbuf *)0, (struct mbuf *)0,
            (struct mbuf *)0);
```

This is the semantics of the **accept** system call:

```
(pr_usrreq)((struct socket *) so, PRU_ACCEPT,
            (struct mbuf *)0, nam, (struct mbuf *)0);
```

This is the semantics of the **connect** system call:

```
(pr_usrreq)((struct socket *) so, PRU_CONNECT,
            (struct mbuf *)0, nam, (struct mbuf *)0);
```

This is the semantics of the **socketpair** system call:

```
(pr_usrreq)(so1, PRU_CONNECT2,
            (struct mbuf *)0, (struct mbuf *)so2,
            (struct mbuf *)0);
```

The semantics for **send**, **sendto**, **sendmsg**, **write** is one of the following forms:

```
(pr_usrreq)((struct socket *) so, PRU_SENDOOB, top,
            addr, control);
/* or */
(pr_usrreq)((struct socket *) so, PRU_SEND, top, addr, control);
```

The semantics for **receive**, **recvfrom**, **recvmsg**, **read** is one of the following forms:

```
(pr->pr_usrreq)((struct socket *) so, PRU_RCVOOB, m,
               (struct mbuf *) (flags & MSG_PEEK), (struct mbuf *)0);
/* or */
(pr_usrreq)((struct socket *) so, PRU_RCVD, (struct mbuf *)0,
            (struct mbuf *) flags, (struct mbuf *)0);
```

This is the semantics of the **shutdown** system call:

```
(pr_usrreq)((struct socket *) so, PRU_SHUTDOWN,
            (struct mbuf *)0, (struct mbuf *)0,
            (struct mbuf *)0);
```

This is the semantics of the **setsockopt** system call:

```
(pr_ctloutput) (PRCO_SETOPT, (struct socket *) so, level,
               optname, &m0);
```

This is the semantics of the **getsockopt** system call:

```
(pr_ctloutput) (PRCO_GETOPT, (struct socket *) so, level,
               optname, mp);
```

This is the semantics of the **ioctl** system call:

```
(pr_usrreq)((struct socket *) so, PRU_CONTROL,
            (struct mbuf *)cmd, (struct mbuf *)data, (struct mbuf *)0));
```

This is the semantics of the **stat** system call:

```
(pr_usrreq)((struct socket *) so, PRU_SENSE,
            (struct mbuf *)ub, (struct mbuf *)0,
            (struct mbuf *)0));
```

This is the semantics of the **getsockname** system call:

```
(pr_usrreq)((struct socket *) so, PRU_SOCKADDR,
            (struct mbuf *)0, m, (struct mbuf *)0);
```

This is the semantics of the **getpeername** system call:

```
(pr_usrreq)((struct socket *) so, PRU_PEERADDR,
            (struct mbuf *)0, m, (struct mbuf *)0);
```

The `pr_flags` field in the protocol switch table describe basic characteristics of the protocol and impact the behavior of the socket interface. Valid values for these flags are listed in **`/usr/include/sys/protosw.h`**. If `PR_CONNREQUIRED` is set then the socket calls will not attempt to transfer data before a connection is established. The `PR_ADDR` flag causes receive data to be preceded by the senders address. The `PR_ATOMIC` flag causes sends to be performed in a single protocol send request. The `PR_WANTRCVD` flag causes the socket routines to notify the protocol when the user has removed data from the socket receive queue. The notification is one of the following socket system calls:

```
(pr->pr_usrreq)((struct socket *) so, PRU_RCVOOB, m,
                (struct mbuf *) (flags & MSG_PEEK),
                (struct mbuf *)0);

(pr_usrreq)((struct socket *) so, PRU_RCVD,
            (struct mbuf *)0, (struct mbuf *)flags,
            (struct mbuf *)0);
```

The `PR_RIGHTS` flag indicates that the protocol supports the passing of access rights.

The Design and Implementation of the 4.3 BSD UNIX Operating System contains additional information on the socket-to-protocol interface.

Protocol-Socket Interface

The AIX kernel establishes a socket structure (refer to **`/usr/include/sys/socketvar.h`**) for all sockets. This structure contains send and receive buffer queues (**`so_snd`** and **`so_rcv`**), a pointer to the protocol's switch table (**`so_proto`**) and a pointer to the protocol's control block (**`so_pcb`**). Protocols establish a back pointer to the socket in the control block (for an example, refer to **`/usr/include/netinet/in_pcb.h`**) which completes the cross linkage. The later three pointers, **`so_proto`**, **`so_pcb`** and the protocol's back pointer are established during the socket system call and the resultant `PRU_ATTACH` `pr_usrreq`. It is through these components of the socket structure that the system's protocol-to-socket interface is largely implemented.

On sending of data, reliable protocols typically use the socket send buffer to hold data until acknowledgment. Data is copied from the send buffer using **`m_copy`** for output. When an acknowledgement is received the protocol removes data from the send buffer with `sbdrop` or `sbdroprecord`.

On receipt of data, protocols employ the **`sbappend`** system services to append data to the appropriate socket's receive buffer. Typically, **`sbappend`** or **`sbappendrecord`** are called after the protocol checks that enough space is available in the receive buffer. This check is performed using the **`sbspace`** kernel service. **`sbappendrecord`** differs from `sbappend` in that the data is treated as being the beginning of a new record. Protocols needed to add either access rights or the sender's address to the receive data employ the **`sbappendaddr`** or **`sbappendrights`** kernel services. For access rights plus data **`sbappendcontrol`** should be used. For sender's address, plus access rights (optional), plus data **`sbappendaddr`** should be employed. Unlike **`sbappend`** or **`sbappendrecord`**, these two kernel services check receive buffer space for the caller. These system services do not wake up waiting receivers, so the protocol must issue an **`sorwakeup`**.

The following sample code illustrates adding data to a socket receive buffer:

```
...
/*
 * Locate pcb for datagram.
 */
nsp = ns_pcblookup(&idp->idp_sna, idp->idp_dna.x_port, NS_WILDCARD);
/*
 * Switch out to protocol's input routine.
 */
nsintr_swtrch++;
...
idp_input(m, nsp);
...
idp_input(m, nsp)
struct mbuf *m;
register struct nspcb *nsp;
{
    register struct idp *idp = mtod(m, struct idp *);
    struct ifnet *ifp = m->m_pkthdr.rcvif;
...
/*
 * Construct sockaddr format source address.
 * Stuff source address and datagram in user buffer.
 */
...
if ( ! (nsp->nsp_flags & NSP_RAWIN) ) {
    m->m_len -= sizeof (struct idp);
    m->m_pkthdr.len -= sizeof (struct idp);
    m->m_data += sizeof (struct idp);
}
if (sbappendaddr(&nsp->nsp_socket->so_rcv, (struct sockaddr *)&idp_ns,
    m, (struct mbuf *)0) == 0)
    goto bad;
sorwakeup(nsp->nsp_socket);
return;
bad:
    m_freem(m);
}
```

Protocol-Network Interface

In AIX Version 4.1, paths through which socket messages can be sent and received are configured into the system via network interfaces. Normally, a hardware or pseudo device is associated with each interface. An interface and its addresses are defined by kernel **ifnet** structures (refer to **/usr/include/net/if.h**). The linked **ifnet** structures provide a list of **socket** interface names and protocol addresses configured in the system. It is through these **ifnet** structures that the system's protocol-to-network interface is largely implemented. This approach provides the socket protocols with a consistent interface to all network hardware devices.

Interface properties, including state information, are conveyed to the protocols through the **ifnet** flags and maximum transmission unit (mtu) fields. The most important interface flags and their meanings are:

IFF_UP Interface is up and available for protocol use
IFF_BROADCAST
 Interface is broadcast capable

IFF_LOOPBACK Interface is a software loopback

IFF_POINTTOPOINT Interface is a point-to-point link (slip)

IFF_RUNNING Interface resources have been allocated

IFF_NOARP Interface should not use address-resolution protocol

IFF_SIMPLEX Interface cannot hear its own transmissions

IFF_DO_HW_LOOPBACK Bypass software loopback

IFF_ALLCAST Token-ring only. Sets all rings broadcast.

IFF_SNAP Ethernet only. Interface is 802.3

Protocols pass data to the network interface employing the **if_output** routine. Each network interface accepts output datagrams of a specified (via the mtu) maximum length. Output occurs when **if_output** is called as follows:

```
(*ifp->if_output)(struct ifnet *ifp, struct mbuf *m, struct
sockaddr *dst, struct rtenry *rt)
```

This has the following parameters:

- `ifp` is the **ifnet** pointer for the interface.
- `m` is the **mbuf** chain to be sent.
- `dst` is the destination address.
- `rt` is an optional routing entry.

The network interface is responsible for encapsulation or decapsulation of any link-layer protocol headers required to deliver the message. This resolution is accomplished as follows:

- If `dst->sa_family` is equal to `AF_UNSPEC`, the link-layer header is copied from `dst->sa_data`. In this case, the protocol provided the link-layer header.
- If `dst->sa_family` is equal to `AF_UNSPEC`, then the network interface calls the protocol's address resolution routine. This routine was registered when the protocol was initialized (see the heading "Initialization" in "Writing or Porting Socket Network Protocols", on page 13-7).

The format of the call to the protocol's address resolution routine is:

```
(*dst->af_family.resolve) (struct arpcom *ac, struct mbuf *m,
struct sockaddr *dst, caddr_t *llh)
```

The call has the following parameters:

- `m` is the mbuf chain to be sent,
- `dst` is the destination address
- `llh` is the link layer header to be filled out by the resolution routine. `llh` must be big enough to hold the largest possible link layer header.

The network interface module generally maintains the **ifnet** data structure as part of a larger data structure (an **arpcom**) that contains interface specific information. Thus `ac` is typically set to `(struct arpcom *) ifp`, where `ifp` is the **ifnet** pointer for the interface. Note

that AIX Version 4.1 provides several generic address resolution routines which may be employed by protocols.

Protocols pass control data to the network interface employing the `if_ioctl` routine. Most importantly, interface addresses are set with IOCTL requests. The IOCTL requests to set interface addresses, (`SIOCSIFADDR`, `SIOCSIFDSTADDR`), and to set and delete multicast addresses (`SIOCADDMULTI`, `SIOCDELMULTI`) generally require work at the interface layer and should be passed along by the protocols with an `if_ioctl`. The specific format of this call is:

```
(*ifp->if_ioctl)(struct ifnet *ifp, int cmd, caddr_t data)
```

This call has the following parameters:

- `ifp` is the `ifnet` pointer for the interface.
- `cmd` is the IOCTL.
- `data` is the IOCTL data.

Network – Protocol Interface

In AIX Version 4.1, receive data is passed directly from the common data link interface to the protocols bypassing the network interface layer. The protocol registers with the system the format and the method of delivery for input packets. This registration is done through the `ns_add_filter` kernel service at the protocol's initialization. By selecting the input packet format, the protocol informs the system whether to strip various portions of the link layer header before receive data is presented. There are two methods of delivery:

- A direct call from the interrupt level to the protocol's input function
- An enqueue of the packet on the protocol's input queue and a schedule of the appropriate software interrupt.

When called directly, the format of the call to the protocol is:

```
(*protocol_input)(struct ndd *ndd, struct mbuf *data, caddr_t *llc, caddr_t *protocookie)
```

The parameters are:

- `ndd` is an `ndd` pointer (see `/usr/include/sys/ndd.h`) to the receiving interface.
- `data` is an `mbuf` pointer to the received data (in the format requested).
- `llc` is a pointer to the link-layer header of the received packet.
- `protocookie` is an address passed in by the protocol when it registered to receive packets of this type. This can be useful for demuxing purposes.

When the packet is enqueued, the protocol receives only the data, in the format requested.

IP Encapsulation/Adding Protocols to the System IP Protocol Switch

The `domain_add` kernel service can be used for adding an entire communications address family with its own protocol switch (see the heading "Initialization" in "Writing or Porting Socket Network Protocols", on page 13-7). To add entries to an existing protocol switch (for example, IP encapsulation or a new protocol within IP) use the `protosw_enable` kernel service. To remove the protocol switch use the `protosw_disable` kernel service. These services currently only support the `AF_INET` communications domain. The user is responsible for pinning the new protocol switch entry.

The following sample code illustrates adding IP protocol switch entries to the system:

```

/* set up protocol switch table in internet protocol */
{
    extern int protosw_enable();
    struct protosw *pr;
    struct protosw xns_sw = { SOCK_RAW,0,IPPROTO_IDP,
                             PR_ATOMIC|PR_ADDR,idpip_input,0,
                             nsip_ctlinput,0,0,0,0,0,0,};
    pr = pffindproto(PF_INET, IPPROTO_RAW, SOCK_RAW);
    if ( pr != 0 ) {
        /* enable the protocol switch so that the IP will handle
         * the incoming xns encapsulating packet and pass along.
         * The route pointers are passed because they are not
         * resolve at load time.
         */
        protosw_enable(&xns_sw);
        /* XXX Should check the rc and do ? */
    }
}

```

Sample Socket Protocol

This sample socket protocol includes the following pieces of sample code:

- Configuration entry point function
- Initialization function
- Packet registration function

Sample Socket Protocol's Configuration Entry Point Function

```

struct protosw nssw[] = {
    { 0,&nsdomain,0,0, 0,idp_output,0,0, 0, ns_init,0,0,0,},
    {SOCK_DGRAM,&nsdomain,0,PR_ATOMIC|PR_ADDR,0,0,idp_ctlinput,idp_ctloutput,
      idp_usrreq,0,0,0,0,},
    {SOCK_STREAM,&nsdomain,NSPROTO_SPP,PR_CONNREQUIRED|PR_WANTRCVD,spp_input,0,
      spp_ctlinput,spp_ctloutput, spp_usrreq,
      spp_init,spp_fasttimo,spp_slowtimo,0,},
    { SOCK_SEQPACKET,&nsdomain,NSPROTO_SPP, PR_CONNREQUIRED|PR_WANTRCVD|PR_ATOMIC,
      spp_input,0,spp_ctlinput,spp_ctloutput, spp_usrreq_sp,0,0,0,0,},
    {SOCK_RAW,&nsdomain,NSPROTO_RAW,PR_ATOMIC|PR_ADDR,idp_input,idp_output,0,
      idp_ctloutput, idp_raw_usrreq,0,0,0,0,},
    { SOCK_RAW,&nsdomain,NSPROTO_ERROR,PR_ATOMIC|PR_ADDR,idp_ctlinput,idp_output,0,
      idp_ctloutput, idp_raw_usrreq,0,0,0,0,},
}
struct domain nsdomain =
    { AF_NS, "network systems", 0, 0, 0,
      nssw, &nssw[sizeof(nssw)/sizeof(nssw[0])],
      0, 0, ns_funnel, ns_funfrfc };
. . .
; /*
 *
 * config_ns - entry point for netns kernel extension
 *
 */
config_ns(cmd, uio)
    int cmd;
    struct uio *uio;
{

```

```

int err, nest;
struct config_proto config_proto;
err = 0;
nest = lockl(&kernel_lock, LOCK_SHORT);
switch (cmd) {
case CFG_INIT:
    /* check if kernel extension already loaded */
    if (extension_loaded)
        goto out;
    ns_lock_init();
/*
 * pin the netns kernel extension
 */
    if (err = pincodex(config_ns))
        goto out; /* Add ns domain */
    domain_add(&nsdomain);
    config_proto.loop = nsintr;
    config_proto.loopq = &nsintrq;
    config_proto.netisr = NETISR_NS;
    config_proto.resolve = ns_arpresolve;
    config_proto.ioctl = NULL;
    config_proto.whohas = NULL;
    nd_config_proto(AF_NS, &config_proto);
    extension_loaded++;
    break;

case CFG_TERM:
default:
    err = EINVAL;
}out:
if (nest != LOCK_NEST)
    unlockl(&kernel_lock);

return(err);
}

```

Sample Socket Protocol's Initialization Function

```

void
ns_init()
{
    struct timestruc_t ct;
    extern void curtime();
    IFQ_LOCK_DECL() ns_broadhost = * (union ns_host *) allones;
    ns_broadnet = * (union ns_net *) allones;
    nspcb.nsp_next = nspcb.nsp_prev = &nspcb;
    nsrawpcb.nsp_next = nsrawpcb.nsp_prev = &nsrawpcb;
    nsintrq.ifq_maxlen = nsqmaxlen;
    curtime(&ct); /* get the current system time */
    ns_pexseq = ct.tv_nsec/1000; /* use microsecond as seq no. */
    ns_netmask.sns_len = 6;
    ns_netmask.sns_addr.x_net = ns_broadnet;
    ns_hostmask.sns_len = 12;
    ns_hostmask.sns_addr.x_net = ns_broadnet;
    ns_hostmask.sns_addr.x_host = ns_broadhost;
    IFQ_LOCKINIT(&nsintrq);
    rtinithead(AF_NS, 16, setnsroutemask);
    (void) netisr_add(NETISR_NS, nsintr, &nsintrq, &nsdomain);
}

```

Sample Socket Protocol's Packet Registration Function

```
/*
 * Generic internet control operations (ioctl's).
 */
ns_control(so, cmd, data, ifp)
    struct socket *so;
    int cmd;
    caddr_t data;
    register struct ifnet *ifp;
{
    register struct ifreq *ifr = (struct ifreq *)data;
    register struct ns_aliasreq *ifra = (struct ns_aliasreq *)data;
    register struct ns_ifaddr *ia;
    struct ifaddr *ifa;
    struct ns_ifaddr *oia;
    struct mbuf *m;
    int error, dstIsNew, hostIsNew;

    . . .

    switch (cmd) {
        . . .
        case SIOCSIFADDR:
            return (ns_ifinit(ifp, ia, (struct sockaddr_ns *)
                &ifr->ifr_addr, 1));
        . . .
    }
    . . .
}
/*
 * Initialize an interface's internet address
 * and routing table entry.
 */
ns_ifinit(ifp, ia, sns, scrub)
    register struct ifnet *ifp;
    register struct ns_ifaddr *ia;
    register struct sockaddr_ns *sns;
{
    struct sockaddr_ns oldaddr;
    register union ns_host *h = &ia->ia_addr.sns_addr.x_host;
    int error;

    . . .

    return(ns_ns_filter(ifp));
}
ns_ns_filter(ifp)
struct ifnet *ifp;
{
    struct ns_user ns_user;
    struct ns_8022 filter;
    struct ndd *nddp;

    char ifname[IFNAMSIZ];
    int rc;
    /*
     * Alloc the ndd. Note that we never free it!!
     */
    sprintf(ifname, "%s%d", ifp->if_name, ifp->if_unit);
    if (rc = ns_alloc(ifname, &nddp))
        return(rc);
    /*
```

```
    * Add 802.3 filter.
    */
    bzero(&filter, sizeof(filter));
    filter.filtertype = NS_8022_LLC_DSAP;
    filter.dsap = DSAP_XNS;
    ns_user.isr = nsintr;
    ns_user.protoq = &nsintrq;
    ns_user.netisr = NETISR_NS;
    ns_user.pkt_format = NS_PROTO;
    ns_user.ifp = ifp;
    rc = ns_add_filter(nddp, &filter, sizeof(filter), &ns_user);
    return(rc);
}
```

Sample Code for Direct Access to Device Driver via STREAMS

Sample client and server source code demonstrating direct user access to device drivers via STREAMS can be found in directory `/usr/sbin/samples/dlpi/`.

Chapter 14. Debugging Tools

This chapter provides information about the available procedures for debugging a device driver which is under development. The procedures discussed include:

- Saving device driver information in a system dump, on page 14-1.
- Using the **crash** command to interpret and format system structures, on page 14-5.
- Using the kernel debugger to set breakpoints and display variables and registers, on page 14-26.
- Error logging to record device-specific hardware or software abnormalities, on page 14-52.
- Using the **trace** facility to monitor entry and exit of device drivers and selectable system events, on page 14-61.

System Dump

The system dump copies selected kernel structures to the dump when an unexpected system halt occurs, when the reset button is pressed, or when the special system dump key sequences are entered. You can also initiate a system dump through the System Management Interface Tool (SMIT).

The dump device can be dynamically configured, which means that either the tape or logical volumes on hard disk can be used to receive the system dump. Use the **sysdumpdev** command to dynamically configure the dump device.

You can also define primary and secondary dump devices. A primary dump device is a dedicated dump device, while a secondary dump device is shared.

The system kernel **dump** routine contains all the vital structures of the running system, such as the process table, the kernel's global memory segment, and the data and stack segment of each process.

Be sure to refer to the system header files in the **/usr/include/sys** directory. The name of the file tells which structure and associated information it contains. For example, the user block is defined in **sys/user.h**. The process block is defined in **sys/proc.h**.

When you examine system data that maps into these structures, you can gain valuable kernel information that can explain why the dump was called.

Initiating a System Dump

A system dump initiated by a kernel panic is written to the primary dump device. If you initiate a system dump by pressing the reset button, with the key in the service position, the system dump is written to the primary dump device.

Use the special key sequences to determine whether the write of a system dump goes to the primary dump device or to the secondary dump device. The key must be in the service position. To write to the primary dump device, use the sequence Ctrl-Alt-NumPad1. To write to the secondary dump device, use the sequence Ctrl-Alt-NumPad2.

To use SMIT, select **Problem Determination** from the main menu, then select **System Dump**. This presents a menu that allows you to initiate a system dump to either the primary or secondary device, and manipulate the dump devices and the system dump files.

If you prefer to initiate the system dump from the command line, use the **sysdumpstart** command. Use the **-p** flag to write to the primary device or the **-s** flag to write to the secondary device.

If you want your device to be a primary or secondary device, the driver must contain a **dddump** routine. For more information, see the “dddump Entry Point” section in Chapter 4.

Note: The system halts after the system dump completes.

Including Device Driver Information in a System Dump

The system dump is table driven. The two parts of the table are:

master dump table

A master dump table entry is a pointer to a function which is provided by the device driver. The function is called by the kernel dump routine when a system dump occurs. The function must return a pointer to a component dump table.

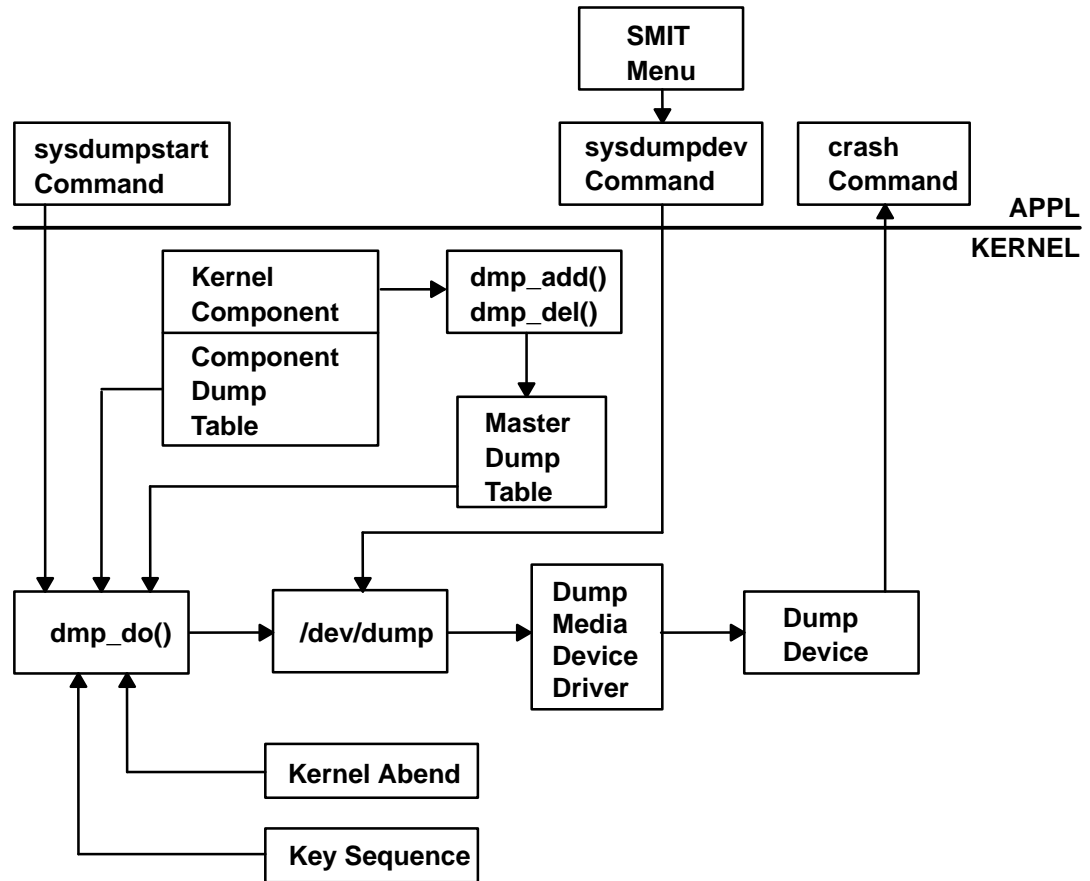
component dump table

Specifies memory areas to be included in a system dump.

Both the master dump table and the component dump table must reside in pinned global memory.

When a dump occurs, the kernel dump routine calls the function pointed to in the master dump table twice. On the first call, an argument of 1 indicates that the kernel dump routine is starting to dump the data specified by the component dump table.

On the second call, an argument of 2 indicates that the kernel dump routine has finished dumping the data specified by the component dump table. The component dump table should be allocated and pinned during initialization. The entries in the component dump table can be filled in later. The function pointed to in the master dump table must not attempt to allocate memory when it is called. The following System Dump Flow figure shows the flow of a system dump.



System Dump Flow

In order 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_add** kernel service to add an entry to the master dump table. Conversely, use the **dmp_del** kernel service to delete an entry from the master dump table. The syntax is as follows:

```
#include <sys/types.h>
#include <sys/errno.h>
#include <sys/dump.h>

int dmp_add(cdt_func) or int dmp_del(cdt_func)
int cdt * ((*cdt_func) ());
```

The **cdt** structure is defined in the **sys/dump.h** header file. A **cdt** structure consists of a fixed-length header (**cdt_head** structure) and an array of one or more **cdt_entry** structures.

The **cdt_head** structure contains a component name field, containing the name of the device driver, and the length of the component dump table. Each **cdt_entry** structure describes a contiguous data area, giving a pointer to the data area, its length, a segment register, and a name for the data area. Use the name supplied for the data area to refer to it when the **crash** command formats the dump. The following Kernel Dump Image figure illustrates a dump image.

Component Dump Table – A
Bitmap for 1st data area
1st data area for component A
Bitmap for 2nd data area
2nd data area for component A
...
Component Dump Table – N
Bitmap for 1st data area
1st data area for component N
Bitmap for 2nd data area
2nd data area for component N

Kernel Dump Image

Formatting a System Dump

Each device driver that includes data in a system dump can install a unique formatting routine in the **/usr/adm/ras/dmprtns** directory. A formatting routine is a command that is called by the **crash** command. The name of the formatting routine must match the component name field of the corresponding component dump table.

The **crash** command forks a child process that runs the formatting routines. If a formatting routine is not provided for a component name, the **crash** command runs the **_default_dmp_fmt** default formatting routine, which prints out the data areas in hex.

The **crash** command calls the formatting routine as a command, passing the file descriptor of the open dump image file as a command line argument. The syntax for this argument is **-file_descriptor**.

The dump image file includes a copy of each component dump table used to dump memory. Before calling a formatting routine, the **crash** command positions the file pointer for the dump image file to the beginning of the relevant component dump table copy.

The dumped memory is laid out in the dump image file with the component dump table and is followed by a bitmap for the first data area, then the first data area itself. A bitmap for the next data area follows, then the next data area itself, and so on.

The bitmap for a given data area indicates which pages of the data area are actually present in the dump image and which are not. Pages that were not in memory when the dump occurred were not dumped. The least significant bit of the first byte of the bitmap is set to 1 if the first page is present. The next least significant bit indicates the presence or absence of the second page, and so on. A macro for determining the size of a bitmap is provided in `sys/dump.h`.

The crash Command

The **crash** command is a particularly useful tool for device driver development and debugging, which interprets and formats the system structures. The **crash** command is interactive and allows you to examine an operating system image or an active system. An operating system image is held in a system dump file, either as a file or on the dump device. When you run the **crash** command, you can optionally specify a system image file and kernel file, as shown in the syntax below:

```
crash [-a] [-i IncludeFile] [ SystemImageFile [ KernelFile ] ]
```

The default *SystemImageFile* is **/dev/mem** and the default *KernelFile* is **/usr/lib/boot/unix**.

To run the **crash** command on the active system, enter:

```
crash
```

Because the command uses **/dev/mem**, you need root permissions.

To invoke the **crash** command on a system image file, enter:

```
crash SystemImageFile
```

where *SystemImageFile* is either a file name or the name of the dump device.

Note that by convention the symbol names for function entry points always begin with a . (period), while symbol names for data areas always begin with an _ (underscore). There is usually a data address corresponding to an external entry point address, and the **od** subcommand displays the data address for a name with no prefix. To be safe, use the proper prefix when looking for addresses.

Use the **-a** flag to generate a list of data structures without using subcommands. The resulting list is large, so you can redirect the output to either a file or to a printer.

Use the **-i** flag to read the given include file, allowing the **print** subcommand to output data according to the include file structures.

You can use a variety of subcommands to view the system structures. These subcommands can have flags that modify the format of the data. If you do not use a flag to specify what you want to see, all valid entries are displayed.

crash Subcommands

Once you initiate the **crash** command, **>** is the prompt character. For a list of the available subcommands, type the **?** character. To exit, type **q**. You can run any shell command from within the crash command by preceding it with an **!** (exclamation mark).

Since the **crash** command only deals with kernel threads, the word **thread** when used alone will be used to mean kernel thread in the **crash** documentation that follows. The default thread for several subcommands is the current thread (the thread currently running). On a multiprocessor system, you can use the **cpu** subcommand to change the current processor: the default thread becomes the running thread on the selected processor.

The parameters *ProcessTableEntry* and *ThreadTableEntry* are used in many subcommands to indicate a process or thread respectively. These parameters are simply numbers for table entry indexes which can be displayed using the **proc** and **thread** subcommands.

Note that many structures displayed are longer than one screen length. Make sure that you can halt scrolling if it is important to view something in detail. To do this, use the **stty** command:

```
stty ixon ixany
```

Use the **Ctrl-S** key sequence to stop scrolling and **Ctrl-Q** to resume scrolling.

buf [*BufferHeaderNumber*]

The **buf** subcommand displays the system buffer headers. A buffer header contains the information required to perform block I/O. If you type the **buf** subcommand with no *BufferHeaderNumber*, a summary of the system buffer headers is displayed.

Aliases = bufhdr, hdr

```
> buf
BUF MAJ MIN BLOCK FLAGS
 0 000a 000b 8 done stale
 1 000a 000b 243 done stale
 2 000a 000b 24 done stale
...
```

If you type the **buf** subcommand with a *BufferHeaderNumber* a single complete header is displayed:

```
> buf 3
BUFFER HEADER 3:
b_forw: 0x014d0528, b_back: 0x014d0160, b_vp: 0x00000000
av_forw: 0x014d0160, av_back: 0x014d0528, b_iodone: 0x000185f8
b_dev: 0x000a000b, b_blkno: 0, b_addr: 0x014e9000
b_bcount: 4096, b_error: 0, b_resid: 0
b_work: 0x80000000, b_options: 0x00000000, b_event: 0xffffffff
b_start.tv_sec: 0, b_start.tv_nsec: 0
b_xmemd.aspace_id: 0x00000000, b_xmemd.subspace_id: 0x00000000
b_flags: read done stale
```

Refer to the **sys/buf.h** header file for the structure definition.

buffer [*Format*] [*BufferHeaderNumber*]

The **buffer** subcommand displays the data in a system buffer according to the *Format* parameter. When specifying a buffer header number, the buffer associated with that buffer header is displayed. If you do not provide a *Format* parameter, the previous *Format* is used. Valid options are **decimal**, **octal**, **hex**, **character**, **byte**, **i-node**, **directory**, and **write**. The **write** option creates a file in the current directory containing the buffer data.

Aliases = b

```
> buffer hex 3

BUFFER FOR BUF_HDR 3
0000: 41495820 4c564342 00006a66 73000000
0020: 00000000 00000000 00000000 00000000
0040: 00000000 00000000 00003030 30303033
...
```

callout

The **callout** subcommand displays all active entries on the active **trblist**. When the **time-out** kernel extension is used in a device driver, this timer request is entered on a system-wide list of active timer requests. This list of timer requests is the **trblist**. Any timer which is active is on this list until it expires.

Aliases = c, call, calls, time, timeout, tout

```
>callout

TRB's On The Active List Of Processor 0.

TRB #1 on Active List
Timer address.....0x0
trb.to_next.....0x0
trb.knext.....0x59aa100
trb.kprev.....0x0
Thread id (-1 for dev drv).....0xffffffffe
Timer flags.....0x12
trb.timerid.....0x0
trb.eventlist.....0xfffffffff
trb.timeout.it_interval.tv_nsec...0x0
trb.timeout.it_interval.tv_sec....0x0
Next scheduled timeout (secs).....0x2d63f6a8
Next scheduled timeout (nanosecs)..0xc849a80
Timeout function.....0x8c748
Timeout function data.....0x59aa040
TRB #2 on Active List
...
```

Refer to **sys/timer.h** for the structure definitions, and to InfoExplorer for a description of the time-out mechanism.

cm [*ThreadT ableEntry SegmentNumber*]

The **cm** subcommand is used by the **od** subcommand to change the current segment map. The **cm** subcommand changes the map of the **crash** command internal pointers for any thread segment not paged out, if you specify the thread *ThreadT ableEntry* and *SegmentNumber*. This allows the **od** subcommand to display data relative to the beginning of the segment desired. The following example sets the map to thread *ThreadT ableEntry*3 to *SegmentNumber* 2, then displays ten words starting from the offset 0:

Aliases = none

```
> cm 3 2
t3,2 >> od 0 10
00000000: 00000000 00000000 00000000 00000000
00000010: 00000000 00000000 00000000 00000000
00000020: 00000000 00000000
t3,2 >>
...
```

Using the **cm** subcommand without any parameters resets the map of internal pointers.

cpu [*ProcessorNumber*]

If no argument is given, the **cpu** subcommand displays the number of the currently selected processor. Initially, the selected processor is processor 0. If the *ProcessorNumber* argument is given, the **cpu** subcommand selects the specified processor as the current processor. By extension, this selects the current kernel thread (the running kernel thread on the selected processor). Processor numbering starts from zero.

Aliases = none

```
> cpu
Selected cpu number : 0
```

dblock [*Address*]

The **dblock** subcommand displays the allocated streams data block headers. The address parameter is required. If the address is not supplied, this command will print an error message stating that the address is required. Refer to the **sys/stream.h** file for the datab structure definitions. The **freep** and **db_size** definitions are not included in **/usr/include/sys/stream.h**. These structure members are described here:

freep address of the free pointer

db_size size of the data block

There is no checking performed on the address passed in as the required parameter. The **dblock** subcommand will accept any address. It is up to the user to be sure that a valid address is specified.

To determine a valid address run the **mblock** subcommand. From the output of the **mblock** subcommand, select a non-zero data block address under the DATABLOCK column heading.

This subcommand can be issued from crash on either a running system or a system dump.

Aliases = dblk

```
> queue 59d5a74
  QUEUE      QINFO      NEXT  PRIVATE  FLAGS      HEAD  OTHERQ      COUNT
59d5a74 1884c1c 59d5474 59d5500 0x003e 59e1c00 59d5a00      4096
> mblk 59e1c00
ADDRESS     NEXT  PREVIOUS      CONT      RPTR      WPTR  DATABLOCK
59e1c00      0      0      0 59e2000 59e3000 59e1c44
> dblk 59e1c44
ADDRESS     FREEP      BASE      LIM  REFCNT      TYPE      SIZE
59e1c44      0 59e2000 59e3000      1      0      1000
```

dlock [*ThreadIdentifier* | **-p** [*ProcessorNumber*]

Displays deadlock analysis information about all types of locks (simple, complex, and lockl). The **dlock** subcommand searches for deadlocks from a given start point. If *ThreadIdentifier* is given, the corresponding kernel thread is the start point. If **-p** is given without a *ProcessorNumber*, the start point is the running kernel thread on the current processor. If **-p** *ProcessorNumber* is given, the running kernel thread on the specified processor is the start point. If no arguments are given, **dlock** searches for deadlocks among all threads on all processors.

The first output line gives information about the starting kernel thread, including the lock which is blocking the kernel thread, and a stack trace showing the function calls which led to the blocking lock request. Each subsequent line shows the lock held by the blocked kernel thread from the previous line, and identifies the kernel thread or interrupt handler which is blocked by those locks. If the information required for a full analysis is not available (paged out), an abbreviated display is shown; in this case, examine the stack trace to locate the locking operations which are causing the deadlock. The display stops when a lock is encountered for a second time, or no blocking lock is found for the current kernel thread.

Aliases = none

```

>dlock 00d3f
Deadlock from tid 00d3f. This tid waits for the first line lock,
owned by Owner-Id that waits for the next line lock, and so on...
      LOCK NAME | ADDRESS | OWNER-ID | WAITING FUNCTION
      lockC1 | 0x001f79e0 | Tid 113d | .lock_write_ppc
              called from : .times + 0000020c
Dump data incomplete.Only 0 bytes found out of 4.
              called from : .file + 0000000b
      lockC2 | 0x001f79e8 | Tid d3f | .lock_write_ppc
              called from : .times + 000001c8
Dump data incomplete.Only 0 bytes found out of 4.
              called from : .file + 0000000b

```

dmodsw

The **dmodsw** subcommand displays the streams drivers switch table. The information printed is contained in an internal structure. The following members of this internal structure are described here:

address	Address of dmodsw
d_next	Pointer to the next driver in the list
d_prev	Pointer to the previous driver in the list
d_name	Name of the driver
d_flags	Flags specified at configuration time
d_sqh	Pointer to synch queue for driver-level synchronization
d_str	Pointer to streamtab associated with the driver
d_sq_level	Synchronization level specified at configuration time
d_refcnt	Number of open or pushed count
d_major	Major number of a driver

The flags structure member, if set, is based one of the following values:

#define	Value	Description
F_MODSW_OLD_OPEN	0x1	Supports old-style (V.3) open/close parameters
F_MODSW_QSAFETY	0x2	Module requires safe timeout/bufcall callbacks
F_MODSW_MPSAFE	0x4	Non-MP-Safe drivers need funneling

The synchronization level codes are described in the `/usr/include/sys/strconf.h` file.

This subcommand can be issued from crash on either a running system or a system dump.

Aliases = none

```

> dmodsw
NAME          ADDRESS      NEXT PREVIOUS FLAG   SYNCHQ  STREAMTAB S-LVL  COUNT  MAJOR
sad           5a0cf40     5a0cf00   5a0c9c0 0x0   5a0ad40  188c600   3     0     12
slog         5a0cf00     5a0cec0   5a0cf40 0x0   5a0ad20  188c8a0   3     0     13
en           5a0cec0     5a0ce80   5a0cf00 0x0   5a0ad00  1893ee0   3     0     27
et           5a0ce80     5a0ce40   5a0cec0 0x0   5a0ace0  1893ee0   3     0     28
tr           5a0ce40     5a0ce00   5a0ce80 0x0   5a0acc0  1893ee0   3     0     29
fi           5a0ce00     5a0cdc0   5a0ce40 0x0   5a0aca0  1893ee0   3     0     30
echo         5a0cdc0     5a0cd80   5a0ce00 0x0     0     18951a0   5     0     31
nuls         5a0cd80     5a0cd40   5a0cdc0 0x0     0     1895190   5     0     32
spx          5a0cd40     5a0cd00   5a0cd80 0x0   5a0ac80  1895d70   3     0     33
unixdg       5a0cd00     5a0ccc0   5a0cd40 0x0   5a0ac60  18a14e0   3     0     34
unixst       5a0ccc0     5a0cc80   5a0cd00 0x0   5a0ac40  18a14e0   3     0     35
udp          5a0cc80     5a0cc40   5a0ccc0 0x0   5a0ac20  18a14e0   3     0     36
tcp          5a0cc40     5a0cb40   5a0cc80 0x0   5a0ac00  18a14e0   3     0     37
rs           5a0cb40     5a0cb00   5a0cc40 0x0     0     18b31d0   5     1     15
pts          5a0cb00     5a0ca40   5a0cb40 0x0     0     18fc930   4     7     24
ptc          5a0ca40     5a0ca00   5a0cb00 0x0     0     18fa5c0   4     2     23
tty          5a0ca00     5a0c9c0   5a0ca40 0x0     0     18fc950   4     0     26
ptyp        5a0c9c0     5a0cf40   5a0ca00 0x0     0     18fc940   4     0     25

```

ds [Address]

The **ds** subcommand returns the symbols closest to the given address. The **ds** subcommand can take either a text address or a data address.

Aliases = none

```
> ds 012345
      .ioctl_systrace + 0x000001b5
```

When a number such as *0x000001b5* is displayed, it shows the number of assembly language instructions by which the given address is offset from the entry point of the routine.

du [ThreadTableEntry]

Displays a combined hex and ASCII dump of the specified thread's uthread structure and of the user structure of the process which owns the thread. If the data is not available (paged out), a message is displayed. The default is the current thread.

Aliases = none

```
> du 3
Uthread structure of thread slot 3
00000000 00000000 00000000 2ff7fec0 00000000 *...../.....*
00000010 00000303 00000000 00030644 000010b0 *.....D....*
00000020 22222828 00030644 00006244 00000009 *"" (...D..bD....*
.
.
.
```

dump

The **dump** subcommand displays the name of each component for which there is data present. After you select a component name from the list, the **crash** program loads and runs the associated formatting routine contained in the **/usr/lib/ras/dmptns** directory.

If there is more than one data area for the selected component, the formatting routine displays a list of the data areas and allows you to select one. The **crash** command then displays the selected data area. You can enter the **quit** subcommand to return to the previously displayed list and make another selection or enter **quit** a second time to leave the **dump** subcommand loop.

Aliases = none

Displays messages in the error log. *Count* is the number of messages to print that have already been read by the **errdemon** process. (The default is 3 messages.) **errpt** always prints all messages that have not yet been read by the **errdemon** process.

Aliases = none

file [FileTableEntry]

The **file** subcommand displays the file table. Unless you request specific file entries, the command displays only those with a nonzero reference.

Aliases = files, f

```
> f 3
SLOT REF      INODE      FLAGS
   3   1 0x018e53f0   read
```

Refer to **sys/file.h** for the structure definition.

fmodsw

The **fmodsw** subcommand displays the streams modules switch table. The information printed is contained in an internal structure. The following members of this internal structure are described here:

address	Address of fmodsw
d_next	Pointer to the next module in the list
d_prev	Pointer to the previous module in the list
d_name	Name of the module
d_flags	Flags specified at configuration time
d_sqh	Pointer to synch queue for module-level synchronization
d_str	Pointer to streamtab associated with the module
d_sq_level	Synchronization level specified at configuration time
d_refcnt	Number of open or pushed count
d_major	-1

The flags structure member, if set, is based one of the following values:

#define	Value	Description
F_MODSW_OLD_OPEN	0x1	Supports old-style (V.3) open/close parameters
F_MODSW_QSAFETY	0x2	Module requires safe timeout/bufcall callbacks
F_MODSW_MPSAFE	0x4	Non-MP-Safe drivers need funneling

The synchronization level codes are described in the `/usr/include/sys/strconf.h` file.

This subcommand can be issued from crash on either a running system or a system dump.

Aliases = none

```
> fmodsw
NAME          ADDRESS      NEXT PREVIOUS FLAG  SYNCHQ  STREAMTAB S-LVL  COUNT MAJOR
bufcall       5a0cf80     5a0cc00  5a0ca80 0x1   5a0ad60  188bf80   3     0    -1
sc             5a0cc00     5a0cbc0  5a0cf80 0x0   5a0abe0  18a29b0   3     0    -1
timod         5a0cbc0     5a0cb80  5a0cc00 0x0   5a0abc0  18a34b0   3     0    -1
tirdwr        5a0cb80     5a0cac0  5a0cbc0 0x0   5a0aba0  18a4010   3     0    -1
ldterm        5a0cac0     5a0ca80  5a0cb80 0x0           0  18ef460   4     8    -1
tioc          5a0ca80     5a0cf80  5a0cac0 0x0           0  18f0e90   4    10    -1
```

fs [*ThreadT ableEntry*]

Traces a kernel stack for the thread specified by *ThreadT ableEntry*. Displays the called subroutines with a hex dump of the stack frame for the subroutine that contains the parameters passed to the subroutine. By default, the current thread is traced. This subcommand will not work on a running system because it uses stack tracing; however, it does work on a dump image.

Aliases = none

```
> fs
STACK TRACE:
**** .et_wait ****
2ff97e78 2FF97ED8 0080D568 00000000 018F4C60 /.^....h.....L`
2ff97e88 2FF97EE8 0080D568 00082BC0 000BA020 /.^....h..+.....
2ff97e98 2FF97ED8 28008044 00082418 2FF98000 /.^.(.D..B./...
2ff97ea8 00000000 000B8468 00000000 00000000 .....h.....
2ff97eb8 2FF97F38 0000000B 00000004 00000004 /.8.....
2ff97ec8 00000005 01DFE258 00000000 E3000600 .....X.....
```

inode [-] [*<Major>* *<Minor>* *<INumber>*]

The **inode** subcommand displays the i-node table and the i-node data block addresses. You can display a specific i-node by specifying the major and minor device numbers of the

device where the i-node resides and the i-node number. The command displays the i-node only if it is currently on the system hash list.

Aliases = ino, i

```
>inode
  ADDRESS    MAJ MIN  INUMB REF LINK  UID  GID  SIZE    MODE  SMAJ SMIN FL
AGS
0x018e4e50  010 0007 11264  0   1   2   2    30 ----777  -  -
0x018f9fd0  010 0009 16384  1   6  201  0    512 d---755  -  -
      addr: 16448  0  0  0  0  0  0  0
0x018ea940  010 0011  0   1   0   0  0    0 ----  0  -  -
...
```

kfp [*FramePointer*]

If you use the **kfp** subcommand without parameters, it displays the last kernel frame pointer address that was set using **kfp**. If you specify a frame pointer address, it sets the kernel frame pointer to the new address. Use this subcommand in conjunction with the **-r** flag of the **trace** subcommand.

Aliases = fp, rl

```
> kfp
```

knlist [*Symbol*]

The **knlist** subcommand displays the addresses of all the specified symbol names. If there is no such symbol, the subcommand displays a `no match` message. Run this subcommand only on an active system.

The **knlist** subcommand runs a subroutine to the active kernel to obtain the address from the system's knlist. The **nm** subcommand provides the same function but searches the symbol table in the Kernel Image File for the address and therefore can be used on a dump.

Aliases = none

```
> knlist open
      open:0x000bbc98
```

le [*Module*] *Address*

The **le** subcommand displays the kernel load list entries. If you specify an address in a kernel extension, the corresponding load list entry is displayed. If you attempt to display a paged out loader entry, the subcommand displays an error message.

Aliases = none

linkblk

The **linkblk** subcommand displays the streams linkblk table. Refer to the **/usr/include/sys/stream.h** file for the `linkblk` structure definitions. If there are no `linkblk` structures found on the system, the **linkblk** subcommand will print a message stating that no structures are found.

This subcommand can be issued from crash on either a running system or a system dump.

Aliases = lblk

This example shows a regular link:

```
> linkblk
      QTOP      QBOT      INDEX
5ab8b74  5ae5074  5ab4200
```

This example shows a persistent link:

```
> linkblk
      QTOP      QBOT      INDEX
      0 5ae5174 5a4ef00
```

mblock *Address*

The **mblock** subcommand displays the allocated streams message block headers. The address parameter is required. If the address is not supplied, this command will print an error message stating that the address is required. Refer to the `/usr/include/sys/stream.h` file for the queue structure definitions.

The **mblock** subcommand's checking of the address parameter is limited to verifying that the address falls on a 128-byte boundary. It is up to the user to be sure that a valid address is specified.

To determine a valid address, run the queue subcommand. From the output of the queue subcommand, select a non-zero address for the head of the message queue under the HEAD column heading for either a read queue or a write queue.

This subcommand can be issued from crash on either a running system or a system dump.

Aliases = mblk

```
> queue 59d5a74
      QUEUE      QINFO      NEXT  PRIVATE  FLAGS      HEAD  OTHERQ      COUNT
59d5a74 1884c1c 59d5474 59d5500 0x003e 59e1c00 59d5a00      4096
> mblk 59e1c00
      ADDRESS      NEXT  PREVIOUS      CONT      RPTR      WPTR  DATABLOCK
59e1c00           0           0           0 59e2000 59e3000 59e1c44
```

mbuf [-] [*Clusters* | *Address*].

The **mbuf** subcommand displays the system **mbuf** structures. These structures are memory buffers that are chained and can be manipulated by the Memory Buffer kernel services. If you specify the `-` flag, the subcommand also displays the data associated with the **mbuf** structures.

The **mbuf** subcommand with no additional arguments displays the chain of **mbuf** structures pointed to by the **mbuf** pointer. If you specify *Clusters*, the subcommand displays the chain of **mbuf** structures pointed to by the kernel `mbclusters` pointer. If you specify *Address*, then the **mbuf** structure at the given address is displayed. Note that valid **mbuf** pointers must be on a 128-byte boundary.

Aliases = mbuf

```
> mbuf
      ADDRESS      SIZE      TYPE      LINK      DATAPTR
0x01a67000      0      free 0x00000000 0x01a67000
      DATA: 0x00000000 0x00000000 0x00000000 0x00000000
```

Refer to the `sys/mbuf.h` header file for the structure definition.

mst [*Address*]

The **mst** subcommand displays the `mstsave` portion of the `uthread` structures at the addresses specified. If you do not specify an address, the subcommand displays information about the last running kernel thread.

Aliases = none

ndb

Displays network kernel data structures either for a running system or a system dump. The **ndb** subcommand, short for network debugger, displays the following options:

?	Provides first-level help information.
help	Provides additional help information.
tcb [<i>Addr</i>]	Shows TCBS. The default is HEAD TCB.
udb [<i>Addr</i>]	Shows UDBs. The default is HEAD UDB.
sockets	Shows sockets from the file table.
mbuf [<i>Addr</i>]	Shows the mbuf at the specified address.
ifnet [<i>Addr</i>]	Shows the ifnet structures at the specified address.
quit	Stops the running option.
xit	Exits the ndb submenu.

Aliases = none

nm [*Symbol*]

The **nm** subcommand displays the symbol value and type found in *KernelFile*

Aliases = none

```
> nm open
00095484 000C70 PR SD <.open>
00095484 PR LD .open
000BBC98 00000C SV SD open
```

od [*SymbolName* | *Address*][*Count*] [*Format*]

The **od** subcommand dumps *Count* number of data values starting at *Symbol* value or *Address* according to *Format*. Possible formats are **octal**, **longoct**, **decimal**, **longdec**, **character**, **hex**, **instruction**, and **byte**. The default is **hex**. Note that if you use the *Format* parameter, you must also use *Count*.

The **od** subcommand is especially useful during program development in order to see structure values at a given point in time.

Aliases = none

```
> od open 10
00095484: 7c0802a6 bf21ffe4 90010008 9421ff30
00095494: 609c0000 832202e0 607b0000 60bd0000
000954a4: 63230000 38800000

> od open 10 byte
00095484: 0174 0010 0002 0246 0277 0041 0377 0344
0009548c: 0220 0001

> od 12345
warning: word alignment performed
00012344: 480001d8
```

ppd [*ProcessorNumber* | *]

Displays per-processor data area (PPDA) structures for the specified processor. If no processor is specified, the current processor selected by the **cpu** subcommand is used. If the asterisk argument is given, the PPDA of every enabled processor is displayed.

Aliases = none

```

> ppd
Per Processor Data Area for processor 0
csa.....2fedf500
mstack.....00315db0
fpowner.....e6001360
curthread.....e6001360
r0.....60000000
r1.....6000068e
r2.....d00089b8
r15.....0000f930
sr0.....d0005a54
sr2.....2feac6e0
iar.....f013c7c4

```

print type Address

Does **dbx**-style printing of structures. The **-i** option must be given on the command line to use this feature.

Aliases = none

proc [-] [-r] [ProcessT ableEntry]

The **proc** subcommand displays the process table, including the kernel thread count (the number of threads in the process) and state of each process. Use the **-r** flag to display only runnable processes. Use the **-** flag to display a longer listing of the process table.

Aliases = **ps, p**

```
>p
```

```

SLT ST   PID   PPID   PGRP   UID   EUID   TCNT   NAME
  0 a     0     0     0     0     0     1   swapper
      FLAGS: swapped_in no_swap fixed_pri kproc
  1 a     1     0     0     0     0     1   init
      FLAGS: swapped_in no_swap
  2 a    204     0     0     0     0     1   wait
      FLAGS: swapped_in no_swap fixed_pri kproc
...

```

```
>p 20
```

```

SLT ST   PID   PPID   PGRP   UID   EUID   TCNT   NAME
 20 a   1406     1   1406     0     0     1   ksh
      FLAGS: swapped_in no_swap

```

```
>p -
```

```

SLT ST   PID   PPID   PGRP   UID   EUID   TCNT   NAME
  0 a     0     0     0     0     0     1   swapper
      FLAGS: swapped_in no_swap fixed_pri kproc

```

```

Links:  *child:0xe3000100  *siblings:0x00000000  *uid1:0xe3001400
        *ganchor:0x00000000
Dispatch Fields:  pevent:0x00000020  wevent:0x00000000
                *p_synch:0xffffffff
Thread Fields:  *threadlist:0xe6000000  threadcount: 1
                active: 1  suspended: 0  local: 0  localsleep: 0
                *synch:0xffffffff

```

```

Scheduler Fields:  fixed pri: 16  repage:0x00000000
scount:0x00000000
Misc:  adspace:0x00000808  *ttyl:0x00000000
      *p_ipc:0x00000000  *p_dblist:0x00000000
*p_dbnext:0x00000000
      *lock:0x00000000  kstackseg:0x007fffff  *pgrp1:0x08x

Signal Information:
  pending:hi 0x00000000,lo 0x00000000
  sigcatch:hi 0x00000000,lo 0x00000000  sigignore:hi
0xffffffff,lo 0xfff7ffff
Statistics:  size:0x00000000(pages)  audit:0x00000000

SLT ST  PID  PPID  PGRP  UID  EUID  TCNT  NAME
  1 a   1    0    0    0    0    1  init
      FLAGS: swapped_in no_swap

  Links:  *child:0xe3001400  *siblings:0x00000000
*uidl:0xe3000100
      *ganchor:0x00000000
Dispatch Fields:  pevent:0x00000020  wevent:0x00000000
      *p_synch:0xffffffff
Thread Fields:  *threadlist:0xe60000a0  threadcount: 1
  active: 1  suspended: 0  local: 0  localsleep: 0
      *synch:0xffffffff
Scheduler Fields:  nice: 20  repage:0x00000000
scount:0x00000000
Misc:  adspace:0x00000505  *ttyl:0x00000000
      *p_ipc:0x00000000  *p_dblist:0x00000000
*p_dbnext:0x00000000
      *lock:0x00000000  kstackseg:0x007fffff  *pgrp1:0x08x
Signal Information:
  pending:hi 0x00000000,lo 0x00000000
  sigcatch:hi 0x00000001,lo 0x18783eff  sigignore:hi
0xfffffffef,lo 0xe787c100
Statistics:  size:0x00000028(pages)  audit:0x00000000
...

```

Refer to the **sys/proc.h** header file for the structure definition.

queue [Address]

The **queue** subcommand displays the STREAMS queue. If the address optional parameter is not supplied, crash will display information for all queues available. Refer to the **/usr/include/sys/stream.h** file for the `queue` structure definitions.

If you wish to see the information stored for a read queue, issue the **queue** subcommand with the read queue address specified as the parameter.

When you issue the **queue** subcommand with the address parameter, the column headings do not distinguish between the read queue and the write queue. One queue address will be displayed under the column heading `QUEUE`. The other queue in the pair will be displayed under the column heading `OTHERQ`. The write queue will have a numerically higher address than the read queue.

This subcommand can be issued from crash on either a running system or a system dump.

Aliases = que

```

> queue
WRITEQ      QINFO      NEXT  PRIVATE  FLAGS      HEAD      READQ      COUNT
59c2a74 188c50c 59c2474 59b1900 0x002a      0 59c2a00      0
59c2474 18f0e50 59c2274 59b6880 0x0028      0 59c2400      0
59c2274 18ef3d8 59c2174 59cc800 0x0028      0 59c2200      0
59c2174 18b31b4      0 54684f8 0x0028      0 59c2100      0
59d5a74 188c50c 5a94874 59d3c00 0x002a      0 59d5a00      0
5a94874 18f0e50 5a9c074 59ec500 0x0028      0 5a94800      0
5a9c074 18fa748      0 59e5b00 0x0028      0 5a9c000      0
59ff074 188c50c 59eab74 59ffe00 0x002a      0 59ff000      0
59eab74 18f0e50 59ff374 59da500 0x0028      0 59eab00      0
59ff374 18fa748      0 59da380 0x0028      0 59ff300      0
5ab4374 188c50c 59ee174 5ab4c00 0x002a      0 5ab4300      0
59ee174 18f0e50 5ab4774 59da100 0x0028      0 59ee100      0
5ab4774 18ef3d8 5ad2874 59d7800 0x0028      0 5ab4700      0

> queue 5ab4700
QUEUE      QINFO      NEXT  PRIVATE  FLAGS      HEAD      OTHERQ      COUNT
5ab4700 18ef3bc 59ee100 59d7800 0x0029      0 5ab4774      0

```

quit

Exit from the **crash** command.

Aliases = q

qrun

The **qrun** subcommand displays the list of scheduled streams queues. If there are no queues found for scheduling, the **qrun** subcommand will print a message stating there are no queues scheduled for service.

This subcommand can be issued from crash on either a running system or a system dump.

Aliases = none

```

> qrun
QUEUE
59d5a74

```

socket [-]

The **socket** subcommand displays the system socket structures. Use the **-** flag to also display the socket buffers.

Aliases = sock

```

> sock
SLOT: 26 type:0x0002 opts:0x0000 linger:0x0000
state:0x0080 pcb:0x01d32d8c proto:0x01c65cf0
q0:0x00000000 qlen: 0 q:0x00000000
qlen: 0 qlimit: 0 head:0x00000000
timeo: 0 error: 0 oobmark: 0 pgrp: 0
...

```

Refer to **sys/socket.h** for structure definitions.

stack [*ThreadT ableEntry*]

The **stack** subcommand displays a dump of the kernel stack of the kernel thread identified by *ThreadT ableEntry*. The addresses are virtual data addresses rather than true physical addresses. If you do not specify an entry, the subcommand displays information about the last running kernel thread. You cannot trace the stack of the current running kernel thread on a running system.

Aliases = s, stk, k, kernel

```

> s 31
KERNEL STACK:
2ff97a50:          8eaa4          16 2ff97ac8          2
2ff97a60:          90b0          8e8b4 2ff97ad8          0
2ff97a70:           1          26 2ff97ac8 2ff98938
...

```

stat

The **stat** subcommand displays statistics found in the dump. These statistics include the panic message (if there is one), time of crash, and system name.

Aliases = none

```

> stat
      sysname: AIX
      nodename: funk
      release: 1
      version: 3
      machine: 000003961000
      time of crash: Fri Sep 28 17:50:38 1990
      age of system: 15 day, 6 hr., 25 min.

```

status [*ProcessorNumber*]

Displays a description of the kernel thread scheduled on the designated processor. If no processor is specified, the **status** subcommand displays information for all processors. The information displayed includes the processor number, kernel thread identifier, kernel thread table slot, process identifier, process table slot, and process name.

Aliases = none

```

> status 0
CPU      TID  TSLOT      PID  PSLOT  PROC_NAME
  0     1fe1     31     1fd8     31     crash

```

stream

The **stream** subcommand displays the stream head table. The information printed is contained in an internal structure. The following members of this internal structure are described here:

address	address of stream head
wq	address of streams write queue
dev	associated device number of the stream
read_error	read error on the stream
write_error	write error on the stream
flags	stream head flag values
push_cnt	number of modules pushed on the stream
wroff	write offset to prepend M_DATA
ioc_id	id of outstanding M_IOCTL request
pollq	list of active polls
sigsq	list of active M_SETSIGs

The flags structure member, if set, is based on combinations of the following values:

#define	Value	Description
F_STH_READ_ERROR	0x0001	M_ERROR with read error received, fail all read calls
F_STH_WRITE_ERROR	0x0002	M_ERROR with write error received, fail all writes
F_STH_HANGUP	0x0004	M_HANGUP received, no more data
F_STH_NDELON	0x0008	Do TTY semantics for ONDELAY handling
F_STH_ISATTY	0x0010	This stream acts a terminal
F_STH_MREADON	0x0020	Generate M_READ messages
F_STH_TOSTOP	0x0040	Disallow background writes (for job control)
F_STH_PIPE	0x0080	Stream is one end of a pipe or FIFO
F_STH_WPIPE	0x0100	Stream is the "write" side of a pipe
F_STH_FIFO	0x0200	Stream is a FIFO
F_STH_LINKED	0x0400	Stream has one or more lower streams linked
F_STH_CTTY	0x0800	Stream controlling tty
F_STH_CLOSED	0x4000	Stream has been closed, and should be freed
F_STH_CLOSING	0x8000	Actively on the way down

This subcommand can be issued from **crash** on either a running system or a system dump.

Aliases = none

```
> stream
ADDRESS  WRITEQ MAJ/MIN RERR  WERR  FLAGS IOCID  WOFF PCNT POLQNEXT SIGQNEXT
59b1900  59c2a74  15,  0    0    0 0x0838    0    0    2        0        0
59d3c00  59d5a74  23,  5    0    0 0x0020    0    0    1        0        0
59ffe00  59ff074  23,  4    0    0 0x0020    0    0    1        0        0
5ab4c00  5ab4374  24,  0    0    5 0x0816    0    0    2        0        0
59d3f00  59ee974  24,  1    0    5 0x0816    0    0    2        0        0
59d3800  59dff74  24,  2    0    5 0x0816    0    0    2        0        0
59d3700  5a9c174  24,  3    0    5 0x0816    0    0    2        0        0
59ff800  59ff774  24,  4    0    0 0x0810    0    0    2        0        0
5a94d00  59ee574  24,  5    0    0 0x0830    0    0    2        0        0
5a94600  5a96c74  24,  6    0    5 0x0816    0    0    2        0        0
```

tcb [*ThreadT ableEntry*]. . .

Displays the mtsave portion of the user structures of the named kernel threads (see the **user.h** and **mtssave.h** header files). If you do not specify an entry, information about the last running kernel thread is displayed. This subcommand replaces the **pcb** subcommand.

Aliases = none

```
> tcb
      UTHREAD AREA FOR SLOT    2
SAVED MACHINE STATE
  curid:0x00000204  m/q:0x00000000  iar:0x00019cfc  cr:0x22000000
  usr:0x00009030  lr:0x00035678  ctr:0x00019c90  xer:0x20000000
  *prevmst:0x00000000  *stackfix:0x00000000  intpri:0x0000000b
  backtrace:0x00  tid:0x00000000  fpeu:0x00  ecr:0x00000000
Exception Struct
  0x00000000  0x00000000  0x00000000  0x00000000  0x00000000
Segment Regs
  0:0x00000000  1:0x007ffffff  2:0x00000408  3:0x007ffffff
  4:0x007ffffff  5:0x007ffffff  6:0x007ffffff  7:0x007ffffff
.
.
.
```

Aliases = none


```

>tcb

        UTHREAD AREA FOR SLOT 25

SAVED MACHINE STATE
  curid:0x0000162e  m/q:0x00000000  iar:0x0187e1dc  cr:0x44224820
  msr:0x000090b0  lr:0x0187e2d8  ctr:0x00040610  xer:0x00000004
  *prevmst:0x00000000  *stackfix:0x00000000  intpri:0x0000000b
  backtrace:0x00  tid:0x00000000  fpeu:0x01  ecr:0x00000000
Exception Struct
  0xd0112540  0x42000000  0x400015b5  0xd0112540  0x00000106
Segment Regs
  0:0x00000000  1:0x007fffff  2:0x000019f9  3:0x007fffff
  ...
General Purpose Regs
  0:0x018abcd8  1:0x2fedefb8  2:0x018ac41c  3:0x018ab6e4
  4:0x018a63b0  5:0x018a63b0  6:0x018a63b0  7:0x00000000
  ...
Floating Point Regs
  Fpscr: 0x00000000
  0:0x00000000  0x00000000  1:0x00000000  0x00000000  2:0x00000000  0x00000000
  3:0x00000000  0x00000000  4:0x00000000  0x00000000  5:0x00000000  0x00000000
  ...

Kernel stack address: 0x2fedf500

```

```

>tcb 10

        UTHREAD AREA FOR SLOT 10

SAVED MACHINE STATE
  curid:0x000009f4  m/q:0x00008003  iar:0x0001ddf8  cr:0x80222822
  msr:0x000010b0  lr:0x0001ddf8  ctr:0x000ee000  xer:0x00000000
  *prevmst:0x00000000  *stackfix:0x2fedf2d8  intpri:0x00000000
  backtrace:0x00  tid:0x00000000  fpeu:0x01  ecr:0x00000000
Exception Struct
  0x30000000  0x40000000  0x00001272  0x30000000  0x00000106
Segment Regs
  0:0x00000000  1:0x007fffff  2:0x00001010  3:0x007fffff
  ...
General Purpose Regs
  0:0x40000707  1:0x2fedf2d8  2:0x00160d2c  3:0x00000420
  4:0x00000001  5:0xe6000644  6:0x000010b0  7:0x00000420
  ...
Floating Point Regs
  Fpscr: 0x00000000
  0:0x00000000  0x00000000  1:0x00000000  0x00000000  2:0x00000000  0x00000000
  3:0x00000000  0x00000000  4:0x00000000  0x00000000  5:0x00000000  0x00000000
  ...

Kernel stack address: 0x2fedf500

```

thread [-] [-r] [-p *ProcessT ableEntry*] [-a *Address* | *ThreadT ableEntry*]

Displays the contents of the kernel thread table. The - (minus) flag displays a longer listing of the thread table. The -r flag displays only runnable kernel threads. The -p flag displays only those kernel threads which belong to the process identified by *ProcessT ableEntry*. The -a flag displays the kernel thread structure at *Address*. If *ThreadT ableEntry* is given, only the corresponding kernel thread is displayed.

Aliases = th

```
> thread 1
SLT ST      TID      PID      CPUID    POLICY  PRI  CPU    EVENT  PROCNAME
  1  s      1e1      1  unbound  other  3c    0      init
      FLAGS: wakeonsig
```

The **trace** subcommand displays a kernel stack trace of the kernel thread identified by *ThreadT ableEntry*. The trace starts at the bottom of the stack and attempts to find valid stack frames deeper in the stack. By default, the current kernel thread is used.

Use the **-r** flag to use the kernel frame pointer set up by the **kfp** subcommand as the starting address instead of the frame pointer found in the *SystemImageFile*. The **trace** subcommand stops and reports an error if an invalid frame pointer is encountered.

Aliases = t

```
> t 31
STACK TRACE:
    .et_wait ()
    .e_sleep ()
    .e_sleep1 ()
    .sleepx ()
    .fifo_read ()
    .fifo_rdwr ()
    .vno_rw ()
    .rwuio ()
    .rdwr ()
    .kreadv ()
```

ts [*TextAddress*]

The **ts** subcommand finds the text symbols closest to the given address.

Aliases = none

```
> ts 012345
    .ioctl_systrace
```

tty [**d**] [**l**] [**e**] [*Name* | *Major* [*Minor*]]

Aliases = term, dz, dh

Refer to the **sys/tty.h** header file for the structure definition.

user [*ThreadT ableEntry*]

Displays the *uthread* structure and the associated user structure of the thread identified by *ThreadT ableEntry*. If you do not specify the entry, the information about the last running kernel thread is displayed.

Aliases = u, uarea, u_area

```
>u 4

      UTHREAD AREA FOR SLOT 4

      SAVED MACHINE STATE
      curid:0x00000408  m/q:0x00008003  iar:0x0001ed98  cr:0x84201000
      msr:0x000010b0  lr:0x0001ed98  ctr:0x00000000  xer:0x20000000
      *prevmst:0x00000000  *stackfix:0x2feaeaa8  intpri:0x00000000
      backtrace:0x00  tid:0x00000000  fpeu:0x00  ecr:0x00000000
      Exception Struct
      0x2feaf688  0x40000000  0x00000c0c  0x2feaf688  0x00000106
```

```

Segment Regs
 0:0x00000000  1:0x007fffff  2:0x00000c0c  3:0x007fffff
...
General Purpose Regs
 0:0x00000000  1:0x2feaeaa8  2:0x00160d2c  3:0x00001000
 4:0x00000001  5:0x2fedf500  6:0x0000000b  7:0x000090b0
...
Floating Point Regs
  Fpscr: 0x00000000
  0:0x00000000 0x00000000  1:0x00000000 0x00000000
...
 30:0x00000000 0x00000000 31:0x00000000 0x00000000

Kernel stack address: 0x2feaeffc

SYSTEM CALL STATE
  errno address:0xc0c0fade  error code:0x00  *kjmpbuf:0x00000000

PER-THREAD TIMER MANAGEMENT
  Real/Alarm Timer (ut_timer.t_trb[TIMERID_ALARM]) = 0x0
  Virtual Timer (ut_timer.t_trb[TIMERID_VIRTUAL]) = 0x0
  Prof Timer (ut_timer.t_trb[TIMERID_PROF]) = 0x0

SIGNAL MANAGEMENT
  *sigsp:0x0  oldmask:hi 0x0,lo 0x0  code:0x0

MISCELLANOUS FIELDS:
  fstid:0x00000000  ioctlrsv:0x00000000  selchn:0x00000000

  USER AREA OF ASSOCIATED PROCESS gil (SLOT 4, PROCTAB 0xe3000400)
  handy_lock:0x00000000  timer_lock:0x00000000
  map:0x00000000  *semundo:0x00000000
  compatibility:0x00000000  lock:0x00000000

SIGNAL MANAGEMENT
  Signals to be blocked (sig#:hi/lo mask,flags,&func)
  1:hi 0x00000000,lo 0x00000000,0x00000000,0x00000000
  2:hi 0x00000000,lo 0x00000000,0x00000000,0x00000000
  3:hi 0x00000000,lo 0x00000000,0x00000000,0x00000000
  ...

USER INFORMATION
  euid:0x0000  egid:0x0000  ruid:0x0000  rgid:0x0000  luid:0x00000000
  suid:0x00000000  ngrps:0x0000  *groups:0x2feacc34  compat:0x00000000
  ref:0x00000004
  acctid:0x00000000  sgid:0x00000000  epriv:0x00000000
  ipriv:0x00000000  bpriv:0x00000000  mpriv:0x00000000

  u_info:

ACCOUNTING DATA
  start:0x2d612cc9  ticks:0x00000002  acflag:0x0000  pr_base:0x00000000
  pr_size:0x00000000  pr_off:0x00000000  pr_scale:0x00000000
  process times:
    user:0x00000000s 0x00000000us
    sys:0x000004f1s 0x14dc9380us
  children's times:
    user:0x00000000s 0x00000000us
    sys:0x00000000s 0x00000000us

```

CONTROLLING TTY

*ttysid:0x00000000 *tty(pgrp):0x00000000
ttyd(evice):0x00000000 ttympx:0x00000000 *ttys(tate):0x00000000
tty id: 0x00000000 *query function: 0x00000000

PINNED PROFILING BUFFER

*pprof: 0x00000000 *mem desc: 0x00000000

RESOURCE LIMITS AND COUNTERS

ior:0x00000000 iow:0x00000000 ioch:0x00000000
text:0x00000000 data:0x00000000 stk:0x01000000
max data:0x08000000 max stk:0x01000000 max file:0x7fffffff
soft core dump:0x7fffffff hard core dump:0x7fffffff
soft rss:0x7fffffff hard rss:0x7fffffff
cpu soft:0x7fffffff cpu hard:0x7fffffff
hard ulimit:0x7fffffff
minflt:0x00000000 majflt:0x00000000

AUDITING INFORMATION

auditstatus:0x00000000

SEGMENT REGISTER INFORMATION

Reg	Flag	Fileno	Pointer
0	0	0	0
1	0	0	0
...			

*adspace:0xa0000000

FILE SYSTEM STATE

*curdir:0x00000000 *rootdir:0x00000000
cmask:0x0000 maxindex:0x0000

FILE DESCRIPTOR TABLE

*ufd: 0x20013b14

> user

UTHREAD AREA FOR SLOT 31

SAVED MACHINE STATE

curid:0x00001fd8 m/q:0x00000000 iar:0x0006ee54 cr:0x2224248a
msr:0x00009030 lr:0x000095d4 ctr:0x00000009 xer:0x00000020
*prevmst:0x00000000 *stackfix:0x00000000 intpri:0x0000000b
backtrace:0x00 tid:0x00000000 fpou:0x01 ecr:0x00000000
Exception Struct
0x10013b7c 0x4000d030 0x60000990 0x10013b7c 0x00000106
Segment Regs
0:0x00000000 1:0x007ffffff 2:0x0000068e 3:0x6000068e
4:0x007ffffff 5:0x007ffffff 6:0x007ffffff 7:0x007ffffff

.
. .

2l_trb[TIMERID_ALARM]) = 0x0
Virtual Timer (ut_timer.t_trb[TIMERID_VIRTUAL]) = 0x0
Prof Timer (ut_timer.t_trb[TIMERID_PROF]) = 0x0

SIGNAL MANAGEMENT

*sigsp:0x0 oldmask:hi 0x0,lo 0x0 code:0x0

MISCELLANEOUS FIELDS:

fstid:0x00000000 ioctlr:0x00000000 selchn:0x00000000

USER AREA OF ASSOCIATED PROCESS crash (SLOT 31, PROCTAB
0xe3001f00)

handy_lock:0x00000000 timer_lock:0x00000000

```

    map:0x00000000 *semundo:0x00000000
    compatibility:0x00000000 lock:0x00000000
SIGNAL MANAGEMENT
    Signals to be blocked (sig#:hi/lo mask,flags,&func)
    1:hi 0x00000000,lo 0x00000000,0x00000000,0x00000000
    2:hi 0x00x00000000
.
.
.

```

Refer to the **sys/user.h** header file for the structure definition.

var

The **var** subcommand displays the tunable system parameters.

Aliases = tune, tunable, tunables

```

> var
buffers          20
files            328
e_files          328
threads         262144
e_threads        51
clists           16384
maxproc          40
iostats          1
locks            200
e_locks         8456344

```

vfs [-] [Vfs SlotNumber]

The **vfs** uses the specified *Vfs SlotNumber* to display an entry in the **vfs** table. Use the **-** flag to display the vnodes associated with the **vfs**. The default displays the entire **vfs** table.

Aliases = m, mnt, mount

```

> vfs 3
VFS ADDRESS TYPE OBJECT STUB NUM FLAGS PATHS
  3 1a62494 jfs 1a6d47c 1a6d650 5 D /dev/hd1 mounted over /u
      flags: C=disconnected D=device I=remote P=removable
             R=readonly S=shutdown U=unmounted Y=dummy

> vfs - 3
VFS ADDRESS TYPE OBJECT STUB NUM FLAGS PATHS
  3 1a62494 jfs 1a6d47c 1a6d650 5 D /dev/hd1 mounted over /u
ADDRESS VFS MVFS VNTYPE FSTYPE COUNT ISLOT INODE FLAGS
1a6e0ac 3 - vreg jfs 1 - 18f82c0
1a6e218 3 - vreg jfs 1 - 18f8770
1a6e24c 3 - vreg jfs 1 - 18f8590
1a6e17c 3 - vdir jfs 3 - 18f7f00
1a6dea4 3 - vreg jfs 2 - 18f65b0
1a6dfa8 3 - vdir jfs 5 - 18f6100
1a6d47c 3 - vdir jfs 1 - 18ea580 vfs_root

```

Refer to the **sys/vfs.h** header file for structure definitions.

vnode [VNodeAddress]

The **vnode** subcommand displays data at the specified *VNodeAddress* as a **vnode**. *VNodeAddress* must be specified in hexadecimal notation. The default is to display all **vnodes** in the **vnode** table.

Aliases = none

```
> vnode 1a6e078
ADDRESS VFS MVFS VNTYPE FSTYPE COUNT ISLOT DATAPTR FLAGS
1a6e078  0    -   vreg   jfs     4    -   18f6790
  Total VNODES printed 1
```

Refer to the **sys/vnode.h** header file for the structure definition.

xmalloc

The **xmalloc** subcommand displays information concerning the allocation and usage of kernel memory, specifically the **pinned_heap** and the **kernel_heap**.

Aliases = xm, malloc

```
>xmalloc

Kernel heap usage
heap size = 242720768 amount used = 79005568
Pinned heap usage
heap size = 242720768 amount used =   342832
Kernel and pinned heap usage
from =  1028ac bytes = 62914560 number =      2
from =   41d58 bytes =  8388608 number =      1
...
```

Low-Level Kernel Debugger

Use the kernel debug program (also known as the low-level debugger) for debugging the kernel and kernel extensions. The kernel debug program provides the following functions:

- Setting breakpoints within the kernel.
- Formatting and displaying selected kernel data structures.
- Viewing and modifying memory for any kernel data.
- Viewing and modifying memory for kernel instructions.
- Modifying the state of the machine by altering system registers.

Entering the Debugger

You can only display the kernel debugger on an ASCII terminal connected to a serial port. The kernel debugger does *not* support any displays connected to any graphics adapters. The debugger has its own device driver for handling the display terminal. It is also possible to connect a serial line between two machines and define the serial line port as the port for the console. In that case, use the **cu** command to connect to the target machine and run the debugger.

It is possible to enter the system debugger through one of the following procedures:

- From a native keyboard, press Ctrl-Alt-Numpad4. The key must be in service position to enter the kernel debug program from the keyboard.
- From the tty keyboard, enter Ctrl-4 (IBM 3151 terminals) or Ctrl-\ (BQ 303, BQ 310C, and WYSE 50). The key must be in service position to enter the kernel debug program from the keyboard.
- The system can enter the debugger if a breakpoint is set. To do this, use the **break** subcommand. See “Setting Breakpoints” on page 14-45 for information on setting a breakpoint.
- The system can also enter the debugger by calling the **brkpoint** subroutine from C code. The syntax for calling this subroutine is:

```
brkpoint();
```

- The system can also enter the debugger if a static debug trap (SDT), is compiled into the code. To do this, place the assembly language instruction:

```
t 0x4, r1 r1
```

at the desired address. One way to do this is to create an assembly language routine that does this, then call it from your driver code.

- The system can also enter the debugger if a system halt is caused by a fatal system error. In such a case, the system creates a log entry in the system log and if the low level debugger is available, calls the low-level debugger. A system dump is generated on exit from the debugger.

If the kernel debug program is not available (nothing happens when you type in the above key sequence), you must load it. To do this, use the **bosboot** command:

```
bosboot -a -d /dev/hdisk0 -D
```

OR

```
bosboot -a -d /dev/hdisk0 -I
```

The **-D** flag causes the kernel debugger program to be loaded. The **-I** flag also causes the kernel debug program to be loaded, but it is also invoked at system initialization. This means that when the system starts, it will trap the kernel debug program.

If the kernel debug program is invoked during initialization, use the **go** subcommand to continue the initialization process. Note that you must shut down and reboot the system after the **bosboot** command in order for the changes to take effect.

You can also use the **crash** command to determine whether the kernel debug program is available. Use the **od** subcommand:

```
# crash
>od dbg_avail
```

If the **od** subcommand returns a 0 or 1, the kernel debug program is available. If it returns 2, the debug program is not available.

Debugger Subcommands and Concepts

When the kernel debugger is invoked, it is the only running process. All the interrupts are disabled and the cache is flushed. After exiting from the kernel debugger, all the processes will continue to run unless you entered the debugger through a system halt. The debugger has its own **mstsave** (machine state save) area. But what is displayed is the **mstsave** area of the running thread when the debugger was entered.

Subcommands

Once in the kernel debugger, use the subcommands to investigate and make alterations. Each command has an alias or a shortened form. This is the minimum number of letters required by the debugger to recognize the alias as unique.

The following table contains a complete list of the kernel debugger subcommands. The table is followed by a summary and examples of each subcommand.

Command	Alias	Description
alter	a	Alter memory
back	b	Decrements the Instruction Address Register (IAR)
break	br	Sets a breakpoint
breaks	breaks	Lists currently set breakpoints
buckets	bu	Displays contents of kmembucket kernel structures
clear	cl	Clears (removes) breakpoints
cpu	cp	Sets the current processor or shows processor states
display	d	Displays a specified amount of memory
dmodsw	dm	Displays the STREAMS driver switch table.
drivers	dr	Displays the contents of the device driver (devsw) table
find	f	Finds a pattern in memory
float	fl	Displays the floating point registers
fmodsw	fm	Displays the STREAMS module switch table
go	g	Starts the program running
? or help	h	Displays the list of valid commands
loop	l	Run until control returns to this point

Command	Alias	Description
map	m	Displays the system loadlist
mblk	mb	Displays the contents of message block structures
next	n	Increments the IAR
origin	o	Sets the origin
ppd	pp	Displays per-processor data
proc	pr	Displays the formatted process table
queue	que	Displays contents of STREAMS queue at specified address
quit	q	Ends a debugging session
reset	r	Releases a user-defined variable
screen	s	Displays a screen containing registers and memory
set	se	Defines or initialize a variable
sregs	sr	Displays segment registers
st	st	Stores a fullword in memory
stack	sta	Displays a formatted kernel stack trace
stc	stc	Stores one byte in memory
step	ste	Performs an instruction single-step
sth	sth	Stores a halfword in memory
stream	str	Displays stream head table
swap	sw	Switches from the current display and keyboard to another RS232 port
thread	th	Displays thread table entries
trace	tr	Displays formatted trace information
trb	trb	Displays the timer request blocks
tty	tt	Displays the ttty structure
user	u	Displays a formatted user area
uthread	ut	Displays the uthread structure
vars	v	Displays a listing of the user-defined variables
vmm	vm	Displays the virtual memory data structure
xlate	x	Translates a virtual address to a real address

alter Address Data

```
alter 1000 ffff Store the 16 bit value ffff at address 1000
a 1000 2C      Store the 8 bit value 2C in the high order
               byte at address 1000
```

The **alter** subcommand modifies the contents of memory, replacing the old value with the new specified value. If data is only one byte, it is stored in the high order byte of the word at address. If data is a halfword, it is stored in the high order halfword at address. Use this subcommand to change data while a program is running.

back [*ByteCount*]

```
back          Decrement the IAR by 4 bytes
b 16         Decrement the IAR by 16 bytes
```

This subcommand decrements the Instruction Address Register (IAR) by the specified amount.

break [*Address*]

```
break          Set a breakpoint at the IAR
break 521a     Set a breakpoint at address 521A
br 8300+A0    Set a breakpoint at A0 + 8300
break +A0     Set a breakpoint at the origin + A0
break lr      Set a breakpoint at address in the link register
```

This subcommand sets a breakpoint which causes the debugger to be entered when the instruction at the target address is next run. There is a maximum of 32 breakpoints.

breaks

```
breaks          List breakpoints
```

This subcommand displays the breakpoints as a combination of segment register values, and offsets into the segment. Since the segment registers can change during a program, an address such as 1000052C can refer to different regions of memory. The segment register value is needed, therefore, to specify a unique address.

buckets

```
buckets          Display kmembucket kernel structure for offset 0 and
                  allocation size of 2.
```

The buckets subcommand displays the contents of the **kmembucket** kernel structures. These structures contain information on the **net_malloc** memory pool by size of allocation.

All output values are printed in hexadecimal format.

clear [*Address*]

```
clear          Clear breakpoint at IAR
cl 10000200    Clear breakpoint at address 10000200
```

This subcommand removes a previously set breakpoint. If more than one breakpoint has the specified offset into the segment (low order 28 bits of address), the command prompts you with the segment values for each of the breakpoints, and asks which to clear.

cpu [*ProcessorNumber*]

```
cpu 2          Set the current processor number to 2
cpu           Shows the kernel debugger state of each processor
```

This subcommand can be used to set the current processor or to show the kernel state of all available processors.

display *Address [ByteCount]*

```
display iar      Display 16 bytes at the IAR
d 152f 12       Display 12 bytes at address 152F
display +B7     Display 16 bytes at the origin + B7
disp r3         Display 16 bytes at the address in r3
d r3>          Display from the address contained in the
               address in r3 (1 level of indirection)
```

This subcommand displays memory in hexadecimal and ASCII (characters). *ByteCount* is the number of bytes to display, and is used to determine how memory is to be accessed. If *ByteCount* is 1, a single byte is loaded. If *ByteCount* is 2, a halfword is loaded and a value of 4 causes a word to be loaded. Any other value (including none) causes memory to be loaded one byte at a time.

dmodsw

The **dmodsw** subcommand displays the internal STREAMS driver switch table, one entry at a time. By pressing the Enter key, you can walk through all the **dmodsw** entries in the table. The contents of the first entry are meaningless except for the **d_next** pointer. When the last entry has been reached, the **dmodsw** command will print the message `This is the last entry.`

The information printed is contained in an internal structure. The structure has the following members:

address	Address of dmodsw
d_next	Pointer to the next driver in the list
d_prev	Pointer to the previous driver in the list
d_name	Name of the driver
d_flags	Flags specified at configuration time
d_sqh	Pointer to synch queue for driver-level synchronization
d_str	Pointer to streamtab associated with the driver
d_sq_level	Synchronization level specified at configuration time
d_refcnt	Number of open or pushed count
d_major	Major number of a driver

The flags structure member, if set, is based on the following values:

#define	Value	Description
F_MODSW_OLD_OPEN	0x1	Supports old-style (V.3) open/close parameters
F_MODSW_QSAFETY	0x2	Module requires safe timeout/bufcall callbacks
F_MODSW_MPSAFE	0x4	Non-MP-safe drivers need funneling

The synchronization level codes are described in the `/usr/include/sys/strconf.h` file.

drivers [*Slot | Address*]

```
drivers          Display device driver (devsw) table
drivers 10       Display slot 10 of the device driver table
dr 130000f      Display last entry point before address 130000F
```

This subcommand displays the contents of a **devsw** table entry. Each **devsw** entry consists of a number of entry points (ddread, ddwrite, and so on) into the specified driver. Each entry consists of a function descriptor, and the address of the function. If you do not specify any parameters, the command displays the entire **devsw** table. If you specify a valid slot number, the command displays the corresponding entry. If the slot number is not valid, the command takes it as an address and displays the slot with the last entry point before the address (along with the name of the entry point).

find *Pattern**[*Start**][*End**][*Align**]]

```

find 7c81          Find first occurrence of 7C81 in virtual
                  memory starting at 0
find "AIX"        Find first occurrence of string AIX
f 7c81 10000      Find first 7C81 after address 10000
f 7c81 0 top      Find first 7C81 between 0 and
                  user-defined variable "top"
find 7c81 *       Find first 7C81 starting at last address used
f 7c81 * * 2     Same as above, but aligned on halfword
f 7c fx+1 * 2    Find next 7c starting at 1 + last address
                  that find stopped at

```

The **find** subcommand locates instances of the specified pattern in the current virtual address space. If you specify an * (asterisk) instead of a value, the previous value is used. For example:

```
find *
```

searches for the last pattern used.

This subcommand sets the variable **fx** to the address of the last item found. Use this feature coupled with the asterisk to continue the search for the pattern. For example:

```
f * fx+1
```

searches for the last pattern starting at the next location. The **find** subcommand remembers the alignment which was used in the previous search.

float

```
float    Display floating point registers
```

This subcommand displays the floating point registers.

fmodsw

The **fmodsw** subcommand displays the internal STREAMS module switch table, one entry at a time. By pressing the Enter key, you can walk through all the **fmodsw** entries in the table. The contents of the first entry are meaningless except for the **d_next** pointer. When the last entry has been reached, the **fmodsw** subcommand will print the message *This is the last entry.*

The information printed is contained in an internal structure. The following members of this internal structure are described here:

address	Address of fmodsw
d_next	Pointer to the next module in the list
d_prev	Pointer to the previous module in the list
d_name	Name of the module
d_flags	Flags specified at configuration time
d_sqh	Pointer to synch queue for module-level synchronization
d_str	Pointer to streamtab associated with the module
d_sq_level	Synchronization level specified at configuration time

d_refcnt	Number of open or pushed count
d_major	-1

The flags structure member, if set, is based one of the following values:

#define	Value	Description
F_MODSW_OLD_OPEN	0x1	Supports old-style (V.3) open/close parameters
F_MODSW_QSAFETY	0x2	Module requires safe timeout/bufcall callbacks
F_MODSW_MPSAFE	0x4	Non-MP-safe drivers need funneling

The synchronization-level codes are described in the `/usr/include/sys/strconf.h` file.

go [*Address*]

```
go          Continue running at the IAR
g 1000     Set the IAR to 1000 and begin running there
```

This subcommand resumes running after a breakpoint. Use it to set the Instruction Address Register (IAR) to a new address and begin running there. If this subcommand is used with no parameters after the debugger was entered via a fatal system error, a system dump will be generated and the machine will halt.

help

```
help      Display the list of valid commands
```

This subcommand displays one line for each debugger command.

loop [*Iterations*]

```
loop 2    Run until the second time the IAR has
          the current value
```

The **loop** subcommand is similar to setting a breakpoint at the current IAR, but allows you to stop on a specified instance that the IAR returns to the current point. For example, if the IAR is at the beginning of a function, entering the command:

```
loop 6
```

causes the program to continue running and allows the IAR to return to the current point five times without entering the debugger. Upon returning for a sixth time, running stops and the debugger is entered.

map [*Address* | *Symbol*]

```
map          Display the system loadlist
m e3000000   Display symbol with value closest to
             E3000000
map execexit Display the value of the symbol "execexit"
```

This subcommand displays information from the system loadlist (the list of symbols exported from the kernel). If you do not specify any parameters, the command displays the entire loadlist one page at a time. If you specify an address, the command displays the symbol value which is closest to, but less than, the address. Since **map** only knows of symbols which were exported from the kernel, this information may not be exact. If you specify a symbol name, the command searches the loadlist for the symbol, and displays any entries (there can be more than one) that match.

The symbol value for a data structure is the address of the data structure. The symbol value for a function is *not* the address of the function. Instead, it is the address of the function descriptor for that function. The first word of the function descriptor is the address of the function.

For example, if:

```
map execexit
```

displays 0x1000, then:

```
display 1000
```

displays the address of the function `execexit`.

mblk [*Address*]

```
mblk          Display M_BLK, MBDATA, and address of mh_freelater
mblk 0005ec80 Display contents of mblock structure at address 0005ec80
```

The **mblk** subcommand displays the contents of the message block structure **msgb** that is defined in the `/usr/include/sys/stream.h` file. If you do not specify an address, the command displays the contents of the message blocks of type `M_MBLK` and `M_MBDATA`, and displays the address of **mh_freelater**. **mh_freelater** is a pointer to the message blocks that have just been freed and are scheduled to be given back to the system, but have not yet been given back.

All output values are printed in hexadecimal format.

next [*ByteCount*]

```
next          Increment the IAR by 4
n 20         Increment the IAR by 20
```

This subcommand increments the Instruction Address Register (IAR) by the specified amount.

origin *Expression*

```
origin 178D   Set the origin to 178D
o 592cc       Set the origin to 592cc
```

This subcommand sets the origin variable to the value of the specified expression. Use the **origin** subcommand to match addresses with assembly language listings, which express addresses as offsets from the start of the file. Set the origin to the address in the settable of the first function in the assembly listing. Offsets from the origin are equivalent to offsets into the assembly listing. This is useful if you wish to know where the instruction you are about to run is in the listing.

ppd [*ProcessorNumber*]

```
ppd          Display the current processor structures
ppd 2        Display processor 2's structures
```

This subcommand displays the per-processor data structures of the given processor, or of the current processor (as specified by the **cpu** subcommand) by default.

proc [*PID*]

```
p           Display the process table
proc 1      Display the process table entry for
            process with process id 1
```

The **proc** subcommand with no arguments displays one line for each process in the process table (similar to the **ps** command), with an * (asterisk) placed next to the currently running process on the processor where the debugger is active. If you specify a process ID, the subcommand displays a long listing of the process table entry. This listing shows every field in use in the process table for that process table entry.

queue *Address*

```
queue 59c1874      Display the queue structure at address 59c1874
```

The **queue** subcommand displays the contents of the STREAMS queue at the specified address. Refer to the `/usr/include/sys/stream.h` file for the queue structure definition.

In the output, values marked with `x` are printed in hexadecimal format.

quit

```
quit      Clear all breakpoints and exit
```

The **quit** subcommand terminates the debug session. Use this subcommand when you have completed debugging and want to clear all breakpoints. The **quit** subcommand performs the following tasks:

- Clears all breakpoints.
- Issues the **go** subcommand, which generates a system dump if the debugger was entered via a fatal system error.

reset *Variable*

```
reset foo      Deletes the user-defined variable "foo"
```

Resetting a variable effectively deletes it, and allows the variable slot to be used again. Currently, 16 user-defined variables are allowed, and when they are all in use, you cannot set any more. Variables that are not user-defined (such as registers) cannot be reset. If you specify a variable that is not defined, the subcommand displays an error message.

screen [*Track*] [+ | - | *Address* | *On Half* | *Off* | *On*]

```
screen +      Display the next 112 bytes of memory
screen -      Display the previous 112 bytes of memory
s 20000ff7    Display memory starting at 20000ff7
s 200>       Display memory at address contained in
              location 200 (one level of indirection)
screen on     Turns on the display
screen off    Turns off the display
screen on half Display format uses about 1/2 the screen
sc track r3   Tracks memory starting at the value in
              general purpose register 3
```

The **screen** subcommand controls which information is displayed, and how much at a time. If you do not specify any parameters, the subcommand shows the current screen format (which may have been scrolled by some other command). Use the parameters to select another area of memory to display, to modify the format of the screen so that only half of the physical screen is used, or even turn off the screen display entirely. The format modification parameters are useful if important information can be scrolled off the screen when the debugger is entered. Restore the default (full) screen by entering:

```
screen on
```

The track option changes the address that the screen displays as the expression that is being tracked changes. This option is useful in a case where, at a breakpoint, the memory to be displayed is addressed by a register.

set *Variable | Register* *Value*

set start 100	Assign value 100 to variable "start"
set r12 0	Set general purpose register 12 to 0
se s3 10000	Set segment register 3 to 10000
set iar 45F0	Assign 45F0 to the IAR
se name "AIX"	Assign string "AIX" to variable "name"

This subcommand sets debugger variables. Use the **set** subcommand to create new variables or modify the value of old variables. Certain debugger variables are symbolic names for machine registers, which you can modify. They are **iar**, **r0** through **r31** (for general purpose registers), **lr** (link register), **s0** through **s15** (segment registers), **frp0** through **frp31** (floating point registers), **fpscr**, **dsisr**, **dar**, **eim0**, **eim1**, **eis0**, **eis1**, **mq**, **msr**, **cr**, **ctr**, **tid**, **xer**, **sdr0**, **sdr1**, **rtcu**, **rtcl**, and **dec**.

sregs

sregs	Display the segment registers
-------	-------------------------------

The **sregs** subcommand creates a display similar to the **screen** subcommand, but shows the segment registers along with certain other registers.

st *Address* *Word*

st 1000 5	Store the 32-bit value 5 at address 1000
-----------	--

The **st** subcommand stores a 32-bit value to memory. This is similar to the **alter** subcommand, but the word size is implicit in the command.

stc *Address* *Byte*

stc 1000 ff	Store the 8 bit value FF at address 1000
-------------	--

The **stc** subcommand stores a single byte at the specified address. This is similar to the **st** subcommand.

sth *Address* *Halfword*

sth 1000 0014	Store the 16 bit value 14 at address 1000
---------------	---

The **sth** subcommand stores a halfword at the specified address. This is similar to the **st** subcommand.

stack [*ThreadID*]

stack	Format any existing stack frames
sta 25	Format stack frames for the thread with ID 25

This subcommand formats the stack frames for the specified (or by default, the currently running) thread. Stack frames show return addresses. Use them to trace the program's calling sequence. Be aware that the first few parameters to the called functions are passed in registers, and are not usually available on the stack. Generally, only the chain (stack back-chain pointer) and return address (caller of the owner of this stack frame) are valid.

To thoroughly interpret the stack, you must use the assembly language listing for a procedure to determine what is stored on the stack. Stack frames for the specified thread are not always accessible.

step [*s* | *Number*]

step	Single step the processor
step s	Single step (skip over a subroutine call)
ste 20	Step for 20 instructions

The **step** subcommand allows the processor to single step (run a single instruction). Use the **s** parameter to skip over subroutine calls and follow the main flow of control. Use an integer parameter to specify the number of instructions to run before returning control to the debugger.

stream [*Address*]

```
stream          Display all entries in the stream head table
str 59b2e00    Display the stream head table at address 59b2e00
```

The **stream** subcommand displays the contents of the stream head table. If no address is specified, the subcommand displays every stream found in the STREAMS hash table. If the address is specified, the subcommand displays the contents of the stream head stored at that address.

The information printed is contained in an internal structure. The following members of this internal structure are described here:

```
sth            Address of stream head
wq            Address of streams write queue
rq            Address of streams read queue
dev           Associated device number of the stream
read mode     Read mode
write mode    Write mode
close_wait_timeout
              Close wait timeout in microseconds
read error    Read error on the stream
write error   Write error on the stream
flags         Stream head flag values
push_cnt     Number of modules pushed on the stream
wroff        Write offset to prepend M_DATA
ioc_id       ID of outstanding M_IOCTL request
ioc_mp       Outstanding ioctl message
next         Next stream head on the link
pollq        List of active polls
sigsq        List of active M_SETSIGs
shsttyp      Pointer to tty information
```

The **read_mode** and **write_mode** values are defined in the **/usr/include/sys/stropts.h** file.

The **read_error** and **write_error** variables are integers defined in the **/usr/include/sys/errno.h** file.

The flags structure member, if set, is based on combinations of the following values:

#define	Value	Description
F_STH_READ_ERROR	0x0001	M_ERROR with read error received, fail all read calls
F_STH_WRITE_ERROR	0x0002	M_ERROR with write error received, fail all writes
F_STH_HANGUP	0x0004	M_HANGUP received, no more data
F_STH_NDELOD	0x0008	Do TTY semantics for ONDELAY handling
F_STH_ISATTY	0x0010	This stream acts a terminal
F_STH_MREADON	0x0020	Generate M_READ messages
F_STH_TOSTOP	0x0040	Disallow background writes (for job control)
F_STH_PIPE	0x0080	Stream is one end of a pipe or FIFO
F_STH_WPIPE	0x0100	Stream is the "write" side of a pipe
F_STH_FIFO	0x0200	Stream is a FIFO
F_STH_LINKED	0x0400	Stream has one or more lower streams linked
F_STH_CTTY	0x0800	Stream controlling tty
F_STH_CLOSED	0x4000	Stream has been closed, and should be freed
F_STH_CLOSING	0x8000	Actively on the way down

In the output, values marked with `x` are printed in hexadecimal format.

swap *PortNumber*

```
swap 1          Switch display to RS-232 port 1
```

The **swap** subcommand allows you to transfer control of the debugger to another terminal. The parameter is the port number. Specify 1 for port 1 (s0) or 2 for port 2 (s1).

thread [*ProcessID* | *ThreadID*]

This subcommand displays the contents of the kernel thread table. If a process ID is given, all the kernel threads belonging to the process are shown. If a thread ID is given, only information about that kernel thread is shown. If no parameter is given, all kernel threads in the kernel thread table are displayed.

trace

```
trace          Display the kernel trace buffers
```

The **trace** subcommand displays the last 128 entries of the kernel trace buffers. Currently there are eight trace channels that can trace any combination of events. Event entries are placed in a trace buffer before they are written out to disk (if they are logged to disk), or the buffers are used in a circular fashion and reused when full. The **trace** subcommand allows you to examine the trace headers which contain pointers into the trace buffers. The subcommand shows the trace entries and gives hook IDs in text form and data in hex.

trb

```
trb           Display timer request block information
```

This subcommand shows a menu of commands to display timer request block information.

tty [-d] [-l] [-e] [*Name* | *Major* [*Minor*]]

Displays the **tty** structures. If no parameters are specified, a short list of all open terminals is displayed. Selected terminals can be displayed by specifying the terminal name, such as **tty1**, or a major number with optional minor number. The flags modify the displayed information: the **-d** flag displays driver information; the **-l** flag displays line discipline information; and the **-e** flag displays information for every module or driver present in the stream for the selected lines.

user [*TID*]

```
user          Display current user area
u 315         Display user area for thread with thread ID 315
```

This subcommand displays the user area for the currently running thread. If you specify a thread ID, the subcommand displays the user area for that thread.

uthread [*ThreadID*]

This subcommand displays the **uthread** structure of the given thread, or of the current thread by default.

vars

```
vars          Display current user-defined variables
```

This subcommand displays a list of the user-defined debugger variables. The subcommand displays the variable name and value, and an indication of what is the base of the value. Since the value 10 can be either decimal or hexadecimal it is displayed as `HEX/DEC`. The subcommand displays string variables with no quotes around the string value.

vmm

vmm Display virtual memory data structures

This subcommand displays a menu of commands for examining the virtual memory data structures. Useful information such as the number of free pages is available.

xlate *VirtualAddress*

xlate 10054000 Translate virtual 10054000 to a real address

This subcommand translates a virtual address to a real address.

Numeric Values and Strings

All subcommands, except **loop** and **step**, require numeric options in hexadecimal. The **loop** and **step** subcommands take decimal values. Decimal numbers must either be decimal constants, variables, or expressions involving constants and variables.

A string is either a hexadecimal constant or a character constant. Use double quotes to delimit the string from other data.

Variables

Variable names must start with a letter and can be up to eight characters long. Variable names cannot contain special symbols. Variables usually represent locations or values which are used again and again. A variable must not represent a valid number. Use the **set** subcommand to define and initialize variables. Variables can contain from 1 to 4 bytes of numeric data or up to 32 characters of string data. You can release a variable with the **reset** subcommand. You cannot use the **reset** subcommand with reserved variables.

For example:

```
set name 1234          Sets your variable called name=1234
set s8 820c00e0       Sets seg reg 8 to point to the IOCC
```

Note that `s8` is a reserved variable.

Reserved Variables

There is a set of variables that have a reserved meaning for the kernel debug program. You can reference and change these variables, but they represent the actual hardware registers. There are also two variables (*fx* and *org*) reserved for use by the kernel debug program, which can be changed or set. If you change the segment registers or the general purpose registers while in the kernel debug program, the change remains in effect when you leave the kernel debug program. The reserved variables are:

bat0l	BAT register 0, lower
bat0u	BAT register 0, upper
bat1l	BAT register 1, lower
bat1u	BAT register 1, upper
bat2l	BAT register 2, lower
bat2u	BAT register 2, upper
cr	Condition register.
ctr	Count register.
dar	Data address register.
dec	Decrementer.
dsier	Data storage interrupt error register.

dsisr	Data storage interrupt status register.
eim0	External interrupt mask (low).
eim1	External interrupt mask (high).
eis0	External interrupt summary (low).
eis1	External Interrupt summary (high).
fp0–fp31	Floating point registers 0 through 31.
fpscr	Floating point status and control register.
fx	Address of the last item found by the find subcommand.
iar	Instruction Address Register (program counter). Points to the current instruction.
lr	Link register.
mq	Multiply quotient.
msr	Machine State register.
org	The current value of origin. It is useful to set this to the program load point.
peis0	Pending external interrupt status register 0.
peis1	Pending external interrupt status register 1.
r0 – r31	General Purpose Registers 0 through 31. These registers have the following usage conventions:
r0	Used on prologs. Not preserved across calls.
r1	Stack pointer. Preserved across calls.
r2	TOC. Preserved across calls.
r3 – r10	Parameter list for a procedure call. The first argument is r3, the second is r4 and so on until r10 is the 8th argument. These registers are not preserved across calls.
r11	Scratch. Pointer to FCN; DSA pointer to <code>int proc(env)</code> .
r12	PL8 exception return. Value preserved across calls.
r13–r31	Scratch. Value preserved across calls.
rtcl	Real Time clock (nano seconds).
rtcu	Real Time clock (seconds).
s0–s15	Segment registers. If a segment register is <i>not</i> in use, it has a value of 007FFFFF.
sdr0	Storage description register 0.
sdr1	Storage description register 1.
srr0	Machine status save/restore 0.
srr1	Machine status save/restore 1.
tbl	Time base register, lower
tbu	Time base register, upper

tid	Transaction register (fixed point).
xer	Exception register (fixed point).
xirr	External interrupt request register.

Expressions

The kernel debug program does not allow full expression processing. Expressions can only contain decimal or hex constants, variables and operators. The variable operators include:

+	addition
-	subtraction
*	multiplication
/	division
>	dereference

The **>** operator indicates that the value of the preceding expression is to be taken as the address of the target value. The contents of the address specified by the evaluated expression are used in place of the expression.

You can enter expressions in the form *Expression(Expression)*. This form causes the two expressions to be evaluated separately and then added together. This form is similar to the base address syntax used in the assembler.

Expressions are processed from left to right only. The type of data specified must be the same for all terms in the expression.

Pointer Dereferences

A pointer dereference can refer indirectly to the contents of a memory location. For example, assume location 0xC50 contains a counter. Use an expression of the form:

```
c50>
```

to refer to the counter. You can put any expression before the **>**, including an expression involving another **>**. In this case, it indicates multiple levels of indirection. To extend the previous example, if location FF7 contains the value C50, the expression:

```
ff7>>
```

refers to the previously discussed counter.

Breakpoints

The debugger creates a table of breakpoints that it internally maintains. The **break** subcommand creates breakpoints. The **clear** subcommand clears breakpoints. When the breakpoint is set, the debugger replaces the corresponding instruction with the trap instruction. A breakpoint can only be set if the instruction is not paged out.

Breakpoints should not be set in any code used by the debugger.

For more information, see "Setting Breakpoints" on page 14-45.

Maps and Listings

The assembler listing and the map files are essential tools for debugging using the kernel debugger. In order to create the assembly list file during compilation use the **-qlist** option while compiling. Also use the **-qsource** option to get the C source listing in the same file:

```
cc -c -DEBUG -D_KERNEL -DIBMR2 demodd.c -qsource -qlist
```

In order to obtain the map file, use the **-bmap:FileName** option on the link editor, enter:

```
ld -o demodd demodd.o -edemoconfig -bimport:/lib/kernex.exp \  
-lsys -lcsys -bmap:demodd.map -bE:demodd.exp
```

You can also create a map file with a slightly different format by using the **nm** command. For example, use the following command to get a map listing for the kernel (**/unix**):

```
nm -xv /unix > unix.m
```

Compiler Listing

The assembler and source listing is used to correlate any C source line with the corresponding assembler lines. The following is a portion of the C source code for a sample device driver. The left column is the line number in the source code:

```
.  
.   
185  
186     if (result = devswadd(devno, &demo_dsw)){  
187         printf("democonfig : failed to add entry points\n");  
188         (void)devswdel(devno);  
189         break;  
190     }  
191     dp->initied = 1;  
192     demos_initied++;  
193     printf("democonfig : CFG_INIT success\n");  
194     break;  
195  
.   
.
```

The following is a portion of the assembler listing for the corresponding C code shown previously. The left column is the C source line for the corresponding assembler statement. Each C source line can have multiple assembler source lines. The second column is the offset of the assembler instruction with respect to the kernel extension entry point.

```

.
.
186| 000218 l      80610098  2  L4Z   gr3=devno(gr1,152)
186| 00021C cal    389F0000  1  LR    gr4=gr31
186| 000220 bl     4BFFFDE1  0  CALL  gr3=devswadd,2,
gr3,(struct_4198576)",gr4,devswadd",gr1,cr[01567],gr0",
gr4"-gr12",fp0"-fp13"
186| 000224 cror   4DEF7B82  1
186| 000228 st     9061005C  2  ST4A  #2357(gr1,92)=gr3
186| 00022C st     9061003C  1  ST4A  result(gr1,60)=gr3
186| 000230 l      8061005C  1  L4A   gr3=#2357(gr1,92)
186| 000234 cmpi   2C830000  2  C4    cr1=gr3,0
186| 000238 bc     41860020  3  BT    CL.16,cr1,0x4/eq
187| 00023C ai     307F01A4  1  AI    gr3=gr31,420
187| 000240 bl     4BFFFDC1  2  CALL  gr3=printf,1,'democonfig :
failed to add entry points",gr3,printf",gr1,cr[01567],gr0",
gr4"-gr12",fp0"-fp13"
187| 000244 cror   4DEF7B82  1
188| 000248 l      80610098  2  L4Z   gr3=devno(gr1,152)
188| 00024C bl     4BFFFDB5  0  CALL  gr3=devswdel,1,gr3,
devswdel",gr1,cr[01567],gr0",gr4"-gr12",fp0"-fp13"
188| 000250 cror   4DEF7B82  1
189| 000254 b      48000104  0  B     CL.6
186|                                     CL.16:
191| 000258 l      80810040  2  L4Z   gr4=dp(gr1,64)
191| 00025C cal    38600001  1  LI    gr3=1
191| 000260 stb    98640004  1  ST1Z  (char)(gr4,4)=gr3
192| 000264 l      8082000C  1  L4A   gr4=.demos_initiated(gr2,0)
192| 000268 l      80640000  2  L4A   gr3=demos_initiated(gr4,0)
192| 00026C ai     30630001  2  AI    gr3=gr3,1
192| 000270 st     90640000  1  ST4A  demos_initiated(gr4,0)=gr3
193| 000274 ai     307F01D0  1  AI    gr3=gr31,464
193| 000278 bl     4BFFFDB9  0  CALL  gr3=printf,1,'democonfig :
CFG_INIT success",gr3,printf",gr1,cr[01567],gr0",gr4"-gr12",
fp0"-fp13"
193| 00027C cror   4DEF7B82  1
194| 000280 b      480000D8  0  B     CL.6
.
.

```

Now with both the assembler listing and the C source listing, you can determine the assembler instruction for a C statement. As an example, consider the C source line at line 191 in the sample code:

```
191          dp->initiated = 1;
```

The corresponding assembler instructions are:

```

191| 000258 l      80810040  2  L4Z   gr4=dp(gr1,64)
191| 00025C cal    38600001  1  LI    gr3=1
191| 000260 stb    98640004  1  ST1Z  (char)(gr4,4)=gr3

```

The offsets of these instructions within the sample device driver (demodd) are 000258, 00025C, and 000260.

Map File

The binder map file is a symbol map in address order format. Each symbol listed in the map file has a storage class (CL) and a type (TY) associated with it.

Storage classes correspond to the **XMC_XX** variables defined in the **syms.h** file. Each storage class belongs to one of the following section types:

.text	Contains read-only data (instructions). Addresses listed in this section use the beginning of the .text section as origin. The .text section can contain one of the following storage class (CL) values:
DB	Debug Table. Identifies a class of sections that has the same characteristics as read only data.
GL	Glue Code. Identifies a section that has the same characteristics as a program code. This type of section has code to interface with a routine in another module. Part of the interface code requirement is to maintain TOC addressability across the call.
PR	Program Code. Identifies the sections that provide executable instructions for the module.
R0	Read Only Data. Identifies the sections that contain constants that are not modified during execution.
TB	Reserved.
TI	Reserved.
XO	Extended Op. Identifies a section of code that is to be treated as a pseudo-machine instruction.
.data	Contains read-write initialized data. Addresses listed in this section use the beginning of the .data section as origin. The .data section can contain one of the following storage class (CL) values:
DS	Descriptor. Identifies a function descriptor. This information is used to describe function pointers in languages such as C and Fortran.
RW	Read Write Data. Identifies a section that contains data that is known to require change during execution.
SV	SVC. Identifies a section of code that is to be treated as a supervisory call.
T0	TOC Anchor. Used only by the predefined TOC symbol. Identifies the special symbol TOC. Used only by the TOC header.
TC	TOC Entry. Identifies address data that will reside in the TOC.
TD	TOC Data Entry. Identifies data that will reside in the TOC.
UA	Unclassified. Identifies data that contains data of an unknown storage class.
.bss	Contains read-write uninitialized data. Addresses listed in this section use the beginning of the .data section as origin. The .bss section contain one of the following storage class (CL) values:
BS	BSS class. Identifies a section that contains uninitialized data.
UC	Unnamed Fortran Common. Identifies a section that contains read write data.

Types correspond to the **XTY_XX** variables defined in the **syms.h** file. The type (TY) can be one of the following values:

ER External Reference
LD Label Definition
SD Section Definition
CM BSS Common Definition

The following is a map file for a sample device driver:

```

1 ADDRESS MAP FOR demodd
2
3 *IE ADDRESS LENGTH AL CL TY Sym# NAME SOURCE-FILE(OBJECT) or
4 -----
5 I ER S1 pinned_heap /lib/kernex.exp{/unix}
6 I ER S2 devswadd /lib/kernex.exp{/unix}
7 I ER S3 devswdel /lib/kernex.exp{/unix}
8 I ER S4 nodev /lib/kernex.exp{/unix}
9 I ER S5 printf /lib/kernex.exp{/unix}
10 I ER S6 uiomove /lib/kernex.exp{/unix}
11 I ER S7 xmalloc /lib/kernex.exp{/unix}
12 I ER S8 xfree /lib/kernex.exp{/unix}
13 00000000 0008B8 2 PR SD S9 <>
/tmp/cliff/demodd/demodd.c(demodd.o)
14 00000000 PR LD S10 .democonfig
15 0000039C PR LD S11 .demoopen
16 000004B4 PR LD S12 .democlose
17 000005D4 PR LD S13 .demoread
18 00000704 PR LD S14 .demowrite
19 00000830 PR LD S15 .get_dp
20 000008B8 000024 2 GL SD S16 <.printf> glink.s(/usr/lib/glink.o)
21 000008B8 GL LD S17 .printf
22 000008DC 000024 2 GL SD S18 <.xmalloc> glink.s(/usr/lib/glink.o)
23 000008DC GL LD S19 .xmalloc
24 00000900 000090 2 PR SD S20 .bzero
noname(/usr/lib/libcsys.a[bzero.o])
25 00000990 000024 2 GL SD S21 <.uiomove> glink.s(/usr/lib/glink.o)
26 00000990 GL LD S22 .uiomove
27 000009B4 000024 2 GL SD S23 <.devswadd> glink.s(/usr/lib/glink.o)
28 000009B4 GL LD S24 .devswadd
29 000009D8 000024 2 GL SD S25 <.devswdel> glink.s(/usr/lib/glink.o)
30 000009D8 GL LD S26 .devswdel
31 000009FC 000024 2 GL SD S27 <.xmalloc> glink.s(/usr/lib/glink.o)
32 000009FC GL LD S28 .xmalloc
33 00000000 000444 4 RW SD S29 <_/tmp/cliff/demodd/demodd$c$>
/tmp/cliff/demodd/demodd.c(demodd.o)
34 00000450 000004 4 RW SD S30 demo_dev
/tmp/cliff/demodd/demodd.c(demodd.o)
35 00000460 000004 4 RW SD S31 demos_inited
/tmp/cliff/demodd/demodd.c(demodd.o)
36 00000470 000080 4 RW SD S32 data
/tmp/cliff/demodd/demodd.c(demodd.o)
37 * E 000004F0 00000C 2 DS SD S33 democonfig
/tmp/cliff/demodd/demodd.c(demodd.o)
38 E 000004FC 00000C 2 DS SD S34 demoopen
/tmp/cliff/demodd/demodd.c(demodd.o)
39 E 00000508 00000C 2 DS SD S35 democlose
/tmp/cliff/demodd/demodd.c(demodd.o)
40 E 00000514 00000C 2 DS SD S36 demoread
/tmp/cliff/demodd/demodd.c(demodd.o)
41 E 00000520 00000C 2 DS SD S37 demowrite
/tmp/cliff/demodd/demodd.c(demodd.o)
42 0000052C 000000 2 T0 SD S38 <TOC>
43 0000052C 000004 2 TC SD S39 <_/tmp/cliff/demodd/demodd$c$>
44 00000530 000004 2 TC SD S40 <printf>
45 00000534 000004 2 TC SD S41 <demo_dev>

```

```

45      00000538 000004 2 TC SD S42 <demos_initied>
46      0000053C 000004 2 TC SD S43 <data>
47      00000540 000004 2 TC SD S44 <pinned_heap>
48      00000544 000004 2 TC SD S45 <xmalloc>
49      00000548 000004 2 TC SD S46 <uiomove>
50      0000054C 000004 2 TC SD S47 <devswadd>
51      00000550 000004 2 TC SD S48 <devswdel>
52      00000554 000004 2 TC SD S49 <xmfree>

```

In the sample map file listed previously, the **.data** section starts from the statement at line 32:

```

32      00000000 000444 4 RW SD S29 <_/tmp/cliff/demodd/demodd$c$>
      /tmp/cliff/demodd/demodd.c(demodd.o)

```

The TOC (Table of Contents) starts from the statement at line 41:

```

41      0000052C 000000 2 T0 SD S38 <TOC>

```

Using the Debugger

This section contains information on setting breakpoints, viewing and modifying global data, displaying registers, and using the stack trace.

Setting Breakpoints

Setting a breakpoint is essential for debugging kernel or kernel extensions. To set a breakpoint, use the following sequence of steps:

1. Locate the assembler instruction corresponding to the C statement.
2. Get the offset of the assembler instruction from the listing.
3. Locate the address where the kernel extension is loaded.
4. Add the address of the assembler instruction to the address where kernel extension is loaded.
5. Set the breakpoint with the **break** command.

The process of locating the assembler instruction and getting its offset is explained in the previous section. The next step is to get the address where the kernel extension is loaded.

Determine the Location of your Kernel Extension

To determine the address where a kernel extension has been loaded, use the following procedure. First, find the load point (the entry point) of the executable kernel extension. This is a label supplied with the **-e** option for the **ld** (links objects) command used while generating the kernel extension. In our example this is the **democonfig** routine.

Then use one of the following six methods to locate the address of this load point. This address is the location where the kernel extension is loaded.

Method 1

If the kernel extension is a device driver, use the **drivers** subcommand to locate the address of the load point routine. The **drivers** subcommand lists all the function descriptors and the function addresses for the device driver (that are in the dev switch table). Usually the **config** routine will be the load point routine. Hence in our example the function address for the **config** (**democonfig**) routine is the address where the kernel extension is loaded.

```

> drivers 255
MAJ#255      Open          Close          Read           Write
func  desc  0x01B131B0    0x01B131BC    0x01B131C8    0x01B131D4
func  addr  0x01B12578    0x01B126A0    0x01B127D4    0x01B12910
          Ioctl          Strategy       Tty            Select
func  desc  0x00019F10    0x00019F10    0x00000000    0x00019F10
func  addr  0x00019A20    0x00019A20    0x00019A20    0x00019A20
          Config       Print          Dump           Mpx
func  desc  0x01B131A4    0x00019F10    0x00019F10    0x00019F10
func  addr  0x01B121EC    0x00019A20    0x00019A20    0x00019A20
          Revoke       Dsdptr        Selptr         Opts
func  desc  0x00019F10    0x00000000    0x00000000    0x00000002
func  addr  0x00019A20

```

Method 2

Another method to locate the address is to use the value of the **kmid** pointer returned by the **sysconfig(SYS_KLOAD)** subroutine when loading the kernel extension. The **kmid** pointer points to the address of the load point routine. Hence to get the address of the load point, print the **kmid** value during the **sysconfig** call from the configuration method. Then go into the low level debugger and display the value pointed to by **kmid**. For clarity, set mnemonics for **kmid**.

```

> set kmid 1b131a4
> vars
Listing of the User-defined variables:
kmid HEX=01B131A4
fx HEX/DEC=01B1256E
org
There are 15 free variable slots.
> d kmid
01B131A4  01B121EC 01B131E0 00000000 01B12578
|..!...1.....%x|
> d kmid>
01B121EC  7C0802A6 BFC1FFF8 90010008 9421FF80
||.....!...|

```

Method 3

If **kmid** is also not known, use the **find** subcommand to locate the load point routine:

```

> find democonfig 1b00000
01B1256E  66616B65 636F6E66 69677C08 02A693E1
|democonfig|.....|

```

The **find** subcommand will locate the specified string. It initiates a search from the starting address specified in the command. The string that is located is at the end of the **democonfig** routine. Now, backup to locate the beginning of the routine.

Usually all procedures have the instruction 7C0802A6 within the first three or four instructions of the procedure (within the first 12 to 16 bytes). See the assembler listing for the actual position of this instruction within the procedure. Use the **screen** subcommand with the **-** flag to keep going back to locate the instruction. You can help speed up your search by using the ASCII section of the screen output to look for occurrences of the pipe symbol (**|**), which corresponds to the hexadecimal value 7C, the first byte of the instruction. Once this instruction is found, you can figure out where the start of the procedure is using the assembler listing as a guide.

```

> screen fx
GPR0  000078E4  2FF7FF70  000C5E78  00000000  2FF7FFF8  00000000  00007910  DEADBEEF
GPR8  DEADBEEF  DEADBEEF  DEADBEEF  7C0802A6  DEADBEEF  DEADBEEF  DEADBEEF  DEADBEEF
GPR16 DEADBEEF  DEADBEEF  DEADBEEF  DEADBEEF  DEADBEEF  DEADBEEF  DEADBEEF  DEADBEEF
GPR24 DEADBEEF  DEADBEEF  DEADBEEF  DEADBEEF  DEADBEEF  DEADBEEF  DEADBEEF  DEADBEEF  00007910
MSR   000090B0  CR    00000000  LR    0002506C  CTR   000078E4
MQ    00000000  XER   00000000  SRR0  000078E4  SRR1  000090B0  DSISR  40000000
DAR   30000000  IAR   000078E4  (ORG+000078E4)  ORG=00000000  Mode:  VIRTUAL
000078E0  00000000  48000000  4E800020  00000000  |...H...N...  ....|
          |
          |  b  0x78E4  (000078E4)
000078F0  000C0000  00000000  00000000  00000000  |.....|
          |
01B12560  80020301  00000000  0000036C  000A6661  |.....l..fa|
01B12570  6B65636F  6E666967  7C0802A6  93E1FFFC  |keconfig|.....|
01B12580  90010008  9421FFA0  83E20000  90610078  |.....!.....a.x|
01B12590  9081007C  90A10080  90C10084  307F0294  |...|.....0...|
01B125A0  48000535  80410014  80610078  5463043E  |H..5.A...a.xTc.>|
01B125B0  90610038  80610078  48000491  9061003C  |.a.8.a.xH....a.<|
01B125C0  28830000  41860020  8061003C  88630004  |(..A.. .a.<.c..|

```

```

> screen -
.
.
>
>
GPR0  000078E4  2FF7FF70  000C5E78  00000000  2FF7FFF8  00000000  00007910  DEADBEEF
GPR8  DEADBEEF  DEADBEEF  DEADBEEF  7C0802A6  DEADBEEF  DEADBEEF  DEADBEEF  DEADBEEF
GPR16 DEADBEEF  DEADBEEF  DEADBEEF  DEADBEEF  DEADBEEF  DEADBEEF  DEADBEEF  DEADBEEF
GPR24 DEADBEEF  DEADBEEF  DEADBEEF  DEADBEEF  DEADBEEF  DEADBEEF  DEADBEEF  DEADBEEF  00007910
MSR   000090B0  CR    00000000  LR    0002506C  CTR   000078E4  MQ    00000000
XER   00000000  SRR0  000078E4  SRR1  000090B0  DSISR  40000000  DAR   30000000
IAR   000078E4  (ORG+000078E4)  ORG=00000000  Mode:  VIRTUAL
000078E0  00000000  48000000  4E800020  00000000  |...H...N...  ....|
          |
          |  b  0x78E4  (000078E4)
000078F0  000C0000  00000000  00000000  00000000  |.....|
          |
01B121E0  00000000  00000000  00000000  7C0802A6  |.....|...|
01B121F0  BFC1FFF8  90010008  9421FF80  83E20000  |.....!.....|
01B12200  90610098  9081009C  90A100A0  307F0040  |.a.....0..@|
01B12210  80810098  480008C1  80410014  307F0058  |...H...A..0..X|
01B12220  83C20008  63C40000  80A2000C  80C20010  |...c.....|
01B12230  480008A5  80410014  63C30000  80810098  |H...A..c.....|
01B12240  5484043E  90810038  38800000  9081003C  |T...>...88.....<|

```

The start of the democonfig routine is at 0x01B121EC.

Method 4

If the load point routine is an exported routine, use the **map** subcommand to locate the appropriate routine:

```
>map <routine name>
```

Method 5

You can also use the **crash** command to locate the load point. After running the **crash** command, run the **le** subcommand to list the load point for all the kernel extensions. The **knlist** subcommand will list the addresses of exported symbols:

```
$ crash
>le
>quit
```

The **le** subcommand shows the module start address. The first procedure in the kernel extension would follow the module header from the module start address. Hence in the case of the example demodd kernel extension, **le** showed the module start address to be 0x01B12000 and the democonfig procedure starts at 0x01B121EC.

You can locate the start of the democonfig procedure by searching for the first instruction of the democonfig procedure which would be usually 0x7C0802A6. Use the assembly listing to determine the first instruction.

First, display memory at 0x01b12000 and then use the **screen** subcommand to search ahead.

```
>screen 01b12000
>screen +
:
.
```

Method 6

Use the **find** subcommand to search for a pattern:

```
> find democonfig 1b00000
01B1256E 66616B65 636F6E66 69677C08 02A693E1
|democonfig|.....|
```

We know that the module starts before 1B1256E. We also know that the “magic” number is 01DF. The loader identifies a file as a load module by looking for 01DF as the first two bytes in the file. So, the greatest address which is less than 1B1256E that contains 01DF, will be the start of the module, provided that it is on a page boundary. This means it has a mask of FFFF000, a 4096 boundary or 0x1000:

```
> find 01df 01900000 * 2
```

Search starting at 1900000 through the kernel storage (the *) for 01DF on a 2-byte boundary.

The greatest address, on a page boundary, that is less than 1B1256E will be the module start. This will be offset 00000000 in the map file.

Change the Origin

Set the origin to the address of the load point. By default this is zero. By changing the origin to the address of the load point, you can directly correlate the address in the assembler listing with the address for the Instruction Address Register (IAR) and break points.

```
>set fkcfg 1B121EC          set a variable called fkcfg

>origin fkcfg
```

Set the Break Point

Now set the break point with the **break** subcommand. Assume that we want to set the breakpoint at the assembler instruction at offset 218 (using the assembler listing):

```
>break +218                If origin has been set to load point
```

OR

```
>break 1B121EC+218
```

Viewing and Modifying Global Data

You can access the global data with two different methods. To understand how to locate the address of a global variable, we use the example of our demodd device driver. Here we try to view and modify the value of the data[] character array in the sample demodd device driver.

Use the first method only when you break in a procedure for the kernel extension to be debugged. You can use the second method at any time.

Method 1

1. After getting into the low level debugger, set a break point at the **demoread** procedure call. You can use any routine in demodd for this purpose.
2. Call the **demoread** routine. When the system breaks in **demoread** and invokes the debugger, the GPR2 (general purpose register 2) points to the TOC address. Now use the offset of the address of any global variable (from the start of TOC) to determine its address. The TOC is listed in the map file.

The map file on page 14-44 shows that the address of the data[] array is at 0x53C while the TOC is at 0x52C. The offset of the address of the data[] array with respect to the start of TOC is $0x53C - 0x52C = 0x10$. Hence the address of the data[] variable is at (r2+10). And the actual data[] variable is located at the address value in (r2 + 10):

```
> d r2
01B131E0 01B12CCC 0004E7D0 01B13114 01B1311C |.....1...1.|
> d r2+10>
01B13124 61626364 65666768 696A6B6C 6D6E6F70 |abcdefghijklmnop|
```

Now we can change the value of the data[] variable. As an example, we change the first four bytes of data[] to "pppp" (p = 70):

```
> st r2+10> 70707070
> d r2+10>
01B13124 70707070 65666768 696A6B6C 6D6E6F70 |ppppefghijklmnop|
```

Method 2

You can use this method at any time. This method requires the map file and the address at which the relevant kernel address has been loaded. This method currently works because of the manner in which a kernel extension is loaded. But it may not work if the procedure for loading a kernel extension changes.

The address of a variable is:

Address of the last function before the variable in the map file	+	Length of the function	+	Offset of the variable
--	---	---------------------------	---	---------------------------

The following is the section of the map file (see page 14-44) showing the data[] variable and the last function (xmfree) in the **.text** section:

```
26 000009B4 000024 2 GL SD S23 <.devswadd> glink.s(/usr/lib/glink.o)
27 000009B4 GL LD S24 .devswadd
28 000009D8 000024 2 GL SD S25 <.devswdel> glink.s(/usr/lib/glink.o)
29 000009D8 GL LD S26 .devswdel
30 000009FC 000024 2 GL SD S27 <.xmfree> glink.s(/usr/lib/glink.o)
31 000009FC GL LD S28 .xmfree
32 00000000 000444 4 RW SD S29 <_/tmp/cliff/demodd/demodd$c$>
/tmp/cliff/demodd/demodd.c(demodd.o)
33 00000450 000004 4 RW SD S30 demo_dev
/tmp/cliff/demodd/demodd.c(demodd.o)
34 00000460 000004 4 RW SD S31 demos_inited
/tmp/cliff/demodd/demodd.c(demodd.o)
35 00000470 000080 4 RW SD S32 data
/tmp/cliff/demodd/demodd.c(demodd.o)
36 * E 000004F0 00000C 2 DS SD S33 democonfig
/tmp/cliff/demodd/demodd.c(demodd.o)
37 E 000004FC 00000C 2 DS SD S34 demoopen
/tmp/cliff/demodd/demodd.c(demodd.o)
```

The last function in the `.text` section is at lines 30–31. The offset address of this function from the map is 0x000009FC (line 30, column 2). The length of the function is 0x000024 (line 30, column 3). The offset address of the `data[]` variable is 0x00000470 (line 35, column 2). Hence the offset of the address of the `data[]` variable is:

```
0x000009FC + 0x000024 + 0x00000470 = 0x00000E90
```

Add this address value to the load point value of the demodd kernel extension. If, as in the case of the sample demodd device handler, this is 0x1B131A4, then the address of the `data[]` variable is:

```
0x1B121EC + 0x00000E90 = 0x1B1307C
```

```
>display 1B1307C
01B1307C 61626364 65666768 696A6B6C 6D6E6F70 |abcdefghijklmnop|
```

Now change the value of the `data[]` variable as in method 1.

Note that in method 1, using the TOC, you found the address of the address of `data[]`, while in method 2 you simply found the address of `data[]`.

Displaying Registers on a Micro Channel Adapter

When you write a device driver for a new Micro Channel adapter, you often want to be able to read and write to registers that reside on the adapter. This is a way of seeing if the hardware is functioning correctly. For example, to examine a register on the Token Ring adapter, first see where this adapter resides in the bus I/O space:

```
$lsdev -C

sys0          Available 00-00 System Object
sysunit0     Available 00-00 RISC System/6000 System Unit
sysplanar0   Available 00-00 CPU Planar
.
.
scsi0        Available 00-01 SCSI I/O Controller
tok0         Available 00-02 Token-Ring High-Performance Adapter
ent0         Available 00-03 Ethernet High-Performance LAN Adapter

$lsattr -l tok0 -E

bus_intr_lvl      3      Bus interrupt level False
intr_priority     3      Interrupt priority  False
.
.
rdto              92      RECEIVE DATA TRANSFER OFFSET      True
bus_io_addr       0x86a0  Bus I/O address                    False
dma_lvl           0x5     DMA arbitration level              False
dma_bus_mem       0x202000 Address of bus memory used DMA      False
```

We now know that the token ring adapter is located at 0x86A0.

To read a specific register, enter the kernel debugger and use the `sregs` subcommand to display the segment registers. Find an unused segment register (=007FFFFF). For this example, assume s9 is not used. Enable the Micro Channel bus addressing with the `set` subcommand:

```
set s9 820c0020
```

Use the `sregs` subcommand to display the segment register values to check that you typed it in correctly.

From the *POWERstation and POWERserver Hardware Technical Information-Options and Devices*, we know that the address of the Adapter Communication and Status register is P6a6. The value of P is based on the Bus I/O address (bus_io_addr) of the adapter. In the above example, this is 86A0. It could have been anything from 86A0 to F6A0 on a 0x1000 byte boundary. Hence P is 8, and the address of the Communication and Status register is 86A6. The **display** subcommand now displays the two-byte register:

```
d 900086a6 2
```

The key is to load a segment register with 820c0020 and then use that segment register to reference registers and memory on your adapter. You can use the same method to access registers resident on the IOCC. In that case, load the segment register with a value of 820c00e0.

Stack Trace

The stack trace gives the stack history which provides the sequence of procedure calls leading to the current IAR. The **Ret Addr** is the address of the instruction calling this procedure. You can use the map file to locate the name of the procedure. Note that the first stack frame shown is almost always useless, since data either has not been saved yet, or is from a previous call. The last function preceding the **Ret Addr** is the function that called the procedure.

You can also use the **map** subcommand to locate the function name if the function was exported. The **map <addr>** command locates the symbol before the given address. The following is a concise view of the stack:

Low Addresses		Stack grows at this end.
Callee's stack -> 0 pointer 4 8 12-16 20	Back chain Saved CR Saved LR Reserved SAVED TOC	<---LINK AREA (callee)
Space for P1-P8 is always reserved	P1 ... Pn Callee's stack area	OUTPUT ARGUMENT AREA <---(Used by callee to construct argument <--- LOCAL STACK AREA
-8*nfprs-4*ngprs --> save	Caller's GPR save area max 19 words	(Possible word wasted for alignment.) Rfirst = R13 for full save R31
-8*nfprs -->	Caller's FPR save area max 18 dblwds	Ffirst = F14 for a full save F31
Caller's stack -> 0 pointer 4 8 12-16 20	Back chain Saved CR Saved LR Reserved Saved TOC	<---LINK AREA (caller)
Space for P1-P8 24 is always reserved	P1 ... Pn Caller's	INPUT PARAMETER AREA <---(Callee's input parameters found here. Is also caller's arg area.)



The following is a sample stack history with a break in the sample **demodd** kernel extension. The breakpoint was set at the start of the **demoread** routine at 0x1B127D4 (Beginning IAR). This was called from an instruction at 0x000824B0 (**Ret Addr**). This in turn is called by the instruction at address 0x00085F54 (**Ret Addr**), and so on.

The low values of the addresses (0x000824B0 and 0x00085F54) suggest that the instructions are in **/unix**. You can use the **crash** program and the **le** subcommand to determine the right kernel extension that is loaded in an address range.

```

0x1b127d4    beginning demoread in demodd
0x000824b0    .rdevread in /unix
0x00085f54    .cdev_rdwr in /unix

> stack
Beginning IAR: 0x01B127D4      Beginning Stack: 0x2FF97C28
Chain:0x2FF97C88  CR:0x24222082  Ret Addr:0x000824B0  TOC:0x000C5E78
P1:0x2003F800  P2:0x2003F800  P3:0x0000008C  P4:0x00000001
P5:0x01B11200  P6:0x00000000  P7:0x2FF97D38  P8:0x00000000
2FF97C60    00000203 00000000 2FF97CF8 2FF7FCD0  |....././...|
2FF97C70    29057E6B 00001000 2FF97DC0 018E8BE0  |)~k..../.}....|
2FF97C80    00FF0000 00000000 2FF97CD8 22222044  |....././." " D|
Returning to Stack frame at 0x2FF97C88
Press ENTER to continue or x to exit:

>

Chain:0x2FF97CD8  CR:0x22222044  Ret Addr:0x00085F54  TOC:0x00000000
P1:0x00000000  P2:0x018C41E0  P3:0x2FF97CF8  P4:0x2FF7FCC8
P5:0x000850E0  P6:0x00000000  P7:0xDEADBEEF  P8:0xDEADBEEF
2FF97CC0    DEADBEEF DEADBEEF 00000000 000BE4F8  |.....|
2FF97CD0    001E70F8 000BE7A4 2FF97D28 000BE5AC  |..p..../.)....|
Returning to Stack frame at 0x2FF97CD8
Press ENTER to continue or x to exit:

...
>

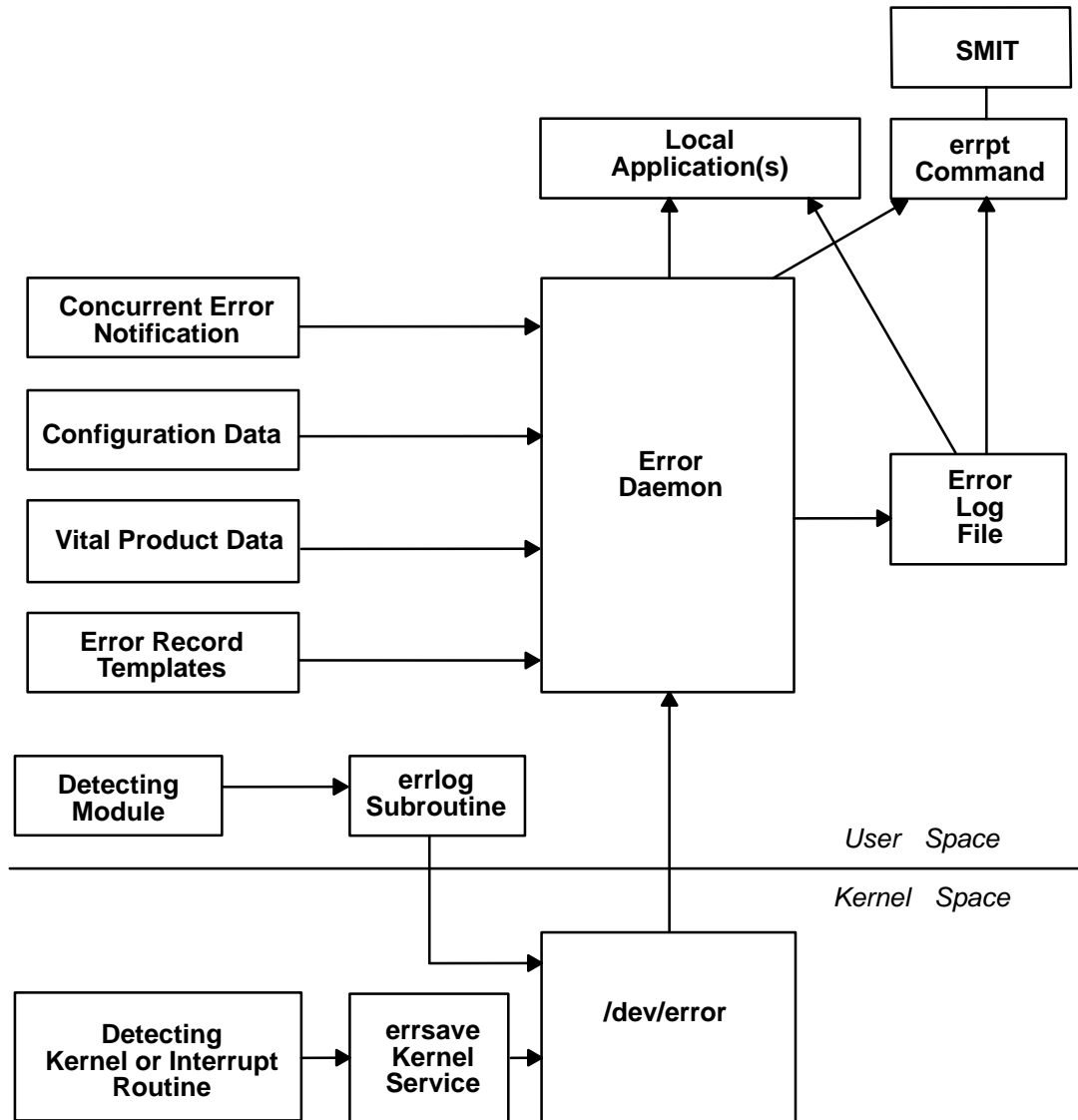
Chain:0x00000000  CR:0x22222022  Ret Addr:0x0000238C  TOC:0x00000000
P1:0x00000003  P2:0x30000000  P3:0x00000800  P4:0x00000000
P5:0x00000000  P6:0x00000000  P7:0x00000000  P8:0x00000000
Returning to Stack frame at 0x0
Press ENTER to continue or x to exit:
> Trace back complete.

```

Error Logging

The error facility allows a device driver to have entries recorded 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 special file **/dev/error**.

The **errdemon** daemon then 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. See the Flow of the Error Logging Facility figure on page 14-53 for an illustration of this.



Flow of the Error Logging Facility

Precoding Steps to Consider

Follow three precoding steps before initiating the error logging process. It is beneficial to understand what services are available to developers, and what the customer, service personnel, and defect personnel see.

Determine the Importance of the Error

The first precoding step is to review the error-logging documentation and determine whether a particular error should be logged. Do not use system resources for logging information that is unimportant or confusing to the intended audience.

It is, however, a worse mistake *not* to log an error that merits logging. You should work in concert 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

The next step is to determine the text of the message. Use the **errmsg** command with the **-w** flag to browse the system error messages file for a list of available messages. If you are developing a product for wide-spread general distribution and do not find a suitable system error message, you can submit a request to your supplier for a new message or follow the procedures that your organization uses to request new error messages. If your product is an in-house application, you can use the **errmsg** command to define a new message that meets your requirements.

Determine the Correct Level of Thresholding

Finally, determine the correct level of thresholding. Each error to be logged, regardless of whether it is a software or hardware error, 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 not unlimited in size. When its size limit is reached, the log wraps. If a particular error is repeated needlessly, existing information is overwritten, possibly causing 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 error log currently equals 1MB. As shipped, it cleans up any entries older than 30 days. In order to ensure that your error log entries are actually informative, noticed, and remain intact, *test your driver thoroughly*

Coding Steps

To begin error logging,

1. Select the error text.
2. Construct error record templates.
3. Add error logging calls into the device driver code.

Selecting the Error Text

The first task is to select the error text. After browsing the contents of the system message file, three possible paths exist for selecting the error text. Either all of the desired messages for the new errors exist in the message file, none of the messages exist, or a combination of errors exists.

- If the messages required already exist in the system message file, make a note of the four-digit hexadecimal identification number, as well as the message-set identification letter. For instance, a desired error description can be:

```
SET E
E859 "The wagon wheel is broken."
```

- If none of the system error messages meet your requirements, and if you are responsible for developing a product for wide-spread general distribution, you can either contact your supplier to allocate new messages or follow the procedures that your organization uses to request new messages. If you are creating an in-house product, use the **errmsg** command to write suitable error messages and use the **errinstall** command to install them. Refer to "Guidelines for Creating Software Product Installation Packages" in *AIX Version 4.1 General Programming Concepts* [Volume 1: Writing Programs](#) for more information. Take care not to overwrite other error messages.

- It is also possible to use a combination of existing messages and new messages within the same error record template definition.

Constructing Error Record Templates

The second step is to construct your *error record templates*. An error record template defines the text that appears in the error report. Each error record template has the following general form:

```

Error Record Template
+LABEL:
    Comment =
    Class =
    Log =
    Report =
    Alert =
    Err_Type =
    Err_Desc =
    Probable_Causes =
    User_Causes =
    User_Actions =
    Inst_Causes =
    Inst_Actions =
    Fail_Causes =
    Fail_Actions =
    Detail_Data = <data_len>, <data_id>, <data_encoding>

```

Each field in this stanza has well-defined criteria for input values. See the **errupdate** command for more information. The fields are:

Label	Requires a unique label for each entry to be added. The label must follow C language rules for identifiers and must not exceed 16 characters in length.						
Comment	Indicates this is a comment field. You must enclose the comment in double quotation marks; and it cannot exceed 40 characters.						
Class	Requires class values of H (hardware), S (software), or U (Undetermined).						
Log	Requires values True or False. If failure occurs, the errors are logged only if this field value is set to True. When this value is False the <code>Report</code> and <code>Alert</code> fields are ignored.						
Report	The values for this field are True or False. If the logged error is to be displayed using error report, the value of this field must be True.						
Alert	Set this field to True for errors that need to be forwarded to the Network Management Alert Manager program to generate an Alert. For errors that are not alertable, set this field to False.						
Err_Type	Describes the severity of the failure that occurred. Possible values are INFO, PEND, PERF, PERM, TEMP, and UNKN where: <table> <tr> <td>INFO</td> <td>The error log entry is informational and was not the result of an error.</td> </tr> <tr> <td>PEND</td> <td>A condition in which it is determined that the loss of availability of a device or component is imminent.</td> </tr> <tr> <td>PERF</td> <td>A condition in which the performance of a device or component was degraded below an acceptable level.</td> </tr> </table>	INFO	The error log entry is informational and was not the result of an error.	PEND	A condition in which it is determined that the loss of availability of a device or component is imminent.	PERF	A condition in which the performance of a device or component was degraded below an acceptable level.
INFO	The error log entry is informational and was not the result of an error.						
PEND	A condition in which it is determined that the loss of availability of a device or component is imminent.						
PERF	A condition in which the performance of a device or component was degraded below an acceptable level.						

	PERM	A permanent failure is defined as a condition that was not recoverable. For example, an operation was retried a prescribed number of times without success.
	TEMP	Recovery from this temporary failure was successful, yet the number of unsuccessful recovery attempts exceeded a predetermined threshold.
	UNKN	A condition in which it is not possible to assess the severity of a failure.
Err_Desc		Describes the failure that occurred. Proper input for this field is the four-digit hexadecimal identifier of the error description message to be displayed from SET E in the message file.
Prob_Causes		Describes one or more probable causes for the failure that occurred. You can specify a list of up to four Prob_Causes identifiers separated by commas. A Prob_Causes identifier displays a probable cause text message from SET P in the message file. List probable causes in the order of decreasing probability. At least one probable cause identifier is required.
User_Causes		Specifies a condition that an operator can resolve without contacting any service organization. You can specify a list of up to four User_Causes identifiers separated by commas. A User_Causes identifier displays a text message from SET U in the message file. List user causes in the order of decreasing probability. Leave this field blank if it does not apply to the failure that occurred. If this field is blank, either the Inst_Causes or the Fail_Causes field must not be blank.
User_Actions		Describes recommended actions for correcting a failure that resulted from a user cause. You can specify a list of up to four recommended User_Actions identifiers separated by commas. A recommended User_Actions identifier displays a recommended action text message, SET R in the message file. You must leave this field blank if the User_Causes field is blank. The order in which the recommended actions are listed is determined by the expense of the action and the probability that the action corrects the failure. Actions that have little or no cost and little or no impact on system operation should always be listed first. When actions for which the probability of correcting the failure is equal or nearly equal, list the least expensive action first. List remaining actions in order of decreasing probability.
Inst_Causes		Describes a condition that resulted from the initial installation or setup of a resource. You can specify a list of up to four Inst_Causes identifiers separated by commas. An Inst_Causes identifier displays a text message, SET I in the message file. List the install causes in the order of decreasing probability. Leave this field blank if it is not applicable to the failure that occurred. If this field is blank, either the User_Causes or the Failure_Causes field must not be blank.
Inst_Actions		Describes recommended actions for correcting a failure that resulted from an install cause. You can specify a list of up to four recommended Inst_actions identifiers separated by commas. A recommended Inst_actions identifier identifies a recommended action text message, SET R in the message file. Leave this field blank if the Inst_Causes field is blank. The order in which the recommended actions are listed is

determined by the expense of the action and the probability that the action corrects the failure. See the `User_Actions` field for the list criteria.

Fail_Causes Describes a condition that resulted from the failure of a resource. You can specify a list of up to four `Fail_Causes` identifiers separated by commas. A `Fail_Causes` identifier displays a failure cause text message, `SET F` in the message file. List the failure causes in the order of decreasing probability. Leave this field blank if it is not applicable to the failure that occurred. If you leave this field blank, either the `User_Causes` or the `Inst_Causes` field must not be blank.

Fail_Actions Describes recommended actions for correcting a failure that resulted from a failure cause. You can specify a list of up to four recommended action identifiers separated by commas. The `Fail_Actions` identifiers must correspond to recommended action messages found in `SET R` of the message file. Leave this field blank if the `Fail_Causes` field is blank. Refer to the description of the `User_Actions` field for criteria in listing these recommended actions.

Detail_Data Describes the detailed data that is logged with the error when the failure occurs. The `Detail_data` field includes the name of the detecting module, sense data, or return codes. Leave this field blank if no detailed data is logged with the error.

You can repeat the `Detail_Data` field. The amount of data logged with an error must not exceed the maximum error record length defined in the `sys/err_rec.h` header file. Save failure data that cannot be contained in an error log entry elsewhere, for example in a file. The detailed data in the error log entry contains information that can be used to correlate the failure data to the error log entry. Three values are required for each detail data entry:

data_len Indicates the number of bytes of data to be associated with the **data_id** value. The **data_len** value is interpreted as a decimal value.

data_id Identifies a text message to be printed in the error report in front of the detailed data. These identifiers refer to messages in `SET D` of the message file.

data_encoding

Describes how the detailed data is to be printed in the error report. Valid values for this field are:

ALPHA The detailed data is a printable ASCII character string.

DEC The detailed data is the binary representation of an integer value, the decimal equivalent is to be printed.

HEX The detailed data is to be printed in hexadecimal.

Sample Error Record Template

An example of an error record template is:

```
+ MISC_ERR:
  Comment = "Interrupt: I/O bus timeout or channel check"
  Class = H
  Log = TRUE
  Report = TRUE
  Alert = FALSE
  Err_Type = UNKN
  Err_Desc = E856
  Prob_Causes = 3300, 6300
  User_Causes =
  User_Actions =
  Inst_Causes =
  Inst_Actions =
  Fail_Causes = 3300, 6300
  Fail_Actions = 0000
  Detail_Data = 4, 8119, HEX      *IOCC bus number
  Detail_Data = 4, 811A, HEX     *Bus Status Register
  Detail_Data = 4, 811B, HEX     *Misc. Interrupt Register
```

Construct the error templates for all new errors to be added in a file suitable for entry with the **errupdate** command. Run the **errupdate** command with the **-h** flag and the input file. The new errors are now part of the error record template repository. A new header file is also created (**file.h**) in the same directory in which the **errupdate** command was run. This header file must be included in the device driver code at compile time. Note that the **errupdate** command has a built-in syntax checker for the new stanza that can be called with the **-c** flag.

Adding Error Logging Calls into the Code

The third step in coding error logging is to put the error logging calls into the device driver 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 <i>buf</i> 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;

    errsava(&log, (uint)sizeof(dderr)); /* run actual loggi
ng */
} /* end errlog_ex */

```

The data to be passed to the **errsava** 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.

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?

Writing to the **/dev/error** Special File

The error logging process begins when a loggable error is encountered and the device driver error logging subroutine sends the error information to the **errsave** kernel service. The error entry is written to the **/dev/error** special file. Once the information arrives at this file, it is time-stamped by the **errdemon** daemon and put in a buffer. The **errdemon** daemon constantly checks the **/dev/error** special file for new entries, and when new data is written, the daemon collects other information pertaining to the resource reporting the error. The **errdemon** daemon then creates an entry in the **/var/adm/ras/errlog** error logging file.

Performance Tracing

The AIX **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.

Care was taken in the design and implementation of this facility to make the collection of **trace** data efficient, 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.

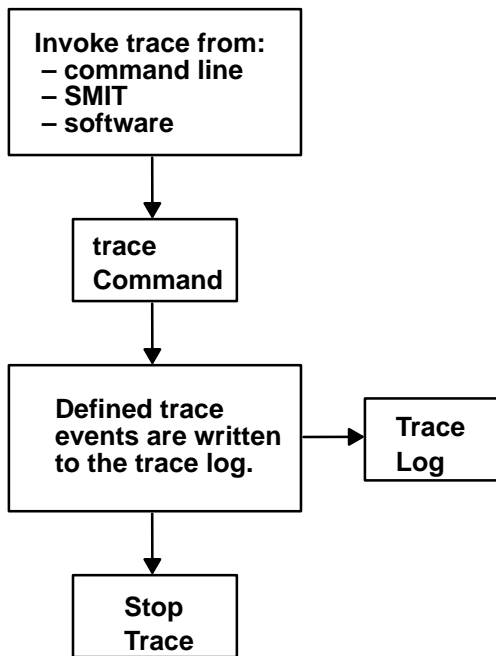
First, the **trace** facility can be triggered in multiple ways to capture small increments of system activity. It is practical to capture seconds to minutes of system activity in this way for post-processing. This is sufficient time to characterize major application transactions or interesting sections of a long task.

Second, the **trace** facility can be configured to direct the event stream to standard output. This allows a realtime 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.

You can start the system trace from:

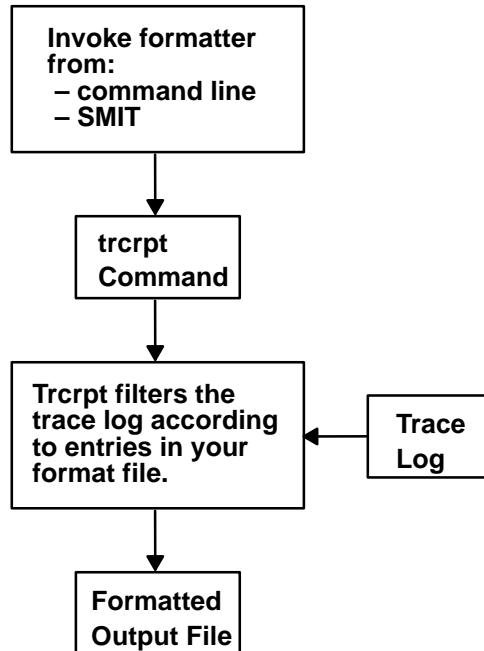
- The command line
- SMIT
- Software

As shown in the following Starting and Stopping Trace figure, 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” on page 14-65. Tracing from a software application is discussed and an example is presented in “Examples of Coding Events and Formatting Events” on page 14-83.



Starting and Stopping Trace

After a trace is started and stopped, you must format it before viewing it. This is illustrated in the following Trace Formatting figure. 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 14-73.



Trace Formatting

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 14-71.

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. The syntax of this command is:

```

trace [-a | -f | -l] [-c] [-d] [-h] [-j Event [,Event]] [-k Event [,Event]]
[-m Message] [-n] [-o OutName] [-o-] [-s] [-L Size] [-T Size] [-1234567]
  
```

The various options of the **trace** command are:

-f or -l Control the capture of trace data in system memory. If you specify neither the **-f** nor **-l** option, the trace facility creates two buffer areas in system memory to capture the trace data. These buffers are alternately written to the log file (or standard output if specified) as they become full. The **-f** or **-l** flag provides you with the ability to prevent data from being written to the file

during data collection. The options are to collect data only until the memory buffer becomes full (**-f** for first), or to use the memory buffer as a circular buffer that captures only the last set of events that occurred before **trace** was terminated (**-l**). The **-f** and **-l** options are mutually exclusive. With either the **-f** or **-l** option, data is not transferred from the memory collection buffers to file until **trace** is terminated.

- a** Run the **trace** collection asynchronously (as a background task), returning a normal command line prompt. Without this option, the **trace** command runs in a subcommand mode (similar to the **crash** command) and returns a **>** prompt. You can issue subcommands and regular shell commands from the **trace** subcommand mode by preceding the shell commands with an **!** (exclamation point).
- c** Saves the previous trace log file adding **.old** to its name. Generates an error if a previous trace log file does not exist. When using the **-o Name** flag, the user-defined trace log file is renamed.
- d** Delay data collection. The trace facility is only configured. Data collection is delayed until one of the collection trigger events occurs. Various methods for triggering data collection on and off are provided. These include the following:
 - **trace** subcommands
 - **trace** commands
 - **ioctl**s to **/dev/systrctl**.
- j events or -k events** Specify a set of events to include (**-j**) or exclude (**-k**) from the collection process. Specify a list of events to include or exclude by a series of three-digit hexadecimal event IDs separated by a space.
- s** Terminate **trace** data collection if the **trace** log file reaches its maximum specified size. The default without this option is to wrap and overwrite the data in the log file on a FIFO basis.
- h** Do not write a **date/sysname/message** header to the **trace** log file.
- m message** Specify a text string (message) to be included in the **trace** log header record. The message is included in reports generated by the **trcrpt** command.
- n** Adds some information to the trace log header: lock information, hardware information, and, for each loader entry, the symbol name, address, and type.
- o outfile** Specify a file to use as the log file. If you do not use the **-o** option, the default log file is **/usr/adm/ras/trcfile**. To direct the trace data to standard output, code the **-o** option as **-o -**. (When **-o-** is specified the **-c** flag is ignored.) Use this technique only to pipe the data stream to another process since the trace data contains raw binary events that are not displayable.
- 1234567** Duplicate the **trace** design for multiple channels. Channel 0 is the default channel and is always used for recording system events. The other channels are generic channels, and their use is not predefined. There are various uses of generic channels in the system. The generic channels are also available to user applications. Each created channel is a separate events data stream. Events recorded to channel 0 are mixed with the

predefined system events data stream. The other channels have no predefined use and are assigned generically.

A program can request that a generic channel be opened by using the **trcstart** subroutine. A channel number is returned, similar to the way a file descriptor is returned when it opens a file. The program can record events to this channel and, thus, have a private data stream. The **trace** command allows a generic channel to be specifically configured by defining the channel number with this option. However, this is not generally the way a generic channel is started. It is more likely to be started from a program using the **trcstart** subroutine, which uses the returned channel ID to record events.

-T size and -L size

Specify the size of the collection memory buffers and the maximum size of the log file in bytes. The trace facility pins the data collection buffers, making this amount of memory unavailable to the rest of the system. It is important to be aware of this, because it means that the trace facility can impact performance in a memory constrained environment. If the application being monitored is not memory constrained, or if the percentage of memory consumed by the trace routine is small compared to what is available in the system, the impact of **trace** "stolen" memory should be small.

If you do not specify a value, trace uses a default size. The trace facility pins a little more than the specified buffer size. This additional memory is required for the trace facility itself. Trace pins a little more than the amount specified for first buffer mode (**-f** option). Trace pins a little more than twice the amount specified for trace configured in alternate buffer or last (circular) buffer mode.

You can also start **trace** from a 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 the sample code on page 14-67.

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

Once **trace** is configured by the **trace** command or the **trcstart** subroutine, controls to **trace** trigger the collection of data on, trigger the collection of data off, and stop the trace facility (stop deconfigures **trace** and unpins buffers). These basic controls exist as subcommands, commands, subroutines, and ioctl controls to the **trace** control device, **/dev/systctl**. These controls are described in the following sections.

Controlling trace in Subcommand Mode

If the **trace** routine is configured without the **-a** option, it runs in subcommand mode. Instead of the normal shell prompt, **->** is the prompt. In this mode the following subcommands are recognized:

- trcon** Triggers collection of **trace** data on.
- trcoff** Triggers collection of **trace** data off.
- q or quit** Stops collection of **trace** data (like **trcoff**) and terminates **trace** (deconfigures).
- !command** Runs the specified shell command.

The following is an example of a trace session in which the trace subcommands are used. First, the system trace points have been displayed. Second, a trace on the system calls have been selected. Of course, you can trace on more than one trace point. Be aware that trace takes a lot of data. Only the first few lines are shown in the following example:

```
# trcrpt -j |pg
004     TRACEID IS ZERO
100     FLIH
200     RESUME
102     SLIH
103     RETURN FROM SLIH
101     SYSTEM CALL
104     RETURN FROM SYSTEM CALL
106     DISPATCH
10C     DISPATCH IDLE PROCESS
11F     SET ON READY QUEUE
134     EXEC SYSTEM CALL
139     FORK SYSTEM CALL
107     FILENAME TO VNODE (lookuppn)
15B     OPEN SYSTEM CALL
130     CREAT SYSTEM CALL
19C     WRITE SYSTEM CALL
163     READ SYSTEM CALL
10A     KERN_PFS
10B     LVM BUF STRUCT FLOW
116     XMALLOC size,align,heap
117     XMFREE address,heap
118     FORKCOPY
11E     ISSIG
169     SBREAK SYSTEM CALL
```

```
# trace -d -j 101 -m "system calls trace example"
-> trcon
-> !cp /tmp/xbugs .
-> trcoff
-> quit
# trcrpt -O "exec=on,pid=on" > cp.trace
# pg cp.trace
pr 3 11:02:02 1991
System: AIX smiller Node: 3
Machine: 000247903100
Internet Address: 00000000 0.0.0.0
system calls trace example
trace -d -j 101 -m -m system calls trace example
```

ID	PROCESS NAME	PID	I	ELAPSED_SEC	DELTA_MSEC	APPL	SYSCALL
001	trace	13939		0.000000000	0.000000		TRACE ON chan 0
101	trace	13939		0.000251392	0.251392		kwritev
101	trace	13939		0.000940800	0.689408		sigprocmask
101	trace	13939		0.001061888	0.121088		kreadv

101 trace	13939	0.001501952	0.440064	kreadv
101 trace	13939	0.001919488	0.417536	kiocctl
101 trace	13939	0.002395648	0.476160	kreadv
101 trace	13939	0.002705664	0.310016	kiocctl

Controlling the trace Facility by Commands

If you configure the **trace** routine to run asynchronously (the **-a** option), you can control the trace facility with the following commands:

trcon Triggers collection of trace data on.

trcoff Triggers collection of trace data off.

trcstop Stops collection of trace data (like **trcoff**) and terminates the **trace** routine.

Controlling the trace Facility by Subroutines

The controls for the **trace** routine are available as subroutines from the **librts.a** library. The subroutines return zero on successful completion. The subroutines are:

trcon Triggers collection of **trace** data on.

trcoff Triggers collection of **trace** data off.

trcstop Stops collection of **trace** data (like **trcoff**) and terminates the **trace** routine.

Controlling the trace Facility with ioctl Calls

The subroutines for controlling **trace** open the trace control device (**/dev/systrctl**), issue the appropriate **ioctl** command, close the control device and return. To control tracing around sections of code, it can be more efficient for a program to issue the **ioctl** controls directly. This avoids the unnecessary, repetitive opening and closing of the trace control device, at the expense of exposing some of the implementation details of **trace** control. To use the **ioctl** call in a program, include **sys/trcctl.h** to define the **ioctl** commands. The syntax of the **ioctl** is as follows:

```
ioctl (fd, CMD, Channel)
```

where:

fd File descriptor returned from opening **/dev/systrctl**

CMD TRCON, TRCOFF, or TRCSTOP

Channel Trace channel (0 for system trace).

The following code sample shows how to start a **trace** from a program and only trace around a specified section of code:

```
#include <sys/trcctl.h>
extern int trcstart(char *arg);
char *ctl_dev = "/dev/systrctl";
int ctl_fd
main()
{
    printf("configuring trace collection \n");
    if (trcstart("-ad")){
        perror("trcstart");
        exit(1);
    }
    if((ctl_fd = open (ctl_dev))<0){
        perror("open ctl_dev");
        exit(1);
    }
}
```



```

printf("turning trace collection on \n");
if(ioctl(ctl_fd,TRCON,0)){
    perror("TRCON");
    exit(1);
}
/* code between here and trcoff ioctl will be traced */
printf("turning trace off\n");
if (ioctl(ctl_fd,TRCOFF,0)){
    perror("TRCOFF");
    exit(1);
}
exit(0);
}

```

Producing a trace Report

A 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.

The trcrpt Command

The syntax of the **trcrpt** command is as follows:

```

trcrpt [-c] [-d List] [-e Date] [-h] [-j] [-k List] [-n Name] [-o File] [-p List]
[-q] [-r] [-s Date] [-t File] [-v] [-O Options] [-T List] [LogFile]

```

Normally the **trcrpt** output goes to standard output. However, it is generally more useful to redirect the report output to a file. The options are:

- c** Causes the **trcrpt** command to check the syntax of the trace format file. The trace format file checked is either the default (**/etc/trcfmt**) or the file specified by the **-t** flag with this command. You can check the syntax of the new or modified format files with this option before attempting to use them.
- d List** Allows you to specify a list of events to be included in the **trcrpt** output. This is useful for eliminating information that is superfluous to a given analysis and making the volume of data in the report more manageable. You may have commonly used event profiles, which are lists of events that are useful for a certain type of analysis.
- e Date** Ends the report time with entries on, or before the specified date. The *Date* parameter has the form *mmddhhmmssyy* (month, day, hour, minute, second, and year). Date and time are recorded in the trace data only when trace data collection is started and stopped. If you stop and restart trace data collection multiple times during a trace session, date and time are recorded each time you start or stop a trace data collection. Use this flag in combination with the **-s** flag to limit the trace data to data collected during a certain time interval.
- h** Omit the column headings of the report.

- j** Causes the **trcrpt** command to produce a list of all the defined events from the specified trace format file. This option is useful in creating an initial file that you can edit to use as an include or exclude list for the **trcrpt** or **trace** command.
- k List** Similar to the **-d** flag, but allows you to specify a list of events to exclude from the **trcrpt** output.
- n Name** Specifies the kernel name list file to be used by **trcrpt** to convert kernel addresses to routine names. If not specified, the report facility uses the symbol table in **/unix**. A kernel name list file that matches the system the data was collected on is necessary to produce an accurate trace report. You can create such a file for a given level of system with the **trcnm** command:

`trcnm /unix > Name`
- o File** Writes the report to a file instead of to standard output.
- p List** Limits the **trcrpt** output to events that occurred during the running of specific processes. List the processes by process name or process ID.
- q** Suppresses detailed output of syntax error messages. This is not an option you typically use.
- r** Produces a raw binary format of the trace data. Each event is output as a record in the order of occurrence. This is not necessarily the order in which the events are in the trace log file since the logfile can wrap. If you use this option, direct the output to a file (or process), since the binary form of the data is not displayable.
- t File** Allows you to specify a trace format file other than the default (**/etc/trcfmt**).
- T List** Limits the report to the kernel thread IDs specified by the *List* parameter. The list items are kernel thread IDs separated by commas. Starting the list with a kernel thread ID limits the report to all kernel thread IDs *in* the list. Starting the list with a ! (exclamation point) followed by a kernel thread ID limits the report to all kernel thread IDs *not in* the list.
- O options** Allows you to specify formatting options to the **trcrpt** command in a comma separated list. Do not put spaces after the commas. These options take the form of option=selection. If you do not specify a selection, the command uses the default selection. The possible options are discussed in the following sections. Each option is introduced by showing its default selection.
 - 2line=off** This option lets the user specify whether the lines in the event report are split and displayed across two lines. This is useful when more columns of information have been requested than can be displayed on the width of the output device.
 - endtime=nnn.nnnnnnnnn** The **starttime** and **endtime** option permit you to specify an elapsed time interval in which the **trcrpt** produces output. The elapsed time interval is specified in seconds with nanosecond resolution.
 - exec=off** Lets you specify whether a column showing the path name of the current process is displayed. This is useful in showing what process (by name) was active at the time of the event. You typically want to specify this option. We recommend that you specify **exec=on** and **pid=on**.

ids=on Lets you specify whether to display a column that contains the event IDs. If the selection is on, a three-digit hex ID is shown for each event. The alternate is off.

pagesize=0 Lets you specify how often the column headings is reprinted. The default selection of 0 displays the column headings initially only. A selection of 10 displays the column heading every 10 lines.

pid=off Lets you specify whether a column showing the process ID of the current process is displayed. It is useful to have the process ID displayed to distinguish between several processes with the same executable name. We recommend that you specify `exec=on` and `pid=on`.

starttime=nnn.nnnnnnnnn

The **starttime** and **endtime** option permit you to specify an elapsed time interval in which the **trcrpt** command produces output. The elapsed time interval is specified in seconds with nanosecond resolution.

svc=off Lets you specify whether the report should contain a column that indicates the active system call for those events that occur while a system call is active.

tid=off Displays the kernel thread IDs in the trace report. Values for this option are **on** and **off**. The default value is **off**.

timestamp=0

The report can contain two time columns. One column is elapsed time since the **trace** command was initiated. The other possible time column is the delta time between adjacent events. The option controls if and how these times are displayed. The selections are:

- 0** Provides both an elapsed time from the start of **trace** and a delta time between events. The elapsed time is shown in seconds and the delta time is shown in milliseconds. Both fields show resolution to a nanosecond. This is the default value.
- 1** Provides only an elapsed time column displayed as seconds with resolution shown to microseconds.
- 2** Provides both an elapsed time and a delta time column. The elapsed time is shown in seconds with nanosecond resolution, and delta time is shown in microseconds with microsecond resolution.
- 3** Omits all time stamps from the report.

logfile The **logfile** is the name of the file that contains the event data to be processed by the **trcrpt** command. The default is the `/usr/adm/ras/trcfile` file.

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 may 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 a

- Hookword
- Data words (optional)
- A TID, or thread identifier
- Timestamp (optional).

The following Format of a Trace Event Record figure illustrates a trace event. A four-bit type is defined for each form the event record can take. The type field is imposed by the recording routine so that the report facility can always skip from event to event when processing the data, even if the formatting rules in the trace format file are incorrect or missing for that event.

12 bit Hook ID	4 bit Type	16 bit Data Field
D1 Optional Data Word 1		
D2 Optional Data Word 2		
D3 Optional Data Word 3		
D4 Optional Data Word 4		
D5 Optional Data Word 5		
TID (Thread ID)		
Optional Time Stamp		

Format of a Trace Event Record

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 of 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

There is a macro to record each possible type of event record. The macros are defined in the **sys/trcmacros.h** header file. The 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)
- **TRCHKL5T** (hw,D1,D2,D3,D4,D5).

Similarly, to record non-time stamped system events (channel 0), 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).

There are only two macros to record events to one of the generic channels (channels 1–7). These 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).

Use of Event IDs (hookids)

Event IDs are 12 bits (or 3-digit hexadecimal), for a possibility of 4096 IDs. Event IDs that are permanently left in and shipped with code need to be permanently assigned by IBM. Permanently assigned event IDs are defined in the **sys/trchkid.h** header file.

To allow you to define events in your environments or during development, a range of event IDs exist for temporary use. The range of event IDs for temporary use is hex 010 through hex 0FF. 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. 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" on page 14-73. 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 pathlength. However, this would generally be an excessive level of instrumentation to ship for a component.

We suggest that you 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 which 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:

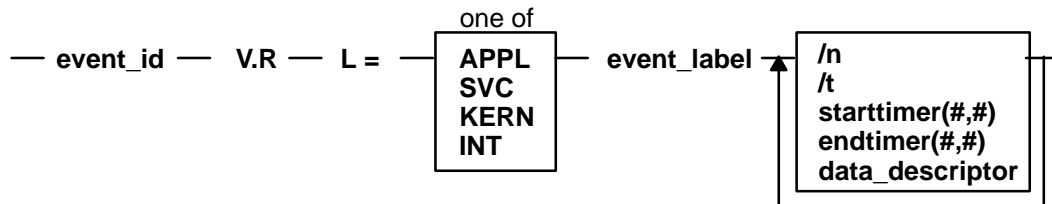
- A trace ID can be used for a group of events by “switching” on one of the data fields. This means that a particular data field can be used to identify from where the trace point was called. The trace format routine can be made to format the trace data for that unique trace point.
- The trace hook is the level at which a group of events can be enabled or disabled.

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.

Refer to the `/etc/tcrfmt` file to see examples of the syntax for stanzas that appear in the trace format file.



Syntax of a Stanza in the Format File

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) character. The fields are:

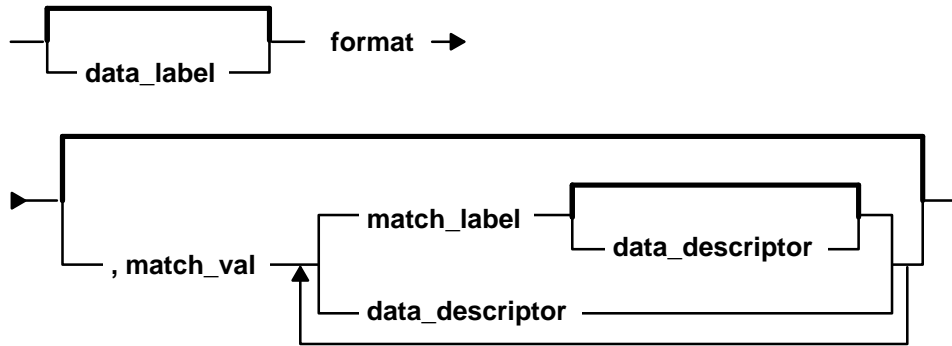
- event_id** Each stanza begins with the three-digit hexadecimal event ID that the stanza describes, followed by a space.
- V.R** This field describes the version (V) and release (R) that the event was first assigned. Any integers work for V and R, and you may want to keep your own tracking mechanism.
- 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 `tcrpt` 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 following Syntax of the `data_descriptor` Field figure illustrates this field's syntax.



Syntax of the data_descriptor Field

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 Review the format of an event record depicted in the figure Format of a trace Event Record. 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
Am.n	This value specifies that <code>m</code> bytes of data are consumed as ASCII text, and that it is displayed in an output field that is <code>n</code> characters wide. The data pointer is moved <code>m</code> bytes.
S1, S2, S4	Left justified string. The length of the field is defined as 1 byte (S1), 2 bytes (S2), or 4 bytes (S4). The data pointer is moved accordingly.
Bm.n	Binary data of <code>m</code> bytes and <code>n</code> bits. The data pointer is moved accordingly.
Xm	Hexadecimal data of <code>m</code> bytes. The data pointer is moved accordingly.
D2, D4	Signed decimal format. Data length of 2 (D2) bytes or 4 (D4) bytes is consumed.
U2, U4	Unsigned decimal format. Data length of 2 or 4 bytes is consumed.
F4, F8	Floating point of 4 or 8 bytes.
Gm.n	Positions the data pointer. It specifies that the data pointer is positioned <code>m</code> bytes and <code>n</code> bits into the data.
Om.n	Skip or omit data. It omits <code>m</code> bytes and <code>n</code> bits.
Rm	Reverse the data pointer <code>m</code> 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:

```
$D1%B2.3
```

This macro causes data word one to be accessed, but to be considered as 2 bytes and 3 bits of binary data.

\$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

```
# Licensed Materials - Property of IBM
#
# US Government Users Restricted Rights - Use, duplication or
# disclosure restricted by GSA ADP Schedule Contract with IBM Corp.

# I. General Information
#
# A. Binary format for the tracehook calls. (1 column = 4 bits)
#   trchk      MMmTDDDD
#   trchkt     MMmTDDDDttttttttt
#   trchkkl    MMmTDDDD11111111
#   trchklt    MMmTDDDD11111111tttttttt
#   trchkg     MMmTDDDD1111111122222222333333334444444455555555
#   trchkg     MMmTDDDD1111111122222222333333334444444455555555tttttttt
#   trcgen     MMmTLLLL11111111vvvvvvvvvvvvvvvvvvvvvvvvvvvvvvvvvxxxxx
#   trcgent    MMmTLLLL11111111vvvvvvvvvvvvvvvvvvvvvvvvvvvvvvvvvxxxxxtttttttt
#
#       legend:
#   MM = major id
#   m  = minor id
#   T  = hooktype
#   D  = hookdata
#   t  = nanosecond timestamp
#   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() type (6,7) is an exception.The DATA_POINTER starts at
#   the fifth byte, ie., at the 'd1' parameter 11111111.
#
# B. 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" ... \
#
# C. 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.
#
# 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
#
# S1, S2, S4
# Left justified ascii string.
# The length of the string is in the first byte(half-word, word)
# of the data. This length of the string does not include this byte.
# The data pointer is advanced by the length value.
#      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
#   The DATA_POINTER is advanced by m.
#
# D2, D4
#   Signed decimal format.
#   The length of the data is 2 (4) bytes.
#   The DATA_POINTER is advanced by 2 (4).
#
# U2, U4
#   Unsigned decimal format.
#   The length of the data is 2 (4) bytes.
#   The DATA_POINTER is advanced by 2 (4).
#
# 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.
#   This is also equal to the 16 bit hookdata.
#   The DATA_POINTER is not changed.
#
# HT
#   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 trchkx() calls with
#   the same template.
#   The DATA_POINTER is not changed.
#
# C. Codes that interpret the data in some way before output.
# T4
#   Output the next 4 bytes as a data and time string,
#   in GMT timezone format. (as in ctime(&seconds))
#   The DATA_POINTER is advanced by 4.
#
# E1,E2,E4
#   Output the next byte (half_word, word) 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, or 4.
#
# P4
#     Use the next word 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.
#
# \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.
#     The DATA_POINTER advances over the \t.
#
# \n
#     Output a newline. \n\n\n outputs 3 newlines.
#     The newline is left-justified according to the INDENTATION LEVEL.
#     The DATA_POINTER advances over the \n.
#
# $macro
#     The value of 'macro' is output as a %04X value. Undefined
#     macros have the value of 0000.
#     The DATA_POINTER is not changed.
#     An optional format can be used with macros:
#     $v1%X4     will output the value $v1 in X4 format.
#     $zz%B0.8   will output the value $v1 in 8 bits of binary.
#     Understood formats are: X, D, U, B. Others default to X2.
#
# "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 is selected to as 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 25 macros can be defined per template.
#
# VI. Special macros:
# $RELLINENO line number for this event. The first line starts at 1.
# $D1 - $D5 dataword 1 through dataword 5. No change to datapointer.
# $HD hookdata (lower 16 bits)
# $SVC Output the name of the current SVC
# $EXECPTH Output the pathname of the executable for current process.
# $PID Output the current process id.
# $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
# $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 may have different starting points into an event.
#
# 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 is 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.
#
# 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', since 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
#
# VII. 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, since it look
# 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 0, the number is octal. 0x or 0X makes the
# number hex.
# A 000 traceid will use this template
# This id is also used to define most of the template functions.
# filemode(omode) expand omode the way ls -l does. The
#   call to setdelim() inhibits spaces between the chars.
#

```

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              ctlfd;
int              i;

main()
{
    printf("configuring trace collection \n");
    if (trcstart("-ad")){
        perror("trcstart");
        exit(1);
    }
    if((ctlfd = open(ctl_file,0))<0){
        perror(ctl_file);
        exit(1);
    }
    printf("turning  trace on \n");
    if(ioctl(ctlfd,TRCON,0)){
        perror("TRCON");
        exit(1);
    }
    /* here is the code that is being traced */
    for(i=1;i<11;i++){
        TRCHKL1T(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(ioctl(ctlfd,TRCSTOP,0)){
        perror("TRCOFF");
        exit(1);
    }
    exit(0);
}
```

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
```

You must have root privilege if you use the default file to collect the trace information (*/usr/adm/ras/trcfile*).

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.

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 -O 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

Include several optional columns of data in the trace output. This causes the output to exceed 80 columns. It is best to view the reports on an output device that supports 132 columns.

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
```

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 100KB 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 -O "exec=on,pid=on" > cp.rpt
```

In the **cp.rpt** file you can see the following activities:

- 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.

The trace output looks a little overwhelming at first. This is a good example to use as a learning aid. If you can discern the activities described, you are well on your way to being able to use the trace facility to diagnose system performance problems.

Effective Filtering of the trace Report

The full detail of the trace data may 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 -O "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 -O "exec=on"
```

This command shows only the opens performed by the **cp** process.

Appendix A. New Interfaces

This appendix contains descriptions of the following new services that can be used by device driver writers but are not documented elsewhere:

- **d_map_clear** kernel service
- **d_map_disable** kernel service
- **d_map_enable** kernel service
- **d_map_init** kernel service
- **d_map_list** kernel service
- **d_map_page** kernel service
- **d_map_slave** kernel service
- **d_unmap_list** kernel service
- **d_unmap_page** kernel service
- **d_unmap_slave** kernel service
- **iomem_att** kernel service
- **iomem_det** kernel service
- **ns_alloc** network service
- **ns_free** network service
- **rmalloc** kernel service
- **rmfree** kernel service

d_map_clear Kernel Service

Purpose

Deallocates resources previously allocated on a **d_map_init** call.

Syntax

```
#include <sys/dma.h>

void d_map_clear (*handle)
struct d_handle *handle
```

Parameters

handle Indicates the unique handle returned by the **d_map_init** kernel service.

Description

The **d_map_clear** kernel service is a bus specific utility routine determined by the **d_map_init** service that deallocates resources previously allocated on a **d_map_init** call. This includes freeing the **d_handle** structure that was allocated by **d_map_init**.

Note: You can use the **D_MAP_CLEAR** macro provided in the **/usr/include/sys/dma.h** file to code calls to the **d_map_clear** kernel service.

Implementation Specifics

The **d_map_clear** kernel service is part of the base device package of your platform.

Related Information

The **d_map_init** kernel service.

d_map_disable Kernel Service

Purpose

Disables DMA for the specified handle.

Syntax

```
#include <sys/dma.h>
int d_map_disable(*handle)
struct d_handle *handle;
```

Parameters

handle Indicates the unique handle returned by **d_map_init**.

Description

The **d_map_disable** kernel service is a bus specific utility routine determined by the **d_map_init** kernel service that disables DMA for the specified *handle* with respect to the platform.

Note: You can use the **D_MAP_DISABLE** macro provided in the */usr/include/sys/dma.h* file to code calls to the **d_map_disable** kernel service.

Return Values

DMA_SUCC Indicates the DMA is successfully disabled

DMA_FAIL Indicates the DMA could not be explicitly disabled for this device or bus.

Implementation Specifics

The **d_map_disable** kernel service is part of the base device package of your platform.

Related Information

The **d_map_init** kernel service.

d_map_enable Kernel Service

Purpose

Enables DMA for the specified handle.

Syntax

```
#include <sys/dma.h>

int d_map_enable(*handle)
struct d_handle *handle;
```

Parameters

handle Indicates the unique handle returned by **d_map_init**.

Description

The **d_map_enable** kernel service is a bus specific utility routine determined by the **d_map_init** kernel service that enables DMA for the specified *handle* with respect to the platform.

Note: You can use the **D_MAP_ENABLE** macro provided in the `/usr/include/sys/dma.h` file to code calls to the **d_map_enable** kernel service.

Return Values

DMA_SUCC Indicates the DMA is successfully enabled

DMA_FAIL Indicates the DMA could not be explicitly enabled for this device or bus.

Implementation Specifics

The **d_map_enable** kernel service is part of the base device package of your platform.

Related Information

The **d_map_init** kernel service.

d_map_init Kernel Service

Purpose

Allocates and initializes resources for performing DMA with PCI and ISA devices.

Syntax

```
#include <sys/dma.h>

struct d_handle* d_map_init (bid, flags, bus_flags, channel)
int bid;
int flags;
int bus_flags;
uint channel;
```

Parameters

<i>bid</i>	Specifies the bus identifier.
<i>flags</i>	Describes the mapping.
<i>bus_flags</i>	Specifies the target bus flags.
<i>channel</i>	Indicates the <i>channel</i> assignment specific to the bus.

Description

The **d_map_init** kernel service allocates and initializes resources needed for managing DMA operations and returns a unique *handle* to be used on subsequent DMA service calls. The *handle* is a pointer to a **d_handle** structure allocated by **d_map_init** from the pinned heap for the device. The device driver uses the function addresses provided in the *handle* for accessing the DMA services specific to its host bus. The **d_map_init** service returns a **DMA_FAIL** error when resources are unavailable or cannot be allocated.

The *channel* parameter is the assigned channel number for the device, if any. Some devices and or buses might not have the concept of *channels*. For example, an ISA device driver would pass in its assigned DMA channel in the *channel* parameter.

Execution Environment

The **d_map_init** kernel service should only be called from the process environment.

Return Values

DMA_FAIL	Indicates that the resources are unavailable. No registration was completed.
struct d_handle *	Indicates successful completion.

Implementation Specifics

The **d_map_init** kernel service is part of the Base Operating System (BOS) Runtime.

Related Information

The **d_map_clear** kernel service, **d_map_page** kernel service, **d_unmap_page** kernel service, **d_map_list** kernel service, **d_unmap_list** kernel service, **d_map_slave** kernel service, **d_unmap_slave** kernel service, **d_map_disable** kernel service, **d_map_enable** kernel service.

d_map_list Kernel Service

Purpose

Performs platform-specific DMA mapping for a list of virtual addresses.

Syntax

```
#include <sys/dma.h>
```

```
int d_map_list (*handle, flags, minxfer, *virt_list, *bus_list)
struct d_handle *handle;
int flags;
int minxfer;
int *virt_list;
int *bus_list;
```

Parameters

<i>handle</i>	Indicates the unique handle returned by the d_map_init kernel service.
<i>flags</i>	Specifies one of the following flags: DMA_READ Transfers from a device to memory. BUS_DMA Transfers from one device to another device. DMA_BYPASS Do not check page access.
<i>minxfer</i>	Specifies the minimum transfer size for the device.
<i>virt_list</i>	Specifies a list of virtual buffer addresses and lengths.
<i>bus_list</i>	Specifies a list of bus addresses and lengths.

Description

The **d_map_list** kernel service is a bus specific utility routine determined by the **d_map_init** kernel service that accepts a list of virtual addresses and sizes and provides the resulting list of bus addresses. This service fills out the corresponding bus address list for use by the device in performing the DMA transfer. This service allows for scatter/gather capability of a device and also allows the device to combine multiple requests that are contiguous with respect to the device. The lists are passed via the **dio** structure. If the **d_map_list** service is unable to complete the mapping due to exhausting the capacity of the provided **dio** structure, an error, **DMA_DIOFULL**, is returned. If the **d_map_list** service is unable to complete the mapping due to exhausting resources required for the mapping, an error, **DMA_NORES**, is returned. In both of these cases, the *bytes_done* field of the **dio** virtual list is set to the number of bytes successfully mapped. This byte count is a multiple of the *minxfer* size for the device as provided on the call to **d_map_list**. The *resid_iov* field is set to the index of the remaining *d_iovec* fields in the list. Unless the **DMA_BYPASS** flag is set, this service verifies access permissions to each page. If an access violation is encountered on a page with the list, an error, **DMA_NOACC**, is returned, and the *bytes_done* field is set to the number of bytes preceding the faulting *iovec*.

Notes:

1. When the **DMA_NOACC** return value is received, no mapping is done, and the bus list is undefined. In this case, the *resid_iov* field is set to the index of the *d_iovec* that encountered the access violation.

2. You can use the **D_MAP_LIST** macro provided in the `/usr/include/sys/dma.h` file to code calls to the **d_map_list** kernel service.

Return Values

DMA_NORES Indicates that resources were exhausted during mapping.

DMA_DIOFULL Indicates that the target bus list is full.

DMA_NOACC Indicates no access permission to a page in the list.

DMA_SUCC Indicates that the entire transfer successfully mapped.

Implementation Specifics

The **d_map_list** kernel service is part of the base device package of your platform.

Related Information

The **d_map_init** kernel service.

d_map_page Kernel Service

Purpose

Performs platform specific DMA mapping for a single page.

Syntax

```
#include <sys/dma.h>
#include <sys/xmem.h>

int d_map_page(*handle, flags, *baddr, *busaddr, *xmp)
struct d_handle *handle;
int flags;
caddr_t *baddr;
uint *busaddr;
struct xmem *xmp;
```

Parameters

<i>handle</i>	Indicates the unique handle returned by the d_map_init kernel service.
<i>flags</i>	Specifies one of the following flags: DMA_READ Transfers from a device to memory. BUS_DMA Transfers from one device to another device. DMA_BYPASS Do not check page access.
<i>baddr</i>	Specifies the buffer address.
<i>busaddr</i>	Points to the <i>busaddr</i> field.
<i>xmp</i>	Cross-memory descriptor for the buffer.

Description

The **d_map_page** kernel service is a bus specific utility routine determined by the **d_map_init** kernel service that performs platform specific mapping of a single 4K or less transfer for DMA master devices. The **d_map_page** kernel service is a fast path version of the **d_map_list** service. The entire transfer amount must fit within a single page in order to use this service. This service accepts a virtual address and completes the appropriate bus address for the device to use in the DMA transfer. Unless the **DMA_BYPASS** flag is set, this service also verifies access permissions to the page.

If the buffer is a global kernel space buffer, the cross-memory descriptor can be set to point to the exported **GLOBAL** cross-memory descriptor, *xmem_global*.

If the transfer is unable to be mapped due to resource restrictions, the **d_map_page** service returns **DMA_NORES**. If the transfer is unable to be mapped due to page access violations, this service returns **DMA_NOACC**.

Note: You can use the **D_MAP_PAGE** macro provided in the */usr/include/sys/dma.h* file to code calls to the **d_map_page** kernel service.

Return Values

DMA_NORES Indicates that resources are unavailable.
DMA_NOACC Indicates no access permission to the page.

DMA_SUCC Indicates that the *busaddr* parameter contains the bus address to use for the device transfer.

Implementation Specifics

The **d_map_page** kernel service is part of the base device package of your platform.

Related Information

The **d_map_init** kernel service, **d_map_list** kernel service.

d_map_slave Kernel Service

Purpose

Accepts a list of virtual addresses and sizes and sets up the slave DMA controller.

Syntax

```
#include <sys/dma.h>

int d_map_slave (*handle, flags, minxfer*vlist, chan_flag)
struct d_handle *handle;
int flags;
int minxfer;
struct dio *vlist;
uint chan_flag;
```

Parameters

<i>handle</i>	Indicates the unique handle returned by the d_map_init kernel service.
<i>flags</i>	Specifies one of the following flags. DMA_READ Transfers from a device to memory. BUS_DMA Transfers from one device to another device. DMA_BYPASS Do not check page access.
<i>minxfer</i>	Specifies the minimum transfer size for the device.
<i>vlist</i>	Specifies a list of buffer addresses and lengths.
<i>chan_flag</i>	Specifies the device and bus specific flags for the transfer.

Description

The **d_map_slave** kernel service accepts a list of virtual buffer addresses and sizes and sets up the slave DMA controller for the requested DMA transfer. This includes setting up the system address generation hardware for a specific slave channel to indicate the specified data buffer(s), and enabling the specific hardware channel. The **d_map_slave** kernel service is not an exported kernel service, but a bus specific utility routine determined by the **d_map_init** kernel service and provided to the caller through the **d_handle** structure.

This service allows for scatter/gather capability of the slave DMA controller and also allows the device driver to coalesce multiple requests that are contiguous with respect to the device. The list is passed with the **dio** structure. If the **d_map_slave** kernel service is unable to complete the mapping due to resource, an error, **DMA_NORES** is returned, and the **bytes_done** field of the **dio** list is set to the number of bytes that were successfully mapped. This byte count is guaranteed to be a multiple of the *minxfer* parameter size of the device as provided to **d_map_slave**. Also, the *resid_iov* field is set to the index of the remaining **d_iovec** that could not be mapped. Unless the **DMA_BYPASS** flag is set, this service will verify access permissions to each page. If an access violation is encountered on a page within the list, an error, **DMA_NOACC** is returned and no mapping is done. The *bytes_done* field of the virtual list is set to the number of bytes preceding the faulting **d_iovec**. Also in this case, the *resid_iov* field is set to the index of the **d_iovec** entry that encountered the access violation.

The virtual addresses provided in the *vlist* parameter may be within multiple address spaces, distinguished by the cross memory structure pointed to for each element of the **dio** list. Each cross memory pointer can point to the same cross memory descriptor for multiple buffers in

the same address space, and for global space buffers, the pointers can be set to the address of the exported GLOBAL cross memory descriptor, *xmem_global*.

The *minxfer* parameter specifies the absolute minimum data transfer supported by the device (the device blocking factor). If the device supports a minimum transfer of 512 bytes (floppy and disks for example), the *minxfer* parameter would be set to 512. This allows the underlying services to map partial transfers to a correct multiple of the device block size.

Notes:

1. The **d_map_slave** kernel service does not support more than one outstanding DMA transfer per channel. Attempts to do multiple slave mappings on a single channel will corrupt the previous mapping(s).
2. You can use the **D_MAP_SLAVE** macro provided in the */usr/include/sys/dma.h* file to code calls to the **d_map_clear** kernel service.
3. The possible flag values for the *chan_flag* parameter can be found in */usr/include/sys/dma.h*. These flags can be logically ORed together to reflect the desired characteristics of the device and channel.

Return Values

DMA_NORES Indicates that resources were exhausted during the mapping

DMA_NOACC Indicates no access permission to a page in the list.

DMA_BAD_MODE

Indicates that the mode specified by the *chan_flag* parameter is not supported.

Implementation Specifics

The **d_map_clear** kernel service is part of the base device package of your platform.

Related Information

The **d_map_init** kernel service.

d_unmap_list Kernel Service

Purpose

Deallocates resources previously allocated on a **d_map_list** call.

Syntax

```
#include <sys/dma.h>
```

```
void d_unmap_list (*handle, *bus_list)
struct d_handle *handle
struct dio *bus_list
```

Parameters

handle Indicates the unique handle returned by the **d_map_init** kernel service.

bus_list Specifies a list of bus addresses and lengths.

Description

The **d_unmap_list** kernel service is a bus specific utility routine determined by the **d_map_init** kernel service that deallocates resources previously allocated on a **d_map_list** call.

The **d_unmap_list** kernel service must be called after I/O completion involving the area mapped by the prior **d_map_list** call. Some device drivers might choose to leave pages mapped for a long-term mapping of certain memory buffers. In this case, the driver must call **d_unmap_list** when it no longer needs the long-term mapping.

Note: You can use the **D_UNMAP_LIST** macro provided in the **/usr/include/sys/dma.h** file to code calls to the **d_unmap_list** kernel service. If not, you must ensure that the **d_unmap_list** function pointer is non-**NULL** before attempting the call. Not all platforms require the unmapping service.

Implementation Specifics

The **d_unmap_list** kernel service is part of the base device package of your platform.

Related Information

The **d_map_init** kernel service.

d_unmap_page Kernel Service

Purpose

Deallocates resources previously allocated on a **d_unmap_page** call.

Syntax

```
#include <sys/dma.h>

void d_unmap_page (*handle, *busaddr)
struct d_handle *handle
uint *busaddr
```

Parameters

handle Indicates the unique handle returned by the **d_map_init** kernel service.

busaddr Points to the *busaddr* field.

Description

The **d_unmap_page** kernel service is a bus specific utility routine determined by the **d_map_init** kernel service that deallocates resources previously allocated on a **d_map_page** call for a DMA master device.

The **d_unmap_page** service must be called after I/O completion involving the area mapped by the prior **d_map_page** call. Some device drivers might choose to leave pages mapped for a long-term mapping of certain memory buffers. In this case, the driver must call **d_unmap_page** when it no longer needs the long-term mapping.

Note: You can use the **D_UNMAP_PAGE** macro provided in the **/usr/include/sys/dma.h** file to code calls to the **d_unmap_page** kernel service. If not, you must ensure that the **d_unmap_page** function pointer is non-**NULL** before attempting the call. Not all platforms require the unmapping service.

Implementation Specifics

The **d_unmap_page** kernel service is part of the base device package of your platform.

Related Information

The **d_map_init** kernel service.

d_unmap_slave Kernel Service

Purpose

Deallocates resources previously allocated on a **d_map_slave** call.

Syntax

```
#include <sys/dma.h>

int d_unmap_slave (*handle)
struct d_handle *handle;
```

Parameters

handle Indicates the unique handle returned by the **d_map_init** kernel service.

Description

The **d_unmap_slave** kernel service deallocates resources previously allocated on a **d_map_slave** call, disables the physical DMA channel, and returns error and status information following the DMA transfer. The **d_unmap_slave** kernel service is not an exported kernel service, but a bus specific utility routine that is determined by the **d_map_init** kernel service and provided to the caller through the **d_handle** structure.

Note: You can use the **D_UNMAP_SLAVE** macro provided in the `/usr/include/sys/dma.h` file to code calls to the **d_unmap_slave** kernel service. If not, you must ensure that the **d_unmap_slave** function pointer is non-NULL before attempting to call. No all platforms require the unmapping service.

The device driver must call **d_unmap_slave** when the IO is complete involving a prior mapping by the **d_map_slave** kernel service.

Note: The **d_unmap_slave** kernel should be paired with a previous **d_map_slave** call. Multiple outstanding slave DMA transfers are not supported. This kernel service assumes that there is no DMA in progress on the affected channel and deallocates the current channel mapping.

Return Values

DMA_SUCC indicates successful transfer. The DMA controller did not report any errors and that the Terminal Count was reached.

DMA_TC_NOTREACHED Indicates a successful partial transfer. The DMA controller reported the Terminal Count reached for the intended transfer as setup by the **d_map_slave** call. Block devices consider this an error, however, for variable length devices this may not be an error.

DMA_FAIL Indicates that the transfer failed. The DMA controller reported an error. The device driver assumes the transfer was unsuccessful.

Implementation Specifics

The **d_unmap_slave** kernel service is part of the base device package of your platform.

Related Information

The **d_map_init** kernel service.

iomem_att Kernel Service

Purpose

Establishes access to memory mapped I/O.

Syntax

```
#include <sys/types.h>
#include <sys/adspc.h>

void *iomem_att(io_map_ptr)
struct io_map *io_map_ptr;

struct io_map {
    int key;
    int flags;
    int size;
    int BID;
    long long busaddress;
}
```

Parameters

The address of the **io_map** structure passes the following parameters to the **iomem_att** kernel service.

<i>key</i>	Set to <code>IO_MEM_MAP</code> .
<i>flags</i>	Describes the mapping.
<i>size</i>	Specifies the number of bytes to map.
<i>bid</i>	Specifies the bus identifier.
<i>busaddress</i>	Specifies the address of the bus.

Description

Note: The **iomem_att** kernel service is only supported on PowerPC based machines. All mappings are done with storage attributes: cache inhibited, guarded, and coherent. It is a violation of the PowerPC architecture to access memory with multiple storage modes. The caller of **iomem_att** must ensure no mappings using other storage attributes exist in the system.

Calling this function on a Power RS machine causes the system to crash.

The **iomem_att** kernel service provides temporary addressability to memory mapped IO. The **iomem_att** kernel service does the following:

- Allocates one segment of kernel address space
- Establishes kernel addressability
- Maps a contiguous region of memory mapped IO into that segment.

The addressability is valid only for the context that called **iomem_att**. The memory is addressable until **iomem_det** is called. IO memory must be mapped each time a context is entered, and freed before returning.

Note: Kernel address space is an exhaustible resource. When exhausted the system crashes. A driver must never map more than two IO regions at once, or call another

driver with an **iomem_att** outstanding. DMA, interrupt, and PIO kernel services may be called with up to two IO regions mapped.

The *size* parameter supports from 4096 bytes to 256 MB. The caller can specify a minimum of size bytes, but may choose to map up to 256 MB. The caller must not reference memory beyond size bytes.

Specifying **IOM_RDONLY** in the *flags* parameter results in a read only mapping. A store to memory mapped in this mode results in a data storage interrupt. If the *flag* parameter is **0** (zero) the memory is mapped read-write. All mappings are read-write on 601 based machines.

Execution Environment

The **iomem_att** kernel service can be called from either the process or interrupt environment.

Return Values

The **iomem_att** kernel service returns the effective address that can be used to address the IO memory.

Implementation Specifics

The **iomem_att** kernel service is part of the Base Operating System (BOS) Runtime.

Related Information

The **iomem_det** Kernel Service.

Kernel Extension and Device Driver Management Kernel Services in *AIX Version 4.1 Kernel Extensions and Device Support Programming Concepts*

iomem_det Kernel Service

Purpose

Releases access to memory mapped IO.

Syntax

```
#include <sys/types.h>
#include <sys/adspc.h>

void iomem_det (ioaddr)
void *ioaddr
```

Parameters

ioaddr Specifies the effective address returned by the **iomem_att** kernel service.

Description

The **iomem_det** kernel service releases memory mapped IO addressability. A call to the **iomem_det** kernel service must be made for every **iomem_att** call, with the address that **iomem_att** returned.

Execution Environment

The **iomem_det** kernel service can be called from either the process or interrupt environment.

Return Values

The **iomem_det** kernel service returns no return values.

Implementation Specifics

The **iomem_det** kernel service is part of the Base Operating System (BOS) Runtime.

Related Information

The **iomem_att** kernel service.

Kernel Extension and Device Driver Management Kernel Services in *AIX Version 4.1 Kernel Extensions and Device Support Programming Concepts*

ns_alloc Network Service

Purpose

Allocates use of a network device driver (NDD).

Syntax

```
#include <sys/ndd.h>

int ns_alloc (nndname, nndpp)
    char *nndname;
    struct ndd **nndpp;
```

Description

The **ns_alloc** network service searches the Network Service (NS) device chain to find the device driver with the specified *nndname* parameter. If the service finds a match, it increments the reference count for the specified device driver. If the reference count is incremented to 1, the **nnd_open** subroutine specified in the **ndd** structure is called to open the device driver.

Parameters

<i>nndname</i>	Specifies the device name to be allocated.
<i>nndpp</i>	Indicates the address of the pointer to a ndd structure.

Examples

The following example illustrates the **ns_alloc** network service:

```
struct ndd      *nndp;
error = ns_alloc("en0", &nndp);
```

Return Values

If a match is found and the **nnd_open** subroutine to the device is successful, a pointer to the **ndd** structure for the specified device is stored in the *nndpp* parameter. If no match is found or the open of the device is unsuccessful, a non-zero value is returned.

0	Indicates the operation was successful.
ENODEV	Indicates an invalid network device.
ENOENT	Indicates no network demuxer is available for this device.

The **nnd_open** routine may specify other return values.

Related Information

The **ns_free** network service.

ns_free Network Service

Purpose

Relinquishes access to a network device.

Syntax

```
#include <sys/ndd.h>
void ns_free (nddp)
    struct ndd *nddp;
```

Description

The **ns_free** network service relinquishes access to a network device. The **ns_free** network service also decrements the reference count for the specified **ndd** structure. If the reference count becomes 0, the **ns_free** network service calls the **ndd_close** subroutine specified in the **ndd** structure.

Parameters

nddp Specifies the **ndd** structure of the network device that is to be freed from use.

Examples

The following example illustrates the **ns_free** network service:

```
struct ndd *nddp
ns_free(nddp);
```

Files

net/cdli.c

Related Information

The **ns_alloc** network service.

rmalloc Kernel Service

Purpose

Allocates an area of memory.

Syntax

```
#include <sys/types.h>
```

```
caddr_t rmalloc (size, align)
int size
int align
```

Parameters

<i>size</i>	Specifies the number of bytes to allocate
<i>align</i>	Specifies alignment characteristics

Description

The **rmalloc** kernel service allocates an area of memory from the contiguous real memory heap. This area is the number of bytes in length specified by the *size* parameter and is aligned on the byte boundary specified by the *align* parameter. The *align* parameter is actually the log base 2 of the desired address boundary. For example, an *align* value of 4 requests that the allocated area be aligned on a 16 byte boundary.

The contiguous real memory heap, **real_heap**, is a heap of contiguous real memory pages located in the low 16MB of real memory. This heap is mapped virtually into the kernel extension segment. By nature, this heap is implicitly pinned, so no explicit pinning of allocated regions is necessary. This heap is provided primarily for device drivers whose devices can only address 0 to 16MB of real memory, so they must “bounce” the I/O in and out of a buffer **rmalloc**'ed from this heap. Also, this heap is useful for devices that require greater than 4K transfers, but do not support scatter/gather.

The **real_heap** is only supported on platforms with an ISA bus. On unsupported platforms, the **rmalloc** service returns **NULL** if the requested memory cannot be allocated.

The **rmfree** kernel service should be called to free allocation from a previous **rmalloc** call. The **rmalloc** kernel service can be called from the process environment only.

Return Values

Upon successful completion, the **rmalloc** kernel service returns the address of the allocated area. A **NULL** pointer is returned if the requested memory cannot be allocated.

Implementation Specifics

The **rmalloc** kernel service is part of the Base Operating System (BOS) Runtime.

Related Information

The **rmfree** kernel service.

rmfree Kernel Service

Purpose

Frees memory allocated by the **rmalloc** kernel service.

Syntax

```
#include <sys/types.h>

int rmfree (pointer, size)
caddr_t pointer
int size
```

Parameters

<i>pointer</i>	Specifies the address of the area in memory to free.
<i>size</i>	Specifies the size of the area in memory to free.

Description

The **rmfree** kernel service frees the area of memory pointed to by the *pointer* parameter in the contiguous real memory heap. This area of memory must be allocated with the **rmalloc** kernel service, and the *pointer* must be the pointer returned from the corresponding **rmalloc** kernel service call. Also, the *size* must be the same size that was used on the corresponding **rmalloc** call.

Any memory allocated in a prior **rmalloc** call must be explicitly freed with an **rmfree** call. This service can be called from the process environment only.

Return Values

0	Indicates successful completion
-1	Indicates one of the following:
	<ul style="list-style-type: none">• The area was not allocated by the rmalloc kernel service.• The heap was not initialized for memory allocation.

Implementation Specifics

The **rmfree** kernel service is part of the Base Operating System (BOS) Runtime.

Related Information

The **rmalloc** kernel service.

Index

A

- accessing device drivers, 1-3
- adapter device attributes, 6-7
- add_input_type kernel service, 12-25
- Address Resolution Protocol (ARP), 12-29, 13-4, data structures, 12-34
- aixgsc system call, 11-43
- allocate memory, rmalloc, A-20
- allocating memory, rmfree, A-21
- arpcom structure, 12-35
- arpreq structure, 12-38
- arptab structure, 12-37
- asynchronous routines, contrasted with synchronous, 1-1

B

- block address translation, 2-2
- block device driver, 1-6
- block I/O, 7-1
- bosboot command, 14-26
- buf structure, 7-3, 7-4, 8-2, 8-13, 8-16
- bufcall utility, 10-20
- bus I/O space, 2-17
- bus memory space (example), 2-18
- busresolve, 6-7

C

- call side, contrasted with interrupt side, 1-1
- canonical mode, 10-7
- cfg_dd structure, 8-14
- change methods, 6-5
- character device driver, 1-7
- character I/O, to block devices, 7-6
- chdisp command, 11-15
- commands
 - bosboot, 14-26
 - chdisp, 11-15
 - crash, 14-4, 14-5–14-25, 14-27, 14-47, 14-52
 - errinstall, 14-54
 - errlogger, 14-59
 - errmsg, 14-54
 - errpt, 14-52, 14-60
 - errupdate, 14-55, 14-58, 14-59
 - ifconfig, 12-39
 - lsdisp, 11-15
 - nm, 14-41
 - odmadd, 9-20
 - odmdelete, 9-20
 - setmaps, 10-1
 - slattach, 10-13
 - sliplogin, 10-13
 - sysdumpdev, 14-1
 - sysdumpstart, 14-2

- trace, 12-39, 14-63
- trcrpt, 14-63, 14-68
- trcstop, 12-39
- compiling
 - when using the kernel debugger, 14-40
 - when using trace, 14-83
- complex locks, 5-9
- configuration, 1-8
- configuration databases, ODM, 6-2
- configuration entry point, 1-10
- configuration method, for network interface driver, 12-39
- configuration methods, 6-1
 - network interface driver, 12-39
 - SCSI, 8-18
 - virtual file system, 9-18
- configuration routine, tty drivers, 10-14
- configure method, purpose of, 1-9
- configure methods, 6-4
- configuring devices with no parent, 6-6
- controller types, I/O, 2-4
- converting, addresses, 2-1
- converting file descriptor to device number, 1-5
- copyin kernel service, 4-5
- copyinstr kernel service, 4-5
- copyout kernel service, 4-5, 11-19
- crash command, 14-4, 14-5–14-25, 14-27, 14-47, 14-52
 - pcb subcommand, 14-19
 - ppd subcommand, 14-14
 - status subcommand, 14-18
 - thread subcommand, 14-20
- crash subcommands
 - buf, 14-6
 - buffer, 14-6
 - callout, 14-7
 - cm, 14-7
 - cpu, 14-7
 - dlock, 14-8
 - ds, 14-10
 - du, 14-10
 - dump, 14-10
 - fs, 14-11
 - inode, 14-11
 - kfp, 14-12, 14-21
 - knlist, 14-12, 14-47
 - le, 14-12, 14-47
 - mbuf, 14-13
 - mst, 14-13
 - ndb, 14-14
 - nm, 14-12, 14-14
 - od, 14-14, 14-27
 - print, 14-15
 - proc, 14-15

crash subcommands (*Continued*)

- quit, 14-17
- socket, 14-17
- stack, 14-17
- stat, 14-18
- trace, 14-12
- ts, 14-21
- tty, 14-21, 14-37
- user, 14-21
- var, 14-24
- vfs, 14-24
- vnode, 14-24
- xmalloc, 14-25

cross-memory services, 4-12

CuDep object class, 11-15

D

d_map_clear kernel service, A-2

d_map_disable kernel service, A-3

d_map_enable, A-4

d_map_init kernel service, A-5

d_map_list kernel service, A-6

d_map_page kernel service, A-8

d_map_slave, A-10

d_unmap_list kernel service, A-12

d_unmap_page kernel service, A-13

d_unmap_slave kernel service, A-14

data packet for Ethernet, 12-26

data structures, network interface driver and ARP,
12-34

deallocates resources

- d_map_clear, A-2

- d_unmap_list, A-12

- d_unmap_slave, A-14

debugger. See kernel debugger

debugging

- network interface driver, 12-39

- network interface drivers, 12-39

- virtual display driver, 11-63

define method, purpose of, 1-9

define methods, 6-3

del_input_type kernel service, 12-25

devdump kernel service, 8-2

device dependent structure, 8-14, 11-16

- example of, 11-33

device driver structure figure, 12-1

device driver types, block, 7-1

device head, 1-10

device methods

- how invoked, 1-9

- types to be provided, 1-9

Device Switch Table, 1-4

devinfo structure, 8-6

devstrat kernel service, 8-2, 8-4, 9-16

devswadd kernel service, 8-14, 11-16

devswdel kernel service, 11-17

devswqry kernel service, 11-15

direct memory access, for Micro Channel, 2-20

disable DMA, d_map_disable, A-3

display device driver, 11-14

- configuration, 11-14

DLPI interfaces, 13-3

DMA

- disable, d_map_disable, A-3

- enable, d_map_enable, A-4

DMA master, 2-20

DMA master devices

- deallocates resources, d_unmap_page, A-13

- mapping, d_map_page, A-8

DMA master transfer, sample call side routine to

- set up, 2-22

DMA operations, allocates and initializes

- resources, d_map_init, A-5

DMA slave, 2-20

dmp_add kernel service, 14-3

dmp_del kernel service, 14-3

dmx_8022_receive function, 12-17

dmx_status function, 12-17

dump

- See *als* system dump

- system, 14-1

E

e_sleep_thread kernel service, 9-15

enable DMA, d_map_enable, A-4

entry points

- configuration, 1-10

- ddclose, 7-3

- ddconfig, 7-2

- dddump, 7-6

- ddopen, 7-3

- ddstrategy, 7-3

- in STREAMS driver, 1-13

- interrupt, 1-21

- open and close, 1-11

- read, 1-11

- strategy, 1-12

- write, 1-12

entry points of a device driver, 1-10

errinstall command, 14-54

errlogger command, 14-59

errmsg command, 14-54

error logging, 14-52–14-60

- adding logging calls, 14-58

- coding steps, 14-54

error record template, 14-55

errpt command, 14-52, 14-60

errsave kernel service, 14-52, 14-58, 14-60

errupdate command, 14-55, 14-58, 14-59

EUC handling, 10-10

event notification, 5-6

F

figures

- Bus I/O Space, 2-17
- Bus Memory Space (Example), 2-18
- CDLI Device Driver Structure, 12-1
- Data Packet for Ethernet, 12-26
- Device Switch Table, 1-4
- Format of a Bus ID, 2-14
- From File Descriptor to Device Number, 1-5
- I/O Subsystem, 2-15
- IOCC Control Space, 2-16
- NID Data Structure Relationships, 12-34
- operating system kernel, 1-3
- SCSI Subsystem, 1-22
- Segmented Virtual Memory, 2-2
- STREAMS Driver Entry Points, 1-13
- System Device Hierarchy, 6-1

file system helper, 9-16

files

- /dev/error, 14-52, 14-60
- /dev/mem, 14-5
- /dev/scsi#, 8-2
- /dev/systrctl, 14-64, 14-65, 14-67
- /etc/filesystems, 9-16, 9-19
- /etc/trcfmt, 14-68, 14-84
- /etc/vfs, 9-16, 9-17, 9-20
- /usr/adm/ras/trcfile, 14-64
- /usr/lib/boot/unix, 14-5
- lp, 10-11
- net/if.h, 12-31, 12-38
- net/if_arp.h, 12-38
- str_tty.h, 10-24
- sys/buf.h, 7-4, 8-2, 8-16, 14-6
- sys/device.h, 11-17
- sys/devinfo.h, 8-6, 12-38
- sys/dir.h, 9-15
- sys/dump.h, 14-3, 14-5
- sys/erec.h, 14-57
- sys/err_rec.h, 14-59
- sys/errids.h, 14-58
- sys/file.h, 14-10
- sys/fshelp.h, 9-16
- sys/gfs.h, 9-5
- sys/mbuf.h, 12-38, 14-13
- sys/proc.h, 14-1, 14-16
- sys/scsi.h, 8-2, 8-3, 8-5, 8-7, 8-11
- sys/socket.h, 12-38, 14-17
- sys/statfs.h, 9-11
- sys/sysconfig.h, 8-14
- sys/timer.h, 14-7
- sys/trcctl.h, 14-67
- sys/trchkid.h, 14-71, 14-72, 14-84
- sys/trcmacros.h, 14-71
- sys/tty.h, 14-21
- sys/user.h, 14-1, 14-24
- sys/vattr.h, 9-13

- sys/vfs.h, 9-5, 14-24

- sys/vmount.h, 9-17

- sys/vnode.h, 9-6, 9-7, 14-25

- find_input_type kernel service, 12-25

- font_data structure, 11-32

- format of a bus ID, 2-14

- fp_close kernel service, 8-2, 8-15

- fp_ioctl kernel service, 8-2, 11-15

- fp_open kernel service, 8-2, 8-15

- fp_opendev kernel service, 11-15

G

- GAI object class, 11-15

- gfs structure, 9-2, 9-5

- gfsadd kernel service, 9-3, 9-5, 9-9

- gfsdel kernel service, 9-5, 9-10

- gnode structure, 9-7

H

- hardware interrupts, 3-2

I

- I/O, 2-1

- I/O address segments, for Micro Channel, 2-15

- I/O controller types, 2-4

- I/O spaces, for Micro Channel, 2-15

- I/O subsystem, 2-15

- if_attach kernel service, 12-23, 12-25

- if_detach kernel service, 12-25

- if_down kernel service, 12-25

- if_nostat kernel service, 12-25

- ifa_ifwithaddr kernel service, 12-25

- ifa_ifwithstaddr kernel service, 12-25

- ifa_ifwithnet kernel service, 12-25

- ifaddr structure, 12-37

- ifconfig command, 12-39

- ifnet structure, 12-23, 12-25, 12-31, 12-35

- ifreq structure, 12-37

- ifunit kernel service, 12-25

- init_heap kernel service, 4-3

- initialization, network device driver, 12-2

- interface protocols, TLI and XTI, 13-5

- interrupt level mapping, 3-6

- interrupt levels, 3-2

- interrupt service times, 3-5

- interrupt side, contrasted with call side, 1-1

- interrupt side routines, 1-21

- interrupts, 3-1

- necessity for device driver to handle, 1-1

- IOCC Control Space, 2-16

- IOCTL calls, supported by network interface driver, 12-31

- iodone kernel service, 7-4, 8-1, 8-4, 8-5, 8-16, 9-16

- iomem_att kernel service, A-15

- iomem_det kernel service, A-17

- ISA bus configuration, 3-3, 6-8

K

- kernel debugger
 - accessing global data, 14-48
 - compiler listings, 14-41
 - compiling options, 14-40
 - displaying registers, 14-50
 - entering, 14-26
 - map files, 14-42
 - setting breakpoints, 14-40, 14-45–14-48
 - stack trace, 14-51
 - subcommands, 14-27–14-40
 - breakpoints, 14-40
 - dereferencing a pointer, 14-40
 - expressions, 14-40
 - reserved variables, 14-38
 - variables, 14-38
- kernel figure, 1-3
- kernel services
 - add_input_type, 12-25
 - copyout, 11-19
 - del_input_type, 12-25
 - devdump, 8-2
 - devstrat, 8-2, 8-4, 9-16
 - devswadd, 8-14, 11-16
 - devswdel, 11-17
 - devswqry, 11-15
 - disable_lock, 10-19
 - dmp_add, 14-3
 - dmp_del, 14-3
 - e_sleep_thread, 9-15
 - errsave, 14-52, 14-58, 14-60
 - find_input_type, 12-25
 - fp_close, 8-2, 8-15
 - fp_ioctl, 8-2, 11-15
 - fp_open, 8-2, 8-15
 - fp_opendev, 11-15
 - gfsadd, 9-3, 9-5, 9-9
 - gfsdel, 9-5, 9-10
 - if_attach, 12-23, 12-25
 - if_detach, 12-25
 - if_down, 12-25
 - if_nostat, 12-25
 - ifa_ifwithaddr, 12-25
 - ifa_ifwithstaddr, 12-25
 - ifa_ifwithnet, 12-25
 - ifunit, 12-25
 - iodone, 7-4, 8-1, 8-4, 8-5, 8-16, 9-16
 - lookupvp, 9-10
 - net_error, 12-39
 - pin, 8-13
 - pincode, 8-13, 8-15
 - pinu, 8-13
 - uiomove, 9-15, 11-16, 11-17
 - unlock_enable, 10-19
 - vfsrele, 9-11, 9-12
 - vm_mount, 9-9
 - vms_create, 9-15
 - vn_free, 9-6
 - vn_get, 9-6, 9-7

L

- lc_sjis module, 10-1
- lock overview, 5-9
- ldterm line discipline, 10-6
- levels, interrupt, 3-2
- library routines, restricted device driver use, 1-2
- loading convention, STREAMS, 13-4
- lock models, 5-9
- locking, simple locks, 8-15
- locking options, for STREAMS modules and drivers, 13-4
- lockl locks, 5-9
- lookupvp kernel service, 9-10
- lsdisp command, 11-15
- ltpin kernel service, 4-3
- ltunpin kernel service, 4-4

M

- main routine, device driver does not have, 1-2
- major number, 1-3
- maps DMA master devices, d_map_page, A-8
- mbuf structure, 12-26
- memory, rmfree, A-21
- memory access services, 4-5
- memory allocation, rmalloc, A-20
- memory allocation services, 4-1
- memory management, 4-1
- memory mapped I/O
 - iomem_att, A-15
 - iomem_det, A-17
- memory pinning services, 4-3
- Micro Channel
 - Direct Memory Access for, 2-20
 - I/O address segments for, 2-15
 - kinds of I/O spaces for, 2-15
 - programmed I/O for, 2-19
- Micro Channel adapters, displaying registers, 14-50
- minor number, 1-4
- mount helper, 9-17
- MP-efficient code, 5-13
- MP-safe interrupt handling, 3-10
- multiprocessing serialization, 5-8
- multiprocessor environment, tty, 10-19
- multiprocessor interrupt concerns, 3-10

N

- nd_add_filter function, 12-14
- nd_add_status function, 12-15
- nd_del_filter function, 12-15
- nd_del_status function, 12-16
- nd_receive function, 12-12, 12-16
- nd_response function, 12-17
- nd_status function, 12-12, 12-17
- ndd_close entry point, 12-6
- ndd_ctl entry point, 12-10
- ndd_open entry point, 12-5
- ndd_output entry point, 12-7

- net_error kernel service, 12-39
- network, interface from protocol, 13-12
- network address translation to hardware address, 12-29
- network device driver, 12-2
 - changes, 12-2
 - initialization and termination, 12-2
- network interface driver
 - basic functions, 12-22
 - changes, 12-22
 - communicating with the device handler, 12-29
 - communicating with the IP, 12-25
 - configuration method for, 12-39
 - configuration methods, 12-39
 - data structures, 12-34
 - debugging, 12-39
 - initializing, 12-22
 - ioctl calls, 12-31
 - loading, 12-22
 - outgoing packets, 12-26
 - output data, 12-29
 - purpose, 12-22
 - terminating, 12-34
 - tracing, 12-39
 - translating network addresses to hardware addresses, 12-29
- network interface driver data structures, 12-34
 - arpcom structure, 12-35
 - arpreq structure, 12-38
 - arptab structure, 12-37
 - ifaddr structure, 12-37
 - ifnet structure, 12-23, 12-25, 12-31, 12-35
 - ifreq structure, 12-37
 - mbuf structure, 12-26
 - sockaddr structure, 12-29
 - xx_softc structure, 12-35
- network interface drivers
 - debugging, 12-39
 - tracing, 12-39
- network interfaces, 13-1
- network protocols, 13-1
 - socket, writing or porting, 13-7
 - streams, writing or porting, 13-3
- network to protocol interface, 13-14
- NID Data Structure Relationships, 12-34
- nls module, 10-1
- nm command, 14-41
- ns_add_filter, 12-18
 - sample DLPI call to, 12-21
- ns_add_status, 12-19
- ns_alloc, sample call to, 12-21
- ns_alloc network service, A-18
- ns_attech kernel service, 12-11
- ns_del_filter, 12-19
- ns_del_status, 12-19

- ns_detach kernel service, 12-11
- ns_free network service, A-19

O

- object classes
 - CuDep, 11-15
 - GAI, 11-15
 - PdAt, 11-14
 - PdDv, 11-14
 - purpose of each, 6-2
- Object Data Manager, 1-9
- object files, pinning, 1-22
- ODM configuration databases, 6-2
- odmadd command, 9-20
- odmdelete command, 9-20
- open and close entry points, 1-11
- outgoing packets, network interface driver, 12-26
- output data, network interface driver, 12-29

P

- page address translation, 2-3
- paging, compared with pinning, 1-2
- parent, configuring devices with no, 6-6
- PCI bus configuration, 6-8
- PdAt object class, 11-14
- PdDv object class, 11-14
- performance tracing. *See* tracing
- phys_displays structure, 11-33
- Physical Volume Identifier. *See* PVID
- pin kernel service, 4-4, 8-13
- pincode kernel service, 4-4, 8-13, 8-15
- pinning, device driver object files, 1-22
- pinning code and data, 1-2
- pinu kernel service, 4-4, 8-13
- preempting, device drivers subject to, 1-2
- presentation space, 11-26
 - height, 11-27
 - width, 11-27
- priorities, interrupt, 3-4
- programmed I/O, for Micro Channel, 2-19
- protocol, interface from network, 13-14
- protocol address resolution, 13-4
- protocol interfaces via DLPI, 13-2
- protocol to network interface, 13-12
- protocol to socket interface, 13-11
- PVID (Physical Volume Identifier), 8-17

R

- raw I/O. *See* character I/O
- read entry point, 1-11
- real time, device drivers required to execute in, 1-2
- real-time timers, 5-4
- rmalloc kernel service, A-20
- rmfree kernel service, A-21

S

- sample code, trace format file, 14-77
- sample device driver, 1-14
- samples
 - call-side routine to set up DMA master transfers, 2-22
 - cross-memory services, 4-12
 - DLPI call to `ns_add_filter`, 12-21
 - ifnet structure, 12-36
 - input device load module, 11-61
 - ioctl routine of network interface driver, 12-31
 - loading and initializing network interface driver, 12-23
 - mapping multicast address in NID, 12-34
 - MP safe code, 5-10
 - network device driver configuration, 12-3
 - `ns_add_filter` fragment, 12-18
 - output routine for network interface driver, 12-27
 - socket protocol, 13-15
 - socket receive buffer, adding data to, 13-12
 - virtual memory management services, 4-11
 - `xmemdma` kernel service, 4-14
- `sc_buf` structure, 8-2–8-4
- `sc_inquiry` structure, 8-7, 8-8
- SCSI adapter device driver, 8-1
- SCSI adapter device driver routines
 - See `als` SCSI device driver routines
 - close, 8-5
 - config, 8-5
 - ioctl, 8-5
 - open, 8-5
 - openx, 8-5
 - strategy, 8-5
- SCSI adapter ioctl operations
 - IOCINFO, 8-6
 - SCIODIAG, 8-11
 - SCIODNLD, 8-11
 - SCIOHALT, 8-10
 - SCIOINQU, 8-7
 - SCIORESET, 8-10
 - SCIOSTART, 8-6
 - SCIOSTOP, 8-7
 - SCIOSTUNIT, 8-8
 - SCIOTRAM, 8-11
 - SCIOTUR, 8-9
- SCSI configuration methods, 8-18
- SCSI device attributes, 8-18
- SCSI device driver, 8-1
- SCSI device driver routines
 - See `als` SCSI adapter device driver routines
 - close, 8-15
 - config, 8-14
 - dump, 8-16
 - ioctl, 8-14
 - open, 8-15
 - read, 8-15
 - strategy, 8-16
 - write, 8-16
- SCSI device driver structure
 - bottom half routines, 8-13
 - top half routines, 8-13
- SCSI Subsystem, 1-22
- SCSI subsystem components
 - SCSI adapter device driver, 8-1
 - SCSI device driver, 8-1
- segment register contents, 2-3
- Segmented Virtual Memory, 2-2
- serialization, 5-1
 - for Streams modules and drivers, 13-4
- serialization services, 5-8
- service times, interrupt guidelines, 3-5
- setmaps command, 10-1
- simple locks, 5-9
- slattach command, 10-13
- slip line discipline, 10-13
- sliplogin command, 10-13
- sockaddr structure, 12-29
- socket network protocols, writing or porting, 13-7
- socket protocol
 - interfaces to support, 13-8
 - sample, 13-15
- socket protocols
 - initializing, 13-7
 - loading, 13-8
- socket receive buffer, sample adding data to, 13-12
- source addresses, DLPI interpretation, 13-4
- spr line discipline, 10-11
- states, of devices, 6-1
- `str_install` utility, 10-14, 10-20
- strategy entry point, 1-12
- stream head, 10-3
- streams, serialization and locking options, 13-4
- STREAMS device driver, 1-7
- STREAMS Driver Entry Points, 1-13
- STREAMS driver routines
 - write-side put, 1-13
 - write-side service or read-side service, 1-13
- STREAMS entry points, 1-13
- STREAMS loading convention, 13-4
- streams network protocols, writing or porting, 13-3
- streams user interfaces, 13-1
- STREAMS-based tty, 10-1
- structures
 - `arpcom`, 12-35
 - `arpreq`, 12-38
 - `arptab`, 12-37
 - `buf`, 7-3, 7-4, 8-2, 8-13, 8-16
 - `cfg_dd`, 8-14
 - `devinfo`, 8-6
 - `font_data`, 11-32
 - `gfs`, 9-2, 9-5
 - `gnode`, 9-7
 - `ifaddr`, 12-37
 - `ifnet`, 12-23, 12-25, 12-31, 12-35
 - `ifreq`, 12-37
 - `mbuf`, 12-26
 - `phys_displays`, 11-33
 - `sc_buf`, 8-2–8-4

- structures (*Continued*)
 - sc_card_diag, 8-11
 - sc_inquiry, 8-7, 8-8
 - sockaddr, 12-29
 - tioc_reply, 10-5
 - uio, 8-14, 8-15
 - vfs, 9-2, 9-5
 - vfsops, 9-2, 9-3, 9-5
 - vmount, 9-3, 9-18
 - vnode, 9-6
 - vnodeops, 9-2, 9-3, 9-5
 - vtt_box_rc_parms, 11-31
 - vtt_cp_parms, 11-31
 - vtt_rc_parms, 11-31
 - xx_softc, 12-35
- subroutines, sysconfig, 8-14, 9-3, 9-18, 9-19, 11-16, 12-39, 14-46
- synchronization, 5-1
- synchronous routines, contrasted with asynchronous, 1-1
- sysconfig subroutine, 8-14, 9-3, 9-18, 9-19, 11-16, 12-39, 14-46
- sysdumpdev command, 14-1
- sysdumpstart command, 14-2
- system dump, 14-1
 - formatting, 14-4
 - including device driver information, 14-2
 - initiating, 14-1
- system dump components
 - component dump table, 14-2
 - master dump table, 14-2

T

- termination, network device driver, 12-2
- timeout utility, 10-20
- timer services, 5-2
- timers
 - real-time, 5-4
 - watchdog, 5-2
- tioc module, 10-4
- tioc_reply structure, 10-5
- TLI addresses, 13-4
- TLI interface protocol, 13-5
- trace command, 12-39, 14-63
- trace events
 - defining, 14-70–14-84
 - event IDs, 14-72–14-84
 - determining location of, 14-72–14-84
 - format file example, 14-77–14-84
 - format file stanzas, 14-73–14-84
 - forms, 14-71–14-84
 - macros, 14-71–14-84
- trace report
 - filtering, 14-86
 - producing, 14-68–14-70
 - reading, 14-85–14-86
- tracing, 14-61–14-86
 - configuring, 14-63

- controlling, 14-65–14-68
 - for network interface drivers, 12-39
- starting, 14-62, 14-63
- translating, addresses, 2-1
- transparent ioctl, 10-4
- trcrpt command, 14-63, 14-68
- trcstop command, 12-39
- tty drivers, 10-13
- tty ioctls, 10-21
- tty line disciplines
 - ldterm, 10-6
 - slip, 10-13
 - sptr, 10-11
- tty multiprocessor environment, 10-19
- tty open disciplines, 10-15
- tty pacing disciplines, 10-15
- tty stream head, 10-3
- tty subsystem, 10-1
- types of device drivers, 1-5

U

- uc_sjis module, 10-1
- uio structure, 8-14, 8-15
- uimove kernel service, 4-5, 9-15, 11-16, 11-17
- unconfigure methods, 6-5
- undefine methods, 6-6
- uniprocessor serialization, 5-8
- unpin kernel service, 4-4
- unpincode kernel service, 4-4
- unpinu kernel service, 4-5
- user software compared with device driver, 1-1

V

- vfs structure, 9-2, 9-5
- vfsops structure, 9-2, 9-3, 9-5
- vfsrele kernel service, 9-11, 9-12
- virtual display driver, debugging, 11-63
- virtual display driver routines
 - device specific routines, 11-20
 - activate, 11-20
 - clear rectangle, 11-22
 - copy full lines, 11-21
 - copy line segment, 11-23
 - deactivate, 11-24
 - define cursor, 11-24
 - draw text, 11-29
 - initialize, 11-25
 - move cursor, 11-27
 - scroll, 11-27
 - terminate, 11-28
- display device driver routines, 11-16
 - close, 11-18
 - configure, 11-16
 - interrupt handling, 11-20
 - ioctl, 11-19
 - open, 11-17

- virtual file system
 - components, 9-8
 - configuration, 9-18
 - creating kernel extensions, 9-9
 - definition of terms, 9-21
 - file system helper, 9-16
 - installing, 9-19
 - loading, 9-20
 - mount helper, 9-17
- virtual file system data structures, 9-3
 - gfs structure, 9-2, 9-5
 - gnode structure, 9-7
 - relationship between, 9-3
 - vfs structure, 9-2, 9-5
 - vmount, 9-3, 9-18
 - vnode structure, 9-6
- virtual file system entry points
 - config, 9-9
 - init, 9-10
 - rootinit, 9-10
- virtual file system operations
 - See *als*/vnode operations
 - vfs_cntl, 9-11
 - vfs_mount, 9-10
 - vfs_root, 9-11
 - vfs_statfs, 9-11
 - vfs_sync, 9-11
 - vfs_unmount, 9-11
 - vfs_vget, 9-11
- virtual memory management services, 4-6
- virtual memory operations, 9-15
- vm_cflush kernel service, 4-9
- vm_det kernel service, 4-9
- vm_mount kernel service, 4-9, 9-9
- vm_move kernel service, 4-9
- vm_release kernel service, 4-10
- vm_releasep kernel service, 4-11
- vm_umount kernel service, 4-9
- vm_write kernel service, 4-10
- vm_writep kernel service, 4-10
- vmount structure, 9-3, 9-18
- vms_att kernel service, 4-8
- vms_create kernel service, 4-8, 9-15
- vms_delete kernel service, 4-8
- vms_handle kernel service, 4-8
- vms_iowait kernel service, 4-10
- vn_free kernel service, 9-6
- vn_get kernel service, 9-6, 9-7
- vnode operations
 - vn_close, 9-13
 - vn_getattr, 9-12
 - vn_hold, 9-12
 - vn_lookup, 9-14
 - vn_open, 9-13
 - vn_rdw, 9-14
 - vn_readdir, 9-14
 - vn_rele, 9-12
 - vn_strategy, 9-14
- vnode structure, 9-6
- vnodeops structure, 9-2, 9-3, 9-5
- vtt_box_rc_parms structure, 11-31
- vtt_cp_parms structure, 11-31
- vtt_rc_parms structure, 11-31

W

- watchdog timers, 5-2
- write entry point, 1-12
- write-side or read-side service routine in STREAMS driver, 1-13
- write-side put routine in STREAMS driver, 1-13

X

- xmalloc kernel service, 4-1
- xmattach kernel service, 4-13
- xmdetach kernel service, 4-13
- xmemdma kernel service, 4-13
- xmemin kernel service, 4-13
- xmemout kernel service, 4-13
- xmfree kernel service, 4-3
- XTI addresses, 13-4
- XTI interface protocol, 13-5

Reader's Comment Form

AIX Version 4.1 Writing a Device Driver

SC23-2593-00

Please use this form only to identify publication errors or to request changes in publications. Your comments assist us in improving our publications. Direct any requests for additional publications, technical questions about systems, changes in programming support, and so on, to your sales representative or to your dealer. You may use this form to communicate your comments about this publication, its organization, or subject matter, with the understanding that we may use or distribute whatever information you supply in any way we believe appropriate without incurring any obligation to you.

Note: To send comments electronically, use this commercial Internet address: aix6kpub@austin.ibm.com.

- If your comment does not need a reply (for example, pointing out a typing error), check this box and do not include your name and address below. If your comment is applicable, we will include it in the next revision of the manual.

- If you would like a reply, check this box. Be sure to print your name and address below.

Page	Comments

Please contact your sales representative or dealer to request additional publications.

Please print

Date _____

Your Name _____

Company Name _____

Mailing Address _____

Phone No. () _____
Area Code

Tape

Please Do Not Staple

Tape

Fold

Fold

Cut or Fold Along Line

Publications Department
Department 997, Internal Zip 9630
11400 Burnet Rd.
Austin, Texas 78758-3493

Fold

Fold

Tape

Please Do Not Staple

Tape

