

# IRIX® 6.3 for O2™ Device Driver Programming Guide

Document Number 007-3443-002

## CONTRIBUTORS

Written by David Cortesi

Illustrated by Dany Galgani

Edited by Christina Cary

Significant engineering contributions by (in alphabetical order): Rich Altmaier, Peter Baran, Brad Eacker, Ben Fathi, Steve Haehnichen, Bruce Johnson, Tom Lawrence, Greg Limes, Ben Mahjoor, Charles Marker, Dave Olson, Bhanu Prakash, James Putnam, Sarah Rosedahl, Brett Rudley, Deepinder Setia, Adam Sweeney, Michael Wang, Daniel Yau.

Beta test contributions by: Jeff Stromberg of GeneSys

St Peter's Basilica image courtesy of ENEL SpA and InfoByte SpA. Disk Thrower image courtesy of Xavier Berenguer, Animatica.

© 1998, Silicon Graphics, Inc.— All Rights Reserved

The contents of this document may not be copied or duplicated in any form, in whole or in part, without the prior written permission of Silicon Graphics, Inc.

## RESTRICTED RIGHTS LEGEND

Use, duplication, or disclosure of the technical data contained in this document by the Government is subject to restrictions as set forth in subdivision (c) (1) (ii) of the Rights in Technical Data and Computer Software clause at DFARS 52.227-7013 and/or in similar or successor clauses in the FAR, or in the DOD or NASA FAR Supplement. Unpublished rights reserved under the Copyright Laws of the United States. Contractor/manufacturer is Silicon Graphics, Inc., 2011 N. Shoreline Blvd., Mountain View, CA 94043-1389.

Silicon Graphics and the Silicon Graphics logo, and the terms CHALLENGE, Crimson, Indigo, Indigo<sup>2</sup>, Indigo<sup>2</sup> Maximum Impact, Indy, IRIX, O2, Onyx, Origin2000, POWER CHALLENGE, POWER Channel, POWER Indigo<sup>2</sup>, and POWER Onyx are trademarks of Silicon Graphics, Inc. MIPS, R4000, R8000, and R10000 are trademarks of MIPS Technologies, Inc. Sun and SunOS are trademarks of Sun Microsystems, Inc. MC6800, MC68000, and VERSAbus are trademarks of Motorola Corporation. IBM is a trademark of International Business Machines. Intel is a trademark of Intel Corporation. UNIX is a registered trademark in the United States and other countries, licensed exclusively through X/Open Company, Ltd. X Window System is a trademark of Massachusetts Institute of Technology.

IRIX® for O2™ Device Driver Programming Guide  
Document Number 007-3443-002

---

# Contents

<b>List of Examples</b>	xvii
<b>List of Figures</b>	xix
<b>List of Tables</b>	xxi
<b>About This Guide</b>	xxv
What You Need to Know	xxv
What This Guide Contains	xxvi
Other Sources of Information	xxvii
Developer Program	xxvii
Internet Resources	xxvii
Standards Documents	xxviii
Important Reference Pages	xxviii
Additional Reading	xxix
Conventions Used in This Guide	xxx

<b>PART I</b>	<b>IRIX Device Integration</b>
<b>1.</b>	<b>Physical and Virtual Memory</b> 3
	Physical Address Space 4
	Physical Device Addresses 4
	Physical Memory Addresses 4

:

---

CPU Access to Memory and Devices	5
CPU Modules	5
CPU Access to Memory	6
Processor Operating Modes	8
Virtual Address Mapping	8
Address Space Creation	9
Address Exceptions	9
CPU Access to Device Registers	10
Direct Memory Access	11
PIO Addresses and DMA Addresses	12
Cache Use and Cache Coherency	15
The 32-Bit Address Space	16
Segments of the 32-bit Address Space	16
Virtual Address Mapping	18
User Process Space—kuseg	18
Kernel Virtual Space—kseg2	19
Cached Physical Memory—kseg0	19
Uncached Physical Memory—kseg1	20
The 64-Bit Address Space	20
Segments of the 64-Bit Address Space	20
Compatibility of 32-Bit and 64-Bit Spaces	22
Virtual Address Mapping	22
User Process Space—xkuseg	23
Supervisor Mode Space—xksseg	23
Kernel Virtual Space—xkseg	24
Cache-Controlled Physical Memory—xkphys	24
Device Driver Use of Memory	26
Allowing for 64-Bit Mode	26
Memory Use in User-Level Drivers	27
Memory Use in Kernel-Level Drivers	28

- 
- 2. **Device Configuration** 31
    - Hardware Inventory 31
      - Using the Hardware Inventory 32
      - Creating an Inventory Entry 34
    - Device Special Files 34
      - Device Representation 35
      - Defining Device Names 37
    - Configuration Files 39
      - Master Configuration Database 40
      - System Configuration Files 41
      - System Tuning Parameters 41
      - X Display Manager Configuration 41
  - 3. **Device Control Software** 43
    - User-Level Device Control 43
      - EISA Mapping Support 44
      - VME Mapping Support 44
      - PCI Mapping Support 45
      - User-Level DMA From the VME Bus 45
      - User-Level Control of SCSI Devices 45
      - Managing External Interrupts 46
      - User-Level Interrupt Management 46
      - Memory-Mapped Access to Serial Ports 47
    - Kernel-Level Device Control 47
      - Kinds of Kernel-Level Drivers 48
      - Typical Driver Operations 48
      - Upper and Lower Halves 56
      - Layered Drivers 58
      - Combined Block and Character Drivers 58
      - Drivers for Multiprocessors 59
      - Loadable Drivers 59

:

---

## **PART II     Device Control From Process Space**

- 4.     User-Level Access to Devices    63**
  - VME Programmed I/O    64
    - Mapping a VME Device Into Process Address Space    64
    - VME PIO Access    67
    - VME PIO Bandwidth    68
  - EISA Programmed I/O    68
    - Mapping an EISA Device Into Memory    69
    - EISA PIO Bandwidth    71
  - VME User-Level DMA    72
    - Using the udmalib Functions    73
    - Advantages of User DMA    75
    - DMA Engine Bandwidth    75
    - Example User DMA Function    76
  - PCI Programmed I/O    78
    - Mapping a PCI Device Into Process Address Space    78
  
- 5.     User-Level Access to SCSI Devices    81**
  - Overview of the dsreq Driver    82
  - Generic SCSI Device Special Files    82
    - Major and Minor Device Numbers in /dev/scsi    83
    - Form of Filenames in /dev/scsi    83
    - Creating Additional Names in /dev/scsi    84
    - Relationship to Other Device Special Files    85
  - The dsreq Structure    85
    - Values for ds\_flags    87
    - Data Transfer Options    89
    - Return Codes and Status Values    89
  - Testing the Driver Configuration    92
  - Using the Special DS\_RESET and DS\_ABORT Calls    93
    - Using DS\_ABORT    93
    - Using DS\_RESET    94

---

	Using dslib Functions	94
	dslib Functions	94
	Using dsopen() and dsclose()	95
	Issuing a Request With doscsireq()	97
	SCSI Utility Functions	97
	Using Command-Building Functions	100
	Example dslib Program	107
<b>6.</b>	<b>Control of External Interrupts</b>	<b>117</b>
	External Interrupts in Challenge and Onyx Systems	118
	Generating Outgoing Signals	118
	Receiving Incoming External Interrupts	119
<b>7.</b>	<b>User-Level Interrupts</b>	<b>125</b>
	Overview of ULI	125
	The User Level Interrupt Handler	126
	Restrictions on the ULI Handler	126
	Planning for Concurrency	127
	Using Multiple Devices	128
	Setting Up	128
	Opening the Device Special File	128
	Locking the Program Address Space	129
	Registering the Interrupt Handler	130
	Interacting With the Handler	131
	Sample Program	133

### **PART III      Kernel-Level Drivers**

<b>8.</b>	<b>Structure of a Kernel-Level Driver</b>	<b>139</b>
	Summary of Driver Structure	140
	Entry Point Naming and lboot	140
	Entry Point Summary	142

:

---

Driver Flag Constant	145
Flag D_MP	145
Flag D_WBACK	146
Flag D_MT	146
Flag D_OLD	146
Initialization Entry Points	147
When Initialization Is Performed	147
Entry Point init()	148
Entry Point edtinit()	148
Entry Point start()	149
Open and Close Entry Points	150
Entry Point open()	150
Entry Point close()	153
Control Entry Point	154
Choosing the Command Numbers	155
Supporting 32-Bit and 64-Bit Callers	155
User Return Value	155
Data Transfer Entry Points	155
Entry Points read() and write()	155
Entry Point strategy()	157
Poll Entry Point	158
Use and Operation of poll(2)	159
Entry Point poll()	160
Memory Map Entry Points	161
Concepts and Use of mmap()	162
Entry Point map()	163
Entry Point mmap()	165
Entry Point unmap()	166
Interrupt Entry Point	167
Associating Interrupt to Driver	167
Interrupt Handler Operation	167

- 
- Support Entry Points 170
    - Entry Point unload() 170
    - Entry Point halt() 171
    - Entry Point size() 172
    - Entry Point print() 172
  - Handling 32-Bit and 64-Bit Execution Models 173
  - Planning for Multiprocessor Use 174
    - The Multiprocessor Environment 174
    - Synchronizing Within Upper-Half Functions 176
    - Coordinating Upper-Half and Interrupt Entry Points 177
    - Converting a Uniprocessor Driver 178
    - Example Conversion Problem 179
  - 9. **Device Driver/Kernel Interface** 181
    - Important Data Types 182
      - The Device Number Types 182
      - Structure `uio_t` 184
      - Structure `buf_t` 185
      - Lock and Semaphore Types 187
    - Important Header Files 188
    - Memory Allocation 189
      - General-Purpose Allocation 190
      - Allocating Objects of Specific Kinds 191
      - Suballocation Functions 193
    - Transferring Data 194
      - General Data Transfer 195
      - Transferring Data Through a `uio_t` Object 197
    - Managing Virtual and Physical Addresses 198
      - Testing Device Physical Addresses 198
      - Managing Mapped Memory 199
      - Working With Page and Sector Units 200
      - Setting Up a DMA Transfer 201
    - User Process Administration 205
      - Sending a Process Signal 206

:

---

Waiting and Mutual Exclusion	206
Mutual Exclusion Compared to Waiting	207
Basic Locks	208
Long-Term Locks	210
Reader/Writer Locks	213
Priority Level Functions	215
Waiting for Time to Pass	216
Waiting for Memory to Become Available	218
Waiting for Block I/O to Complete	219
Waiting for a General Event	221
Semaphores	224
<b>10. Building and Installing a Driver</b>	<b>227</b>
Defining Device Numbers	228
Selecting a Major Number	228
Selecting Minor Numbers	229
Defining Device Special Files	229
Static Definition of Device Special Files	229
Dynamic Definition of Device Special Files	229
Compiling and Linking	230
Using /var/sysgen/Makefile.kernio	230
Compiler Variables	231
Compile Options, 32-Bit Kernel	232
Compile Options, 64-Bit Kernel	233
Configuring a Nonloadable Driver	234
How Names Are Used in Configuration	234
Placing the Object File in /var/sysgen/boot	235
Describing the Driver in /var/sysgen/master.d	235
Configuring a Kernel	238
Generating a Kernel	239

- 
- Configuring a Loadable Driver 239
    - Public Global Variables 239
    - Compile Options for Loadable Drivers 240
    - Master File for Loadable Drivers 240
    - Registration 241
    - Loading 241
    - Unloading 242
  - Configuring for a Dynamic Major Number 243
  - 11. Testing and Debugging a Driver 245**
    - Preparing the System for Debugging 245
      - Placing symmon in the Volume Header 246
      - Enabling Debugging in irix.sm 247
      - Generating a Debugging Kernel 249
      - Specifying a Separate System Console 249
      - Verifying the Debugging Tools 250
    - Producing Diagnostic Displays 251
      - Using cmn\_err 251
      - Using printf() 253
      - Using ASSERT 254
    - Using symmon 254
      - How symmon Is Entered 255
      - Commands of symmon 257
      - Syntax of Command Elements 257
      - Commands for Symbol Conversion and Lookup 258
      - Commands to Control Execution Flow 259
      - Commands to Manage Virtual Memory 261
      - Commands to Display Memory 262
      - Utility Commands 263

:

---

Using idbg	264
Loading and Invoking idbg	264
Commands of idbg	266
Commands to Display Memory and Symbols	267
Commands to Display Process Information	267
Commands to Display Locks and Semaphores	269
Commands to Display I/O Status	269
Commands to Display buf_t Objects	270
Commands to Display STREAMS Structures	270
Commands to Display Network-Related Structures	271
Using icrash	271
<b>12. Driver Example</b>	<b>273</b>
Installing the Example Driver	273
Obtaining the Source Files	274
Compiling the Example Driver	274
Configuring the Example Driver	275
Creating Device Special Files	277
Verifying Driver Operation	277
Example Driver Source Files	280
Descriptive File	280
System File	281
Header File	281
Source File	283

## **PART IV      SCSI Device Drivers**

<b>13. SCSI Device Drivers</b>	<b>299</b>
SCSI Support in Silicon Graphics Systems	300
SCSI Hardware Support	300
IRIX Kernel SCSI Support	301

---

Host Adapter Facilities	302
Purpose of the Host Adapter Driver	302
Host Adapter Concepts	303
Overview of Host Adapter Functions	305
How the Host Adapter Functions Are Found	306
Using scsi_info()	308
Using scsi_alloc()	309
Using scsi_free()	310
Using scsi_command()	311
Using scsi_abort()	316
Using scsi_reset()	317
Designing a SCSI Driver	317
SCSI Driver Initialization	317
Opening a SCSI Device	318
Accessing a SCSI Device	318
Configuring a SCSI Driver	318
Example SCSI Device Driver	318
Designing a Host Adapter Driver	323
Overview of Host Adapter Driver Architecture	323
Host Adapter Initialization	323
SCSI Reference Data	325
SCSI Error Messages	325
SCSI Error Message Tables	326
WD93 States and Phases	332

## **PART V      Network Drivers**

<b>14.</b>	<b>Network Device Drivers</b>	<b>337</b>
	Overview of Network Drivers	338
	Application Interfaces	339
	Protocol Stack Interfaces	339
	Device Driver Interfaces	340

:

---

Network Driver Interfaces	340
Kernel Facilities	341
Principal ifnet Header Files	341
Debugging Facilities	342
Information Sources	342
Network Inventory Entries	344
Multiprocessor Considerations	344
Ineffective spl() Functions	345
Multiprocessor Locking Macros	345
Compiler Flags for MP TCP/IP	345
Mutual Exclusion Macros	346
Example ifnet Driver	348

## **PART VI      PCI Drivers**

15. PCI Device Drivers	375
PCI Bus in Silicon Graphics Workstations	376
PCI Bus and System Bus	376
Buses, Slots, Cards, and Devices	378
PCI Implementation in O2 Workstations	378
Unsupported PCI Signals	379
64-bit Address and Data Support	379
Configuration Register Initialization	380
Address Spaces Supported	380
Slot Priority and Bus Arbitration	381
Interrupt Signal Distribution	382

---

Driver/Kernel Interface for PCI Access	383
Overview of PCI Driver Structure	383
Initializing and Registering the Driver	384
Attaching a Device	386
Establishing Logical Devices	393
Normal Operation	395
Detaching A Device	401
Unloading	402
PCI Function Summary	403
Example Driver	405

## **PART VII      STREAMS Drivers**

16.	<b>STREAMS Drivers</b>	499
	Driver Exported Names	500
	Streamtab Structure	500
	Driver Flag Constant	500
	Initialization Entry Points	501
	Entry Point open()	501
	Entry Point close()	502
	Put Functions wput() and rput()	502
	Service Functions rsrv() and wsrsv()	503
	Building and Debugging	504
	Special Considerations for Multiprocessing	505
	Special Considerations for IRIX	507
	Extension of Poll and Select	507
	Support for Pipes	507
	Service Scheduling	508
	Supplied STREAMS Modules	508
	No #ifdefs	509
	Different I/O Hardware Model	509
	Different Network Model	509
	Support for CLONE Drivers	510
	Summary of Standard STREAMS Functions	512

:

---

	STREAMS Modules for X Input Devices	514
	The X Input Subsystem	514
	Shared Memory Input Queue	515
	IDEV Interface	516
	Input Device Naming	516
	Opening Input Devices	517
	Device Controls	518
<b>A.</b>	<b>Silicon Graphics Driver/Kernel API</b>	<b>521</b>
	Driver Exported Names	522
	Kernel Data Structures and Declarations	523
	Kernel Functions	525
<b>B.</b>	<b>New and Updated Reference Pages</b>	<b>539</b>
	Address/Length List Reference Pages	539
	PCI Infrastructure Reference Pages	547
	<b>Glossary</b>	<b>577</b>
	<b>Index</b>	<b>591</b>

---

## List of Examples

<b>Example 2-1</b>	Testing the Hardware Inventory in a Shell Script	33
<b>Example 2-2</b>	Function Returning Type Code for CPU Module	33
<b>Example 4-1</b>	Opening and Using a Hypothetical VME Device	67
<b>Example 4-2</b>	User-Level DMA Access to VME	76
<b>Example 5-1</b>	Testing the Generic SCSI Configuration	93
<b>Example 5-2</b>	Code of the <code>testunitread00()</code> Function	106
<b>Example 5-3</b>	Program That Uses <code>dslib</code> Functions	107
<b>Example 6-1</b>	Function to Test and Set External Interrupt Pulse Width	121
<b>Example 7-1</b>	Hypothetical ULI Program	133
<b>Example 8-1</b>	Hypothetical <code>pxread()</code> entry in a Character/Block Driver	157
<b>Example 8-2</b>	<code>pxpoll()</code> Code for Hypothetical Driver	160
<b>Example 8-3</b>	Edited Fragment of <code>flash_map()</code>	164
<b>Example 8-4</b>	Hypothetical Call to <code>pollwakeup()</code>	170
<b>Example 8-5</b>	Entry Point <code>pxprint()</code>	172
<b>Example 8-6</b>	Uniprocessor Upper-Half Wait Logic	179
<b>Example 8-7</b>	Uniprocessor Interrupt Logic	179
<b>Example 9-1</b>	LIFO Queue Using Basic Locks	209
<b>Example 9-2</b>	Skeleton Code for Use of <code>SV_WAIT</code>	223
<b>Example 10-1</b>	Defining Variables in Master Descriptive File	238
<b>Example 11-1</b>	Verifying Presence of <code>symmon</code>	246
<b>Example 11-2</b>	Setting Kernel <code>putbuf</code> Size	252
<b>Example 11-3</b>	Debugging Macros Using <code>cmn_err()</code>	253
<b>Example 11-4</b>	More Elaborate Debugging Macro	253
<b>Example 11-5</b>	Invoking <code>idbg</code> Interactively	264
<b>Example 11-6</b>	Invoking <code>idbg</code> with a Log File	265
<b>Example 11-7</b>	Invoking <code>idbg</code> for a Single Command	265
<b>Example 12-1</b>	Compiling the Example Driver for a 32-bit Kernel	274

<b>Example 12-2</b>	Displaying Simulated Volume Header Using <code>idbg</code>	276
<b>Example 12-3</b>	Install Command to Create Device Special File	277
<b>Example 12-4</b>	Applying <code>prtvtoc</code> to a RAM Drive of 2 MB	278
<b>Example 12-5</b>	Making a Filesystem on a RAM Drive	278
<b>Example 12-6</b>	Mounting a RAM Drive Filesystem	279
<b>Example 13-1</b>	Storing the Adapter Type Number in <code>pfxedtinit()</code>	307
<b>Example 13-2</b>	Extracting an Adapter Number From a Minor Device Number	308
<b>Example 13-3</b>	Macro to Encapsulate a Call to <code>scsi_alloc()</code>	308
<b>Example 14-1</b>	Skeleton <code>ifnet</code> Driver	348
<b>Example 15-1</b>	Driver Registration	385
<b>Example 15-2</b>	Allocation of PCI PIO Map	389
<b>Example 15-3</b>	Reading PCI Configuration Space	391
<b>Example 15-4</b>	Setting Up a PCI Interrupt Handler	392
<b>Example 15-5</b>	Creating Logical Devices for a PCI Device	394
<b>Example 15-6</b>	Retrieving Device Information	395
<b>Example 15-7</b>	Example PCI Driver for IRIX 6.3—Descriptive File	405
<b>Example 15-8</b>	Example PCI Driver for IRIX 6.3—Configuration File	405
<b>Example 15-9</b>	Example PCI Driver for IRIX 6.3—Driver Header File	406
<b>Example 15-10</b>	Example PCI Driver for IRIX 6.3—Driver Source Code	412
<b>Example 16-1</b>	Testing Pipe Configuration	507
<b>Example B-1</b>	<code>alenlist(d4x)</code>	539
<b>Example B-2</b>	<code>alenlist_ops(d3x)</code>	542
<b>Example B-3</b>	<code>pciio(d3)</code>	547
<b>Example B-4</b>	<code>pciio_config(d3)</code>	550
<b>Example B-5</b>	<code>pciio_dma(d3)</code>	555
<b>Example B-6</b>	<code>pciio_error(d3)</code>	561
<b>Example B-7</b>	<code>pciio_get(d3)</code>	563
<b>Example B-8</b>	<code>pciio_intr(d3)</code>	567
<b>Example B-9</b>	<code>pciio_pio(d3)</code>	570

---

## List of Figures

<b>Figure 1-1</b>	CPU Access to Memory	7
<b>Figure 1-2</b>	CPU Access to Device Registers	10
<b>Figure 1-3</b>	Device Access to Memory	11
<b>Figure 1-4</b>	Device Access Through a Bus Adapter	12
<b>Figure 1-5</b>	The 32-Bit Address Space	17
<b>Figure 1-6</b>	MIPS 32-Bit Virtual Address Format	18
<b>Figure 1-7</b>	Main Parts of the 64-Bit Address Space	21
<b>Figure 1-8</b>	MIPS 64-Bit Virtual Address Format	23
<b>Figure 1-9</b>	Address Decoding for Physical Memory Access	24
<b>Figure 3-1</b>	Overview of Device Open	49
<b>Figure 3-2</b>	Overview of Device Control	50
<b>Figure 3-3</b>	Overview of Programmed Kernel I/O	52
<b>Figure 3-4</b>	Overview of Memory Mapping	53
<b>Figure 3-5</b>	Overview of DMA I/O	55
<b>Figure 5-1</b>	Bit Assignments in SCSI Device Minor Numbers	83
<b>Figure 14-1</b>	Overview of Network Architecture	338
<b>Figure 15-1</b>	PCI Bus In Relation to System Bus	377



---

## List of Tables

<b>Table 1-1</b>	CPU Modules and System Names	5
<b>Table 1-2</b>	Number of TLB Entries by Processor Type	9
<b>Table 1-3</b>	Cache Algorithm Selection	25
<b>Table 4-1</b>	VME Bus PIO Bandwidth	68
<b>Table 4-2</b>	EISA Bus PIO Bandwidth (32-Bit Slave, 33-MHz GIO Clock)	72
<b>Table 4-3</b>	EISA Bus PIO Bandwidth (16-Bit Slave, 33-MHz GIO Clock)	72
<b>Table 4-4</b>	VME Bus Bandwidth, DMA Engine, D32 Transfer	75
<b>Table 5-1</b>	Fields of the dsreq Structure	86
<b>Table 5-2</b>	Flag Values for ds_flags	87
<b>Table 5-3</b>	Return Codes From SCSI Operations	89
<b>Table 5-4</b>	SCSI Status Codes	91
<b>Table 5-5</b>	SCSI Message Byte Values	91
<b>Table 5-6</b>	Fields of the dsconf Structure	92
<b>Table 5-7</b>	dslib Function Summary	94
<b>Table 5-8</b>	Lookup Tables in dslib	99
<b>Table 6-1</b>	Functions for Outgoing External Signals	118
<b>Table 6-2</b>	Functions for Incoming External Interrupts	119
<b>Table 8-1</b>	Entry Points in Alphabetic Order	142
<b>Table 8-2</b>	Use of Driver Entry Points	143
<b>Table 9-1</b>	Functions to Manipulate Device Numbers	182
<b>Table 9-2</b>	Accessible Fields of buf_t Objects	186
<b>Table 9-3</b>	Header Files Often Used in Device Drivers	188
<b>Table 9-4</b>	Functions for Kernel Virtual Memory	190
<b>Table 9-5</b>	Functions for Allocating pollhead Structures	192
<b>Table 9-6</b>	Functions for Allocating buf_t Objects and Buffers	193
<b>Table 9-7</b>	Functions for Suballocation	193
<b>Table 9-8</b>	Functions for General Data Transfer	195

<b>Table 9-9</b>	Functions Moving Data Using <code>uio_t</code>	197
<b>Table 9-10</b>	Functions to Test Physical Addresses	198
<b>Table 9-11</b>	Functions to Manipulate a <code>vhandl_t</code> Object	199
<b>Table 9-12</b>	Functions to Convert Bytes to Sectors or Pages	201
<b>Table 9-13</b>	Functions Related to Physical Memory	202
<b>Table 9-14</b>	Functions to Map Buffer Pages	202
<b>Table 9-15</b>	Functions Related to Cache Coherency	203
<b>Table 9-16</b>	Functions for User Process Management	205
<b>Table 9-17</b>	Functions for Basic Locks	208
<b>Table 9-18</b>	Functions for Mutex Locks	210
<b>Table 9-19</b>	Functions for Sleep Locks	212
<b>Table 9-20</b>	Functions for Reader/Writer Locks	214
<b>Table 9-21</b>	Functions to Set Interrupt Levels	215
<b>Table 9-22</b>	Functions for Timed Delays	216
<b>Table 9-23</b>	Functions for Synchronizing Block I/O	219
<b>Table 9-24</b>	Functions for Synchronization: sleep/wakeup	221
<b>Table 9-25</b>	Functions for Synchronization: Synchronization Variables	222
<b>Table 9-26</b>	Functions for Semaphores	224
<b>Table 10-1</b>	Compiler Variables Tested by System Header Files	231
<b>Table 10-2</b>	Compiler Options for 32-Bit Kernel Modules	232
<b>Table 10-3</b>	Compiler Options for 64-Bit Kernel Modules	233
<b>Table 11-1</b>	Commands for Symbol Conversion and Lookup	258
<b>Table 11-2</b>	Commands to Control Execution	259
<b>Table 11-3</b>	Commands to Manage Virtual Memory	261
<b>Table 11-4</b>	Commands to Display Memory	262
<b>Table 11-5</b>	Utility Commands	263
<b>Table 11-6</b>	Commands to Display Memory and Symbols	267
<b>Table 11-7</b>	Commands to Display Process Information	267
<b>Table 11-8</b>	Commands to Display Locks and Semaphores	269
<b>Table 11-9</b>	Commands to Display I/O Status	269
<b>Table 11-10</b>	Commands to Display <code>buf_t</code> Objects	270
<b>Table 11-11</b>	Commands to Display STREAMS Structures	270
<b>Table 11-12</b>	Commands to Display Network-Related Structures	271

---

<b>Table 13-1</b>	Host Adapter Driver Classes	303
<b>Table 13-2</b>	Host Adapter Function Summary	305
<b>Table 13-3</b>	Input Fields of the <code>scsi_request</code> Structure	311
<b>Table 13-4</b>	Values for the <code>sr_flags</code> Field of a <code>scsi_request</code>	312
<b>Table 13-5</b>	Values Returned From a SCSI Command	314
<b>Table 13-6</b>	Software Status Values From a SCSI Request	314
<b>Table 13-7</b>	SCSI Status Bytes	315
<b>Table 13-8</b>	Host Adapter Status After a SCSI Request	316
<b>Table 13-9</b>	Adapter Error Codes	326
<b>Table 13-10</b>	Primary Sense Key Error Table	327
<b>Table 13-11</b>	Additional Sense Code Table	328
<b>Table 13-12</b>	SCSI State Error Messages	332
<b>Table 14-1</b>	Important Reference Pages Related to Network Drivers	343
<b>Table 14-2</b>	Mutual Exclusion Macros for <code>ifnet</code> Drivers	346
<b>Table 15-1</b>	PCI Interrupt Distribution to System Interrupt Numbers	382
<b>Table 15-2</b>	Functions for PIO Maps for PCI	396
<b>Table 15-3</b>	Least Significant Address Bytes for Short PIO	398
<b>Table 15-4</b>	Functions for Simple DMA Maps for PCI	399
<b>Table 15-5</b>	Functions for DMA Using Address-Length Lists	400
<b>Table 15-6</b>	Functions for Managing PCI Interrupt Handlers	401
<b>Table 15-7</b>	PCI-Related Kernel Functions	403
<b>Table 16-1</b>	Multiprocessing STREAMS Functions	506
<b>Table 16-2</b>	Kernel Entry Points	512
<b>Table A-1</b>	Driver Exported Names	522
<b>Table A-2</b>	Device Driver Interface Objects	523
<b>Table A-3</b>	STREAMS Driver Interface Objects	524
<b>Table A-4</b>	Kernel Functions	525



---

## About This Guide

This guide describes the ways in which hardware devices are integrated into and controlled from a Silicon Graphics® computer system running the IRIX™ operating system version 6.3 for O2.

**Note:** This edition applies only to IRIX 6.3 for O2, and discusses only hardware supported by that system version. If your device driver will work with a different release, you should use the version of this manual appropriate to that release (see “Internet Resources” on page xxvii).

Three general classes of device-control software exist in an IRIX system: process-level drivers, kernel-level drivers, and STREAMS drivers.

- A process-level driver executes as part of a user-initiated process. An example is the use of the *dslib* library to control a SCSI device from a user program.
- A kernel-level driver is loaded as part of the IRIX kernel and executes in the kernel address space, controlling devices in response to calls to its read, write, and ioctl (control) entry points.
- A STREAMS driver is dynamically loaded into the kernel address space to monitor or modify a stream of data passing between a device and a user process.

All three classes are discussed in this guide, although the greatest amount of attention is given to kernel-level drivers.

### What You Need to Know

In order to write a process-level driver you must be an experienced C programmer with a thorough understanding of the use of IRIX system services and, of course, detailed knowledge of the device to be managed.

In order to write a kernel-level driver or a STREAMS driver you must be an experienced C programmer who knows UNIX® system administration, and especially IRIX system administration, and who understands the concepts of UNIX device management.

## What This Guide Contains

This guide is divided into the following major parts.

Part I, “IRIX Device Integration”	How devices are attached to Silicon Graphics computers, configured to IRIX, and initialized at boot time.
Part II, “Device Control From Process Space”	Details of user-level handling of PCI devices, and SCSI control using <i>dslib</i> .
Part III, “Kernel-Level Drivers”	How kernel-level drivers are designed, compiled, loaded, and tested. Survey of kernel services for drivers.
Part IV, “SCSI Device Drivers”	Kernel-level drivers for the SCSI bus.
Part V, “Network Drivers”	Kernel-level drivers for network interfaces.
Part VI, “PCI Drivers”	Kernel-level drivers for the PCI bus.
Part VII, “STREAMS Drivers”	Design of STREAMS drivers.
Appendix A, “Silicon Graphics Driver/Kernel API”	Summary of kernel functions with compatibility notes.

In the printed book, you can locate these parts using the part-tabs printed in the margins. Using IRIS InSight™, each part is a top-level division in the clickable table of contents, or you can jump to any part by clicking the blue cross-references in the list above.

## Other Sources of Information

### Developer Program

Information and support are available through the Silicon Graphics Developer Program. The Developer Toolbox CDROM contains numerous code examples. To join the program, contact the Developer Response Center at (800) 770-3033 or send e-mail to [devprogram@sgi.com](mailto:devprogram@sgi.com).

### Internet Resources

A great deal of useful material can be found on the internet. Some starting points are in the following list.

Earlier versions of this book as well as all other SGI technical manuals to read or download.	<a href="http://www.sgi.com/Technology/TechPubs/">http://www.sgi.com/Technology/TechPubs/</a>
SGI patches, examples, and other material.	<a href="ftp://ftp.sgi.com">ftp://ftp.sgi.com</a>
Network of pages of information about Silicon Graphics and MIPS® products	<a href="http://www.sgi.com">http://www.sgi.com</a>
Text of all Internet RFC documents.	<a href="ftp://ds.internic.net/rfc/">ftp://ds.internic.net/rfc/</a>
Computer graphics pointers at the UCSC Perceptual Science Laboratory.	<a href="http://mambo.ucsc.edu/psl/cg.html">http://mambo.ucsc.edu/psl/cg.html</a>
Pointers to binaries and sources at The National Research Council of Canada's Institute For Biодiagnostics.	<a href="http://zeno.ibd.nrc.ca:80/~sgi/">http://zeno.ibd.nrc.ca:80/~sgi/</a>
A Silicon Graphics "meta page" at the Georgia Institute of Technology College of Computing.	<a href="http://www.cc.gatech.edu/service/sgimeta.html">http://www.cc.gatech.edu/service/sgimeta.html</a>
Dazzling Silicon Graphics "meta page" at NASA in Huntsville, AL.	<a href="http://chernobog.msfc.nasa.gov/SGI/html/SGI.html">http://chernobog.msfc.nasa.gov/SGI/html/SGI.html</a>
Complete SCSI-2 standard in HTML.	<a href="http://abekas.com:8080/SCSI2/">http://abekas.com:8080/SCSI2/</a>
IEEE Catalog and worldwide ordering information.	<a href="http://stdsbbs.ieee.org:70/0/pub/htmlfiles/stctoc.htm">http://stdsbbs.ieee.org:70/0/pub/htmlfiles/stctoc.htm</a>
MIPS processor manuals in HTML form.	<a href="http://www.mips.com/">http://www.mips.com/</a>
Home page of the PCI bus standardization organization	<a href="http://www.pcisig.com">http://www.pcisig.com</a>

## Standards Documents

The following documents are the official standard descriptions of buses:

- *PCI Local Bus Specification, Version 2.1*, available from the PCI Special Interest Group, P.O. Box 14070, Portland, OR 97214 (fax: 503-234-6762)
- *ANSI/IEEE standard 1014-1987 (VME Bus)*, available from IEEE Customer Service, 445 Hoes Lane, PO Box 1331, Piscataway, NJ 08855-1331 (but see also “Internet Resources” on page xxvii).

## Important Reference Pages

The following reference pages contain important details about software tools and practices that you need.

getinvent(3)	The interface to the inventory database
hinvt(1)	The use of the inventory display command
intro(7)	The conventions used for special device filenames
MAKEDEV(1)	The use of the program that creates device special files
master(4)	Syntax of files in <i>/var/sysgen/master.d</i>
prom(1)	Commands of the “miniroot” and other features of the boot PROM, which you use to bring up the system when testing a new device driver
system(4)	Syntax of files in <i>/var/sysgen/system/*.sm</i>
udmalib(3)	Functions for performing user-level DMA from VME.
uli(3)	Functions for registering and using a user-level interrupt handler.
usrvme(7)	Naming conventions for mappable VME device special files.

## Additional Reading

The following books, obtainable from Silicon Graphics, can be helpful when designing or testing a device driver.

- *MIPSpro Compiling and Performance Tuning Guide*, document number 007-2360-*nnn*, tells how to use the C compiler and related tools.
- *MIPSpro Assembly Language Programmer's Guide*, document number 007-2418-*nnn*, tells how to compile assembly-language modules.
- *MIPSpro 64-Bit Porting and Transition Guide*, document number 007-2391-*nnn*, documents the implications of the 64-bit execution mode for user programs.
- *MIPSpro N32 ABI Handbook*, document number 007-2816-*nnn*, gives details of the code generated when the *-n32* compiler option is used.
- *Topics in IRIX Programming*, document number 008-2478-*nnn*, documents some of the sophisticated services offered by the IRIX kernel to user-level programs.
- *MIPS R4000 User's Manual* (2nd ed.) by Joe Heinrich, document number 007-2489-001, gives detailed information on the MIPS instruction set and hardware registers for the processor used in many Silicon Graphics computer systems (also available as HTML on <http://www.mips.com/>).
- *MIPS R10000 User's Manual* by Joe Heinrich gives detailed information on the MIPS instruction set and hardware registers for the processor used in certain high-end systems. Available only in HTML form from <http://www.mips.com/>.
- *IRIX Administration: System Configuration and Operation*, document number 007-2859-*nnn*, describes the basic administrative tools for configuring, operating, and tuning IRIX.
- *IRIX Administration: Disks and File Systems*, document number 007-2825-*nnn*, describes the configuration of new disk subsystems and the management of logical volumes and file systems.
- *IRIX Administration: Peripheral Devices*, document number 007-2861-*nnn*, describes the administration of tapes, printers, and other devices.

The following books, obtainable from bookstores or libraries, can also be helpful.

- Egan, Janet I., and Thomas J. Teixeira. *Writing a UNIX Device Driver*. John Wiley & Sons, 1992.
- Leffler, Samuel J., et alia. *The Design and Implementation of the 4.3BSD UNIX Operating System*. Palo Alto, California: Addison-Wesley Publishing Company, 1989.

- A. Silberschatz, J. Peterson, and P. Galvin. *Operating System Concepts*, Third Edition. Addison Wesley Publishing Company, 1991.
- Heath, Steve. *VMEbus User's Handbook*. CRC Press, Inc, 1989. ISBN 0-8493-7130-9.
- *Device Driver Reference, UNIX SVR4.2*, UNIX Press 1992.
- *UNIX System V Release 4 Programmer's Guide*, UNIX SVR4.2. UNIX Press, 1992.
- *STREAMS Modules and Drivers, UNIX SVR4.2*, UNIX Press 1992. ISBN 0-13-066879.

## Conventions Used in This Guide

Special terms and special kinds of words are indicated with the following typographical conventions:

Data structures, variables, function arguments, and macros.	The <i>dsiovec</i> structure has members <i>iov_base</i> and <i>iov_len</i> . Use the <i>IOVLEN</i> macro to access them.
Kernel and library functions and functions in examples.	When successful, <b>v_mapphys()</b> returns 0.
Driver entry point names that must be completed with a unique prefix string.	The <b>munmap()</b> system function calls the <i>pfxunmap()</i> entry point.
Files and directories.	Device special files are in <i>/dev</i> , and are created using the <i>/dev/MAKEDEV</i> script.
First use of terms defined in the glossary (see "Glossary" on page 577).	The <i>inode</i> of a <i>device special file</i> contains the <i>major device number</i> .
Literal quotes of code examples.	The SCSI driver's prefix is <i>scsi_</i> .

PART ONE

## IRIX Device Integration

**Chapter 1:** Physical and Virtual Memory

An overview of physical memory, virtual address space management, and device addressing in Silicon Graphics/MIPS systems.

**Chapter 2:** Device Configuration

How IRIX locates devices, and how devices are represented in software.

**Chapter 3:** Device Control Software

A survey of the ways in which you can control devices under IRIX, from user-level processes and from kernel-level drivers of different kinds.



---

## Physical and Virtual Memory

This chapter gives an overview of the management of physical and virtual memory in the MIPS® R4x00®, R5000™, R8000™, and R10000™ processors. Access to physical devices is included in this topic, because device registers and bus attachments are accessed using physical memory addresses.

This information is only of academic interest if you intend to control a device from a user-level process. When you are designing a kernel-level driver, this information helps you understand the operation of the kernel functions that you call on, and the constraints on their operations. ( See Chapter 3, “Device Control Software,” for the difference between these two types of drivers.)

The following main topics are covered in this chapter.

- “Physical Address Space” on page 4 describes the range and meaning of address numbers on the hardware bus.
- “CPU Access to Memory and Devices” on page 5 summarizes the hardware architecture by which the CPU accesses memory.
- “The 32-Bit Address Space” on page 16 describes the divisions of the 32-bit virtual address space and their uses.
- “The 64-Bit Address Space” on page 20 describes the divisions of the 64-bit virtual address space and their uses.
- “Device Driver Use of Memory” on page 26 describes the techniques and rules for how kernel-level device drivers allocate and use memory.

**Note:** This chapter tells only enough about memory access and cache management to explain the rules of the driver/kernel interface. For complete details on the MIPS hardware processors, see the hardware manuals listed under “Additional Reading” on page xxix.

## Physical Address Space

The CPU emits physical addresses in order to select RAM, ROM, device registers, and bus attachments. Physical addresses start at 0 and can (in some systems) go as high as  $2^{40}$ . This range includes  $1.1 \times 10^{12}$  unique numbers, or 1,024 *gigabytes* (GB), or 1 *terabyte* (TB).

Software never uses physical addresses directly. Kernel-level software can access physical memory and devices using indirect addressing discussed later.

## Physical Device Addresses

The MIPS processor architecture has no I/O instructions. Certain ranges of physical addresses are reserved as device addresses. That is, when the CPU emits one of these addresses, the hardware decodes it as an access to a particular device or bus attachment, instead of an access to memory.

Each Silicon Graphics computer model has a particular set of device addresses. The choice of device addresses is part of the architecture of the whole computer system; it is not designed into the processor chip.

For example, the relationship between physical address space and the PCI bus is discussed under “Address Spaces Supported” on page 380.

## Physical Memory Addresses

Some physical addresses are decoded to select memory hardware. Physical memory includes ROM as well as RAM. Each block of physical memory has a range of physical addresses. The physical addresses where RAM or ROM can be found depend on the particular computer system.

Physical memory does not necessarily occupy sequential addresses. There can be (and often are) gaps, ranges of physical addresses that do not relate to either memory or devices, between ROM addresses and RAM addresses. In most systems, all RAM is given a single sequential span of physical addresses. However, this is not a requirement. Blocks of RAM addresses can also be separated by gaps that are not populated with memory. Since all software uses virtual addresses, software always sees a sequential range of addresses without gaps.

## CPU Access to Memory and Devices

Each Silicon Graphics computer system has one or more CPU modules. The CPU reads memory or a device by placing an address on a system bus, and receiving data back from the addressed memory or device. Access to memory can pass through multiple levels of cache.

### CPU Modules

A CPU is a hardware module containing a MIPS processor chip such as the R8000, together with system interface chips and possibly a secondary cache. Silicon Graphics CPU modules have model designation of the form IP $nn$ ; for example, the IP22 module is used in the Indy™ workstation. The CPU modules supported by IRIX 6.3 for O2 are listed in Table 1-1.

**Table 1-1** CPU Modules and System Names

Module	MIPS Processor	System Families
IP17	R4000	Crimson™
IP19	R4x00	Challenge (other than S model), Onyx
IP20	R4x00	Indigo®
IP21	R8000	POWER Challenge™, POWER Onyx™
IP22	R4x00	Indigo <sup>2</sup> , Indy, Challenge S
IP25	R10000	POWER Challenge R10000
IP26	R8000	POWER Indigo <sup>2</sup> ™
IP32	R10000	O2

Modules with the same IP designation can be ordered in a variety of clock speeds, and they can differ in other ways. Also, the choice of graphics hardware is independent of the CPU model. However, all these CPUs are identical as seen from software.

### Interrogating the CPU Type

At the interactive command line, you can determine which CPU module a system uses with the command

```
hinv -c processor
```

Within a shell script, it is more convenient to process the terse output of

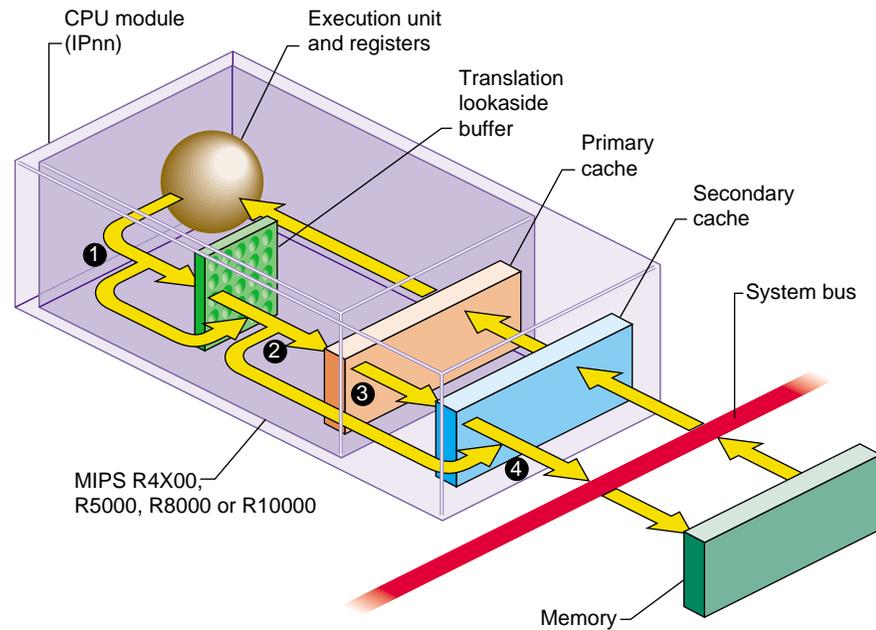
```
uname -m
```

(See the `uname(1)` and `hinv(1)` reference pages.)

Within a program, you can get the CPU model using the `getinvent()` function. For an example, see “Testing the Inventory In Software” on page 33.

### CPU Access to Memory

The CPU generates the address of data that it needs—the address of an instruction to fetch, or the address of an operand of an instruction. It requests the data through a mechanism that is depicted in simplified form in Figure 1-1.



**Figure 1-1** CPU Access to Memory

1. The address of the needed data is formed in the processor execution or instruction-fetch unit. Most addresses are then mapped from virtual to real through the Translation Lookaside Buffer (TLB). Certain ranges of addresses are not mapped, and bypass the TLB.
2. Most addresses are presented to the *primary cache*, the cache in the processor chip. If a copy of the data with that address is found, it is returned immediately. Certain address ranges are never cached; these addresses pass directly to the bus.
3. When the primary cache does not contain the data, the address is presented to the secondary cache. If it contains a copy of the data, the data is returned immediately. The size and the architecture of the secondary cache differ from one CPU model to another, and some CPUs do not have a secondary cache.
4. The address is placed on the system bus. The memory module that recognizes the address places the data on the bus.

## Processor Operating Modes

The MIPS processor under IRIX operates in one of two modes: kernel and user. The processor enters the more privileged kernel mode when an interrupt, a system instruction, or an exception occurs. It returns to user mode only with a “Return from Exception” instruction.

Certain instructions cannot be executed in user mode. Certain segments of memory can be accessed only in kernel mode, and other segments only in user mode.

## Virtual Address Mapping

The MIPS processor contains an array of Translation Lookaside Buffer (TLB) entries that map, or translate, virtual addresses to physical ones. Most memory accesses are first mapped by reference to the TLB. This permits the IRIX kernel to implement *virtual memory* for user processes, and permits it to relocate parts of the kernel itself. The translation scheme is summarized in the following sections and covered in detail in the hardware manuals listed under “Additional Reading” on page xxix.

### TLB Misses and TLB Sizes

Each TLB entry describes a segment of memory containing two adjacent *pages*. When the input address falls in a page described by a TLB entry, the TLB supplies the physical memory address for that page. The translated address, now physical instead of virtual, is passed on to the cache, as shown in Figure 1-1 on page 7.

When the input address is not covered by any active TLB entry, the MIPS processor takes a “TLB miss” interrupt to an IRIX kernel routine. The kernel routine inspects the address. When the address has a valid translation to some page in the address space, the kernel loads a TLB entry to describe that page, and restarts the instruction.

The size of the TLB is important for performance. The size of the TLB in different processors is shown in Table 1-2.

**Table 1-2** Number of TLB Entries by Processor Type

Processor Type	Number of TBL Entries
R4x00	96
R5000	96
R8000	384
R10000	128

## Address Space Creation

There are not sufficient TLB entries to describe all the address space of every process. The IRIX kernel creates a page table for each process, containing one entry for each virtual memory page in the address space of that process. Whenever an executing program refers to an address for which there is no current TLB entry, the processor traps to the handler for the TLB miss exception. The exception handler loads one TLB entry from the appropriate page table entry of the current process, in order to describe the needed virtual address. Then it resumes execution with the failed instruction.

The kernel maintains a page table in kernel memory for each process, and a page table for the kernel virtual address space as well. In order to extend a virtual address space, the kernel takes the following two steps.

- It allocates unused page table entries to describe the needed pages. This defines the virtual addresses the pages will have.
- It allocates page frames in memory to contain the pages themselves, and puts their physical addresses in the page table entries.

## Address Exceptions

When the CPU requests an invalid address—because the processor is in the wrong mode, or an address does not translate to a valid location in the address space, or an address refers to hardware that does not exist in the system—an addressing exception occurs. The processor traps to a particular address in the kernel.

An addressing exception can also be detected while handling a TLB miss. If there is no page table entry assigned for the desired address, that address is not part of the address space of the process.

When a user-mode process caused the addressing exception, the kernel sends the process a SIGSEGV (see the signal(5) reference page), usually causing a segmentation fault. When kernel-level code such as a device driver causes the exception, the kernel executes a “panic,” taking a crash dump and shutting down the system.

### CPU Access to Device Registers

The CPU accesses a device register using the mechanism illustrated in Figure 1-2. Access to device registers is always uncached. It is not affected by considerations of cache coherency in any system (see “Cache Use and Cache Coherency” on page 15).

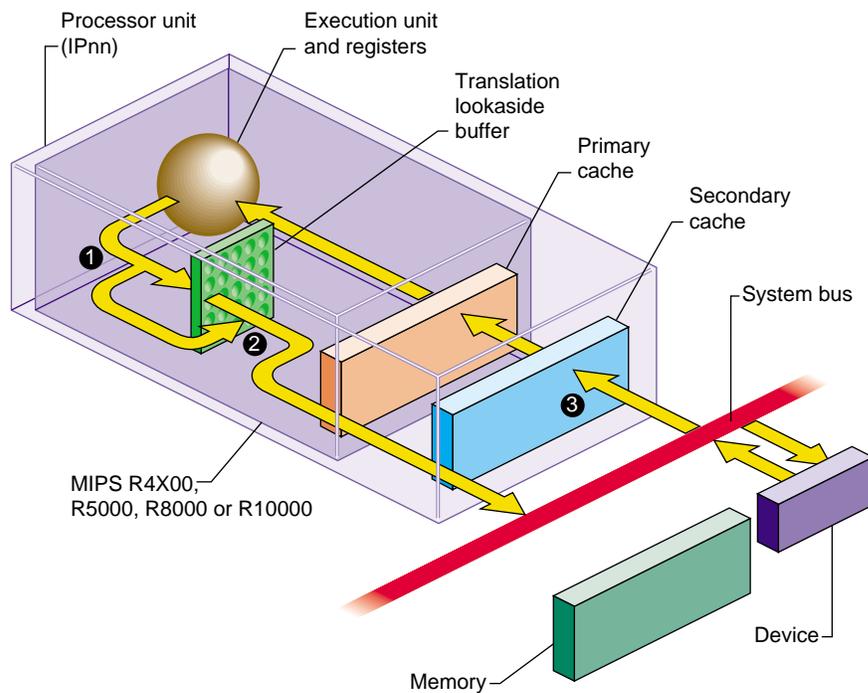


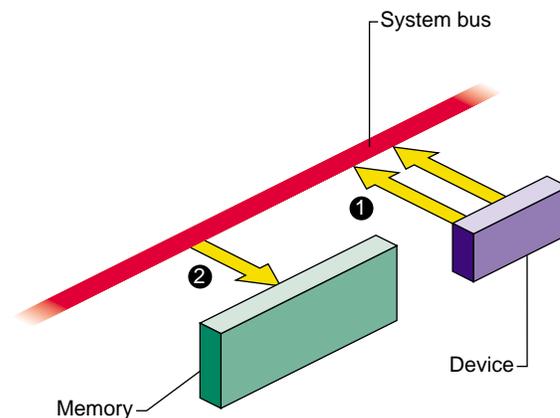
Figure 1-2 CPU Access to Device Registers

1. The address of the device is formed in the Execution unit. It may or may not be an address that is mapped by the TLB.
2. A device address, after mapping if necessary, always falls in one of the ranges that is not cached, so it passes directly to the system bus.
3. The device or bus attachment recognizes its physical address and responds with data.

### Direct Memory Access

Some devices can perform *direct memory access (DMA)*, in which the device itself, not the CPU, reads or writes data into memory. A device that can perform DMA is called a *bus master* because it independently generates a sequence of bus accesses without help from the CPU.

In order to read or write a sequence of memory addresses, the bus master has to be told the proper physical address range to use. This is done by storing a bus address number into a device registers from the CPU. When the device's DMA address registers are loaded, it can access memory through the system bus, as shown in Figure 1-3.



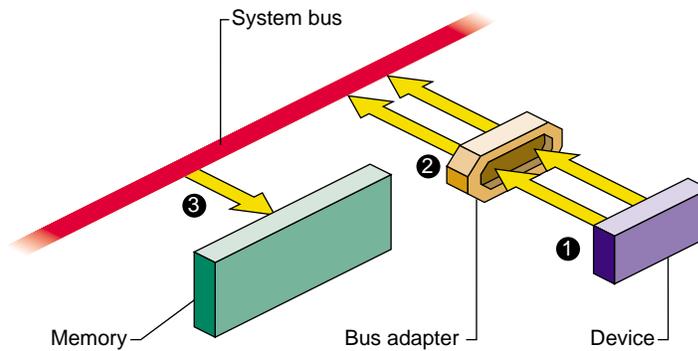
**Figure 1-3** Device Access to Memory

1. The device places the next physical address, and data, on the system bus.
2. The memory module stores the data.

When a device is programmed with an invalid physical address, the result is a bus error interrupt.

### PIO Addresses and DMA Addresses

Figure 1-3 is too simple for some devices that are attached through a bus adapter. A bus adapter connects a bus of a different type to the system bus, as shown in Figure 1-4.



**Figure 1-4** Device Access Through a Bus Adapter

For example, the PCI adapter connects a PCI bus to the system bus. Multiple PCI devices can be plugged into the PCI bus, and can use the PCI bus to read and write. The bus adapter translates the PCI bus protocol into the system bus protocol. (For details on the PCI bus adapter, see Chapter 15, “PCI Device Drivers.”)

Each bus has address lines that carry the address values used by devices on the bus. These bus addresses are not related to the physical addresses used on the system bus. The issue of bus addressing is made complicated by three facts:

- Bus-master devices independently generate memory-read and memory-write commands that are intended to access system memory.
- The bus adapter can translate addresses between addresses on the bus it manages, and different addresses on the system bus it uses.
- The translation done by the bus adapter can be programmed dynamically, and can change from one I/O operation to another.

This subject can be simplified by dividing it into two distinct subjects: PIO addressing, used by the CPU to access a device, and DMA addressing, used by a bus master to access memory. These addressing modes need to be treated differently.

### **PIO Addressing**

Programmed I/O (PIO) is the term for a load or store instruction executed by the CPU that names an I/O device as its operand. As described earlier (“CPU Access to Device Registers”), the CPU places a physical address on the system bus. The bus adapter repeats the read or write command on its bus, but not necessarily using the same address bits as the CPU put on the system bus.

One task of a bus adapter is to translate between the physical addresses used on the system bus and the addressing scheme used within the proprietary bus. The address placed on the target bus is not necessarily the same as the address generated by the CPU. The translation is done differently with different bus adapters and in different system models.

With some bus types in some systems, the translation is hard-wired. For a simple example, the address translation from the Indigo<sup>2</sup> system bus to the EISA bus is hardwired. In an Indigo<sup>2</sup>, CPU access to a physical address of 0x0000 4010 is always translated to location 0x0010 in the I/O address space of slot 4 of the EISA bus.

With the more sophisticated PCI buses, the translation is dynamic. This bus supports bus address spaces that are as large or larger than the physical address space of the system bus. It is impossible to hard-wire a translation of the entire bus address space.

In order to use a dynamic PIO address, a device driver creates a software object called a PIO map that represents that portion of bus address space that contains the device registers the driver uses. When the driver wants to use the PIO map, the kernel dynamically sets up a translation from an unused part of physical address space to the needed part of the bus address space. The driver extracts an address from the PIO map and uses it as the base for accessing the device registers. PIO maps are discussed in Chapter 15, “PCI Device Drivers.”

### **DMA Addressing**

A bus-master device on the PCI bus can be programmed to perform transfers to or from memory independently and asynchronously. A bus master is programmed (using PIO access) with a starting bus address and a length. The bus master generates a series of

memory-read or memory-write operations to successive addresses. But what *bus* addresses should it use in order to store into the proper *memory* addresses?

The bus adapter translates the addresses used on the proprietary bus to corresponding addresses on the system bus. Considering Figure 1-4, the operation of a DMA device is as follows:

1. The device places a bus address and data on the PCI bus (or the EISA, VME, or GIO bus in other hardware architectures).
2. The bus adapter translates the address to a meaningful physical address, and places that address and the data on the system bus.
3. The memory modules stores the data.

The translation of bus virtual to physical addresses is fixed for some bus types in some systems. In most systems, however, the kernel can program the bus adapter to translate bus addresses to different physical addresses. Dynamic translation is necessary because the bus address space is as large or larger than physical address space, and only some portions of bus address space can be mapped at any one time—different portions depending on what bus masters are active.

For example, the VME bus protocol used in the Silicon Graphics Challenge systems defines several different address spaces: A16, 16-bit addresses; A32, 32-bit addresses; and so on. These addresses have no direct relationship to the physical addresses used on the system bus. The VME bus adapter in a Challenge or Onyx system can be programmed to place 15 different “windows” of VME address space at different locations in physical address space at any time.

In order to create a mapping for DMA, a device driver creates a software object called a DMA map. Using kernel functions, the driver establishes the range of memory addresses that the bus master wants to access—typically the address of an I/O buffer. When the driver activates the DMA map, the kernel sets up a dynamic mapping between some range of bus addresses and the desired range of memory space. The driver extracts from the DMA map the starting bus address, and (using PIO) programs that bus address into the bus master device.

The management of DMA maps is discussed in Chapter 15, “PCI Device Drivers.”

## Cache Use and Cache Coherency

The primary and secondary caches shown in Figure 1-1 on page 7 are essential to CPU performance. There is an order of magnitude difference in the speed of access between cache memory and main memory. Execution speed remains high only as long as a very high proportion of memory accesses are satisfied from the primary or secondary cache.

The use of caches means that there are often multiple copies of data: a copy in main memory, a copy in the secondary cache (when one is used) and a copy in the primary cache. Moreover, a multiprocessor system has multiple CPU modules like the one shown, and there can be copies of the same data in the cache of *each* CPU.

The problem of *cache coherency* is to ensure that all cache copies of data are true reflections of the data in main memory. Different Silicon Graphics systems use different hardware designs to achieve cache coherency.

In most cases, cache coherence is achieved by the hardware, without any effect on software. In a few cases, specialized software, such as a kernel-level device driver, must take specific steps to maintain cache coherency.

### Cache Coherency in Multiprocessors

Multiprocessor systems have more complex cache coherency protection because it is possible to have data in multiple caches. In a multiprocessor system, the hardware ensures that cache coherency is maintained under all conditions, including DMA input and output, without action by the software. However, in some systems the cache coherency hardware works correctly only when a DMA buffer is aligned on a cache-line-sized boundary. You ensure this by using the `KM_CACHEALIGN` flag when allocating buffer space with `kmem_alloc()` (see the `kmem_alloc(D3)` reference page).

### Cache Coherency in Uniprocessors

In some uniprocessor systems, it is possible for the CPU cache to have newer information than appears in memory. This is a problem only when a bus master device is going to perform DMA. If the bus master reads memory, it can get old data. If it writes memory, the input data can be destroyed when the CPU writes the modified cache line back to memory.

In systems where this is possible, a device driver calls a kernel function to ensure that all cached data has been written to memory prior to DMA output (the `dki_cache_wb(D3)` reference page). The device driver calls a kernel function to ensure that the CPU receives the latest data following a DMA input (see the `dki_cache_inval(D3)` reference page). In a multiprocessor these functions do nothing, but it is always safe to call them.

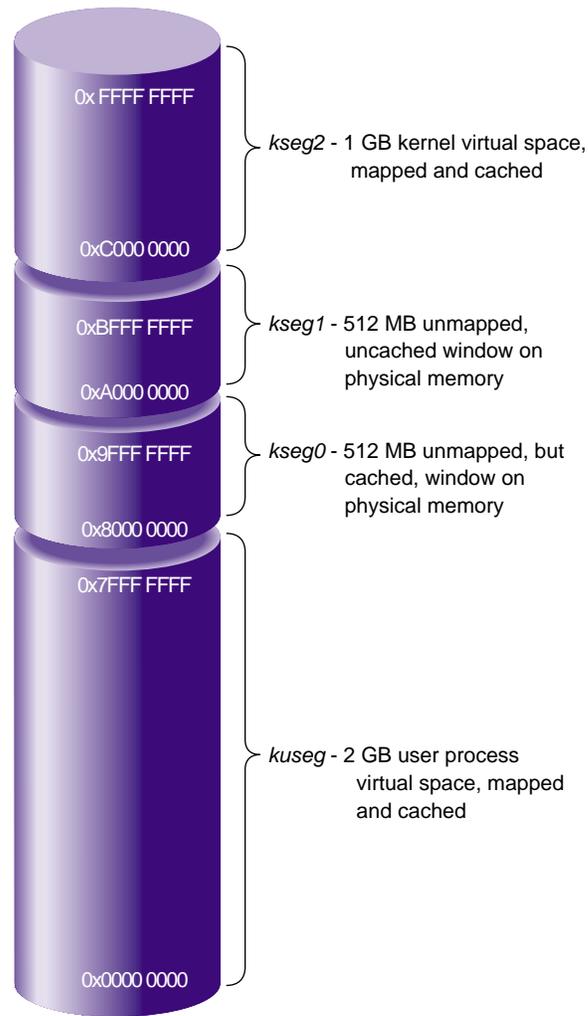
## The 32-Bit Address Space

The MIPS processors can operate in one of two address modes: 32-bit and 64-bit. The choice of address mode is independent of other features of the instruction set architecture such as the number of available registers and the precision of integer arithmetic. For example, programs compiled to the n32 binary interface use 32-bit addresses but 64-bit integers. The implications for user programs are documented in manuals listed under “Additional Reading” on page xxix.

The addressing mode can be switched dynamically; for example, the IRIX kernel can operate with 64-bit addresses, but the kernel can switch to 32-bit address when it dispatches a user program that was compiled for that mode. The 32-bit address space is the range of all addresses that can be used when in 32-bit mode. This space is discussed first because it is simpler and more familiar than the 64-bit space.

### Segments of the 32-bit Address Space

When operating in 32-bit mode, the MIPS architecture uses addresses that are 32-bit unsigned integers from `0x0000 0000` to `0xFFFF FFFF`. However, this address space is not uniform. The MIPS hardware divides it into segments, and treats each segment differently. The ranges are shown graphically in Figure 1-5.



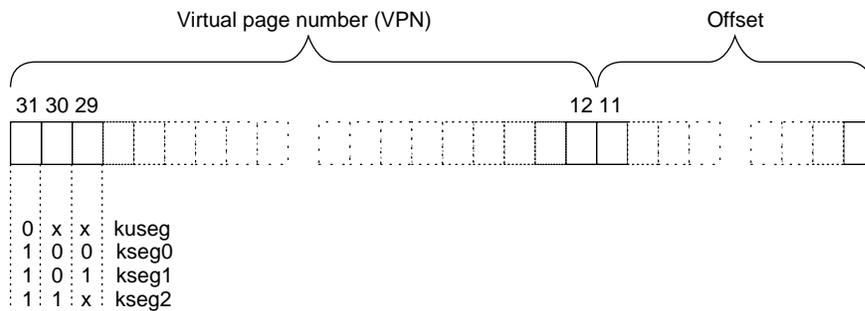
**Figure 1-5** The 32-Bit Address Space

The address segments differ in three characteristics:

- whether access to an address is mapped; that is, passed through the translation lookaside buffer (TLB)
- whether an address can be accessed when the CPU is operating in user mode or in kernel mode
- whether access to an address is cached; that is, looked up in the primary and secondary caches before it is sent to main memory

### Virtual Address Mapping

In the mapped segments, each 32-bit address value is treated as shown in Figure 1-6.



**Figure 1-6** MIPS 32-Bit Virtual Address Format

The three most significant bits of the address choose the segment among those drawn in Figure 1-5. When bit 31 is 0, bits 30:12 select a *virtual page number* (VPN) from  $2^{19}$  possible pages in the address space of the current user process. When bits 31:30 are 11, bits 29:12 select a VPN from  $2^{18}$  possible pages in the kernel virtual address space.

### User Process Space—kuseg

The total 32-bit address space is divided in half. Addresses with a most significant bit of 0 constitute the 2 GB user process space. When executing in user mode, only addresses in *kuseg* are valid; an attempt to use an address with bit 31=1 causes an addressing exception.

Access to *kuseg* is always mapped through the TLB. The kernel creates a unique address space for each user process. Of the  $2^{19}$  possible pages in an address space, most are typically unassigned—few processes ever occupy more than a fraction of *kuseg*—and many are shared pages of program text from dynamic shared objects (DSOs) that are mapped into the address space of every process that needs them.

### **Kernel Virtual Space—*kseg2***

When bits 31:30 are 11, access is to kernel virtual memory. Only code that is part of the kernel can access this space. References to this space are translated through the TLB. The kernel uses the TLB to map kernel pages in memory as required, possibly in noncontiguous locations. Although pages in kernel space are mapped, they are always associated with real memory. Kernel memory is never paged to secondary storage.

This is the space in which the IRIX kernel allocates such objects as stacks, user page tables, and per-process data that must be accessible on context switches. This area contains automatic variables declared by loadable device drivers. It is the space in which kernel-level device drivers allocate memory. Since kernel space is mapped, addresses in *kseg2* that are apparently contiguous need not be contiguous in physical memory. However, a device driver can allocate space that is both logically and physically contiguous, when that is required (see for example the `kmem_alloc(D3)` reference page).

### **Cached Physical Memory—*kseg0***

When address bits 31:29 contain 100, access is directed to physical memory through the cache. If the addressed location is not in the cache, bits 28:0 are placed on the system bus as a physical memory address, and the data presented by memory or a device is returned. *Kseg0* contains the exception address to which the MIPS processor branches it when it detects an exception such as an addressing exception or TLB miss.

Since only 29 bits are available for mapping physical memory, only 512 MB of physical memory space can be accessed through this segment in 32-bit mode. Some of this space must be reserved for device addressing. It is possible to gain cached access to wider physical addresses by mapping through the TLB into *kseg2*, but systems that need access to more physical memory typically run in 64-bit mode (see “Cache-Controlled Physical Memory—`xkphys`” on page 24).

### Uncached Physical Memory—*kseg1*

When address bits 31:29 contain 101, access is directly to physical memory, bypassing the cache. Bits 28:0 are placed on the system bus for memory or device transfer.

The kernel refers to *kseg1* when performing PIO to devices because loads or stores from device registers should not pass through cache memory. The kernel also uses *kseg1* when operating on certain data structures that might be *volatile*. Kernel-level device drivers sometimes need to write to uncached memory, and must take special precautions when doing so (see “Uncached Memory Access in the IP26 CPU” on page 29).

Portions of *kseg0* or *kseg1* can be mapped into *kuseg* by the **mmap()** function. This is covered at more length under “Memory Use in User-Level Drivers” on page 27.

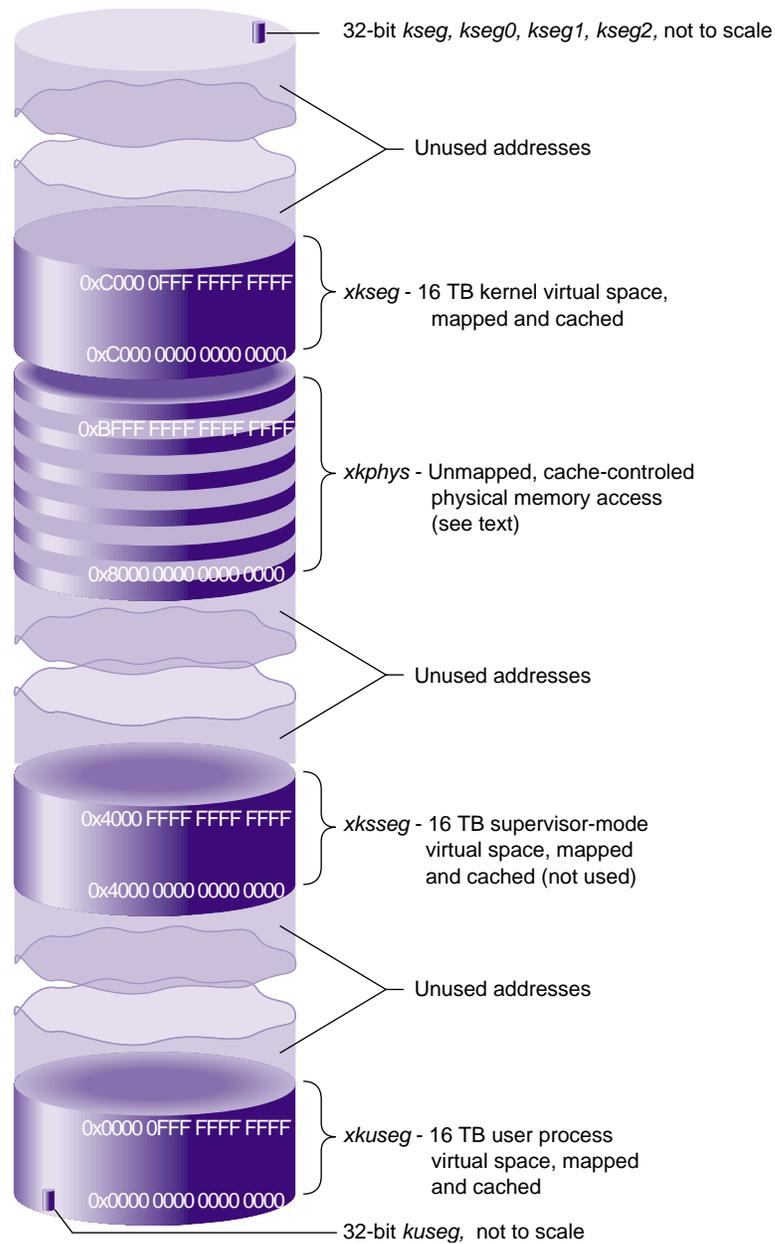
## The 64-Bit Address Space

The 64-bit mode is an upward extension of 32-bit mode. All MIPS processors from the R4000 on support 64-bit mode. However, this mode was not used in Silicon Graphics software until IRIX 6.0 was released.

### Segments of the 64-Bit Address Space

When operating in 64-bit mode, the MIPS architecture uses addresses that are 64-bit unsigned integers from 0x0000 0000 0000 0000 to 0xFFFF FFFF FFFF FFFF. This is an immense span of numbers—if it were drawn to a scale of 1 millimeter per terabyte, the drawing would be 16.8 kilometers long (just over 10 miles).

The MIPS hardware divides the address space into segments based on the most significant bits, and treats each segment differently. The ranges are shown graphically in Figure 1-7. These major segments define only a fraction of the 64-bit space. Most of the possible addresses are undefined and cause an addressing exception (segmentation fault) if used.



**Figure 1-7** Main Parts of the 64-Bit Address Space

As in the 32-bit space, these major segments differ in three characteristics:

- whether access to an address is mapped; that is, passed through the translation lookaside buffer (TLB)
- whether an address can be accessed when the CPU is operating in user mode or in kernel mode.
- whether access to an address is cached; that is, looked up in the primary and secondary caches before it is sent to main memory

### Compatibility of 32-Bit and 64-Bit Spaces

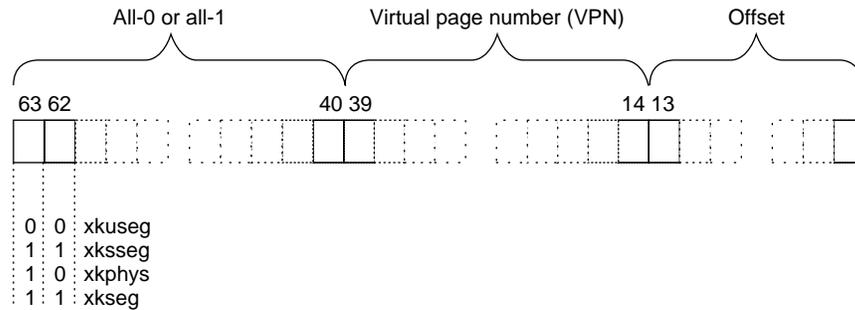
The MIPS-3 instruction set (which is in use when the processor is in 64-bit mode) is designed so that when a 32-bit instruction is used to generate or to load an address, the 32-bit operand is automatically sign-extended to fill the high-order 32 bits.

As a result, any 32-bit address that falls in the user segment *kuseg*, and which must have a sign bit of 0, is extended to a 64-bit integer with 32 high-order 0 bits. This automatically places the 32-bit *kuseg* in the bottom of the 64-bit *xkuseg*, as shown in Figure 1-7.

A 32-bit kernel address, which must have a sign bit of 1, is automatically extended to a 64-bit integer with 32 high-order 1 bits. This places all kernel segments shown in Figure 1-5 at the extreme top of the 64-bit address space. However, these 32-bit kernel spaces are not used by a kernel operating in 64-bit mode.

### Virtual Address Mapping

In the mapped segments, each 64-bit address value is treated as shown in Figure 1-8.



**Figure 1-8** MIPS 64-Bit Virtual Address Format

The two most significant bits select the major segment (compare these to the address boundaries in Figure 1-7). Bits 61:40 must all be 0. (In principle, references to 32-bit kernel segments would have bits 61:40 all 1, but these segments are not used in 64-bit mode.)

The size of a page of virtual memory is a compile-time parameter when the kernel is created. In IRIX 6.2, the page size in a 32-bit kernel is 4 KB and in a 64-bit kernel is 16 KB. (Either size could change in later releases, so always determine it dynamically. In a user-level program, call the `getpagesize()` function (see the `getpagesize(2)` reference page). In a kernel-level driver, use the `ptob()` kernel function (see the `ptob(D3)` reference page) or the constant `NBPP` declared in `sys/immu.h`.)

When the page size is 16 KB, bits 13:0 of the address represent the offset within the page, and bits 39:14 select a VPN from the  $2^{26}$ , or 64 M, pages in the virtual segment..

### User Process Space—`xkuseg`

The first 16 TB of the address space are devoted to user process space. Access to `xkuseg` is always mapped through the TLB. The kernel creates a unique address space for each user process. Of the  $2^{26}$  possible pages in a process's address space, most are typically unassigned, and many are shared pages of program text from dynamic shared objects (DSOs) that are mapped into the address space of every process that needs them.

### Supervisor Mode Space—`xksseg`

The MIPS architecture permits three modes of operation: user, kernel, and supervisor. When operating in kernel or supervisor mode, the 2 TB space beginning at

0x4000 0000 0000 0000 is accessible. IRIX does not employ the supervisor mode, and does not use *xksseg*. If *xksseg* were used, it would be mapped and cached.

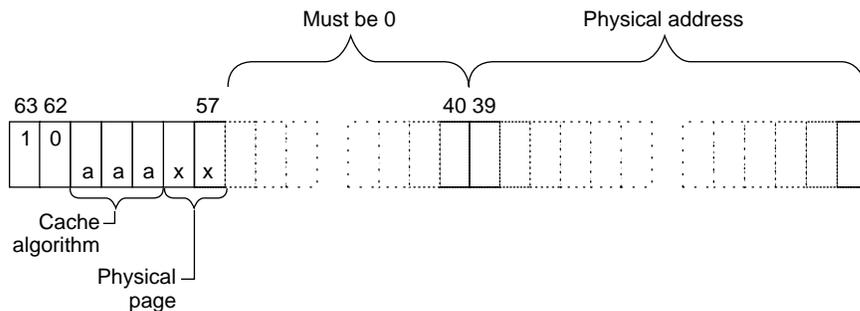
### Kernel Virtual Space—*xkseg*

When bits 63:62 are 11, access is to kernel virtual memory. Only code that is part of the kernel can access this space, a 2 TB segment starting at 0xC000 0000 0000 0000. References to this space are translated through the TLB, and cached. The kernel uses the TLB to map kernel pages in memory as required, possibly in noncontiguous locations. Although pages in kernel space are mapped, they are always associated with real memory. Kernel pages are never paged to secondary storage.

This is the space in which the IRIX kernel allocates such objects as stacks, per-process data that must be accessible on context switches, and user page tables. This area contains automatic variables declared by loadable device drivers. It is the space in which kernel-level device drivers allocate memory. Since kernel space is mapped, addresses in *kseg2* that are apparently contiguous need not be contiguous in physical memory. However, a device driver can allocate space that is both logically and physically contiguous, when that is required (see for example the *kmem\_alloc(D3)* reference page).

### Cache-Controlled Physical Memory—*xkphys*

One-quarter of the 64-bit address space—all addresses with bits 63:62 containing 10—are devoted to special access to the 1 TB physical address space. In 64-bit mode this space replaces the *kseg0* and *kseg1* spaces used in 32-bit mode. Addresses in this space are interpreted as shown in Figure 1-9.



**Figure 1-9** Address Decoding for Physical Memory Access

Bits 39:0 select a physical address in a 1 TB range. As a result, a system operating in 64-bit mode can access a much larger physical address space than the 512 MB space allowed by *kseg0*. This permits more physical memory to be installed, and it gives more freedom in assigning device and bus addresses.

Bits 57:40 must always contain 0. Bits 61:59 select the hardware cache algorithm to be used. The only values defined for these bits are summarized in Table 1-3.

**Table 1-3** Cache Algorithm Selection

Address 61:59	Algorithm	Meaning
010	Uncached	This is the 64-bit equivalent of <i>kseg1</i> in 32-bit mode—uncached access to physical memory.
110	Cacheable coherent exclusive on write	This is the 64-bit equivalent of <i>kseg0</i> in 32-bit mode—cached access to physical memory, coherent access in a multiprocessor.
011	Cacheable non-coherent	Data is cached; on a cache miss the processor issues a non-coherent read (one without regard to other CPUs).
100	Cacheable coherent exclusive	Data is cached; on a read miss the processor issues a coherent read exclusive.
101	Cacheable coherent update on write	Same as 110, but updates memory on a store hit in cache.
111	Uncached Accelerated	Same as 010, but the cache hardware is permitted to defer writes to memory until it has collected a larger block, improving write utilization.

Only the 010 (uncached) and 110 (cached) algorithms are implemented on all systems. The others may or may not be implemented on particular systems.

Bits 58:59 must be 00 unless the cache algorithm is 010 (uncached) or 111(uncached accelerated). Then bits 58:59 can in principle be used to select four other properties to qualify the operation. No present Silicon Graphics computer system supports these properties, so bits 58:59 always contain 00 at this time.

Portions of *xkphys* and *xkseg* can be mapped to user process space by the **mmap()** function. This is covered in more detail under “Memory Use in User-Level Drivers” on page 27.

## Device Driver Use of Memory

Memory use by device drivers is simpler than the details in this chapter suggest. The primary complication for the designer is the use of 64-bit addresses, which may be unfamiliar.

### Allowing for 64-Bit Mode

You must take account of a number of considerations when porting an existing C program to an environment where 64-bit mode is used, or might be used. This can be an issue for all types of drivers, kernel-level and user-level alike. For detailed discussion, see the *MIPSpro 64-Bit Porting and Transition Guide* listed on page xxix.

The most common problems arise because the size of a pointer and of a long int changes between a program compiled with the `-64` option and one compiled `-32`. When you use pointers, longs, or types derived from longs, in structures, the field offsets differ between the two modes.

When all programs in the system are compiled to the same mode, there is no problem. This is the case for a system in which the kernel is compiled to 32-bit mode: only 32-bit user programs are supported. However, a kernel compiled to 64-bit mode executes user programs in 32-bit or 64-bit mode. A structure prepared by a 32-bit program—a structure passed as an argument to `ioctl()`, for example—does not have fields at the offsets expected by a 64-bit kernel device driver. For more on this specific problem, see “Handling 32-Bit and 64-Bit Execution Models” on page 173.

The basic strategy to make your code portable between 32-bit and 64-bit kernels is to be extremely specific when declaring the types of data. You should almost never declare a simple “int” or “char.” Instead, use a data type that is explicit as to the precision and the sign of the variable. The header files `sgidefs.h` and `sys/types.h` define type names that you can use to declare structures that always have the same size. The type `__psint_t`, for example, is an integer the same size as a pointer; you can use it safely as alias for a pointer. Similarly, the type `__uint32_t` is guaranteed to be an unsigned, 32-bit, integer in all cases.

## Memory Use in User-Level Drivers

When you control a device from a user process, your code executes entirely in user process space, and has no direct access to any of the other spaces described in this chapter.

Depending on the device and other considerations, you may use the **mmap()** function to map device registers into the address space of your process (see the `mmap(2)` reference page). When the kernel maps a device address into process space, it does it using the TLB mechanism. From **mmap()** you receive a valid address in process space. This address is mapped through a TLB entry to an address in segment that accesses uncached physical memory. When your program refers to this address, the reference is directed to the system bus and the device.

Portions of kernel virtual memory (*kseg0* or *xkseg*) can be accessed from a user process. Access is based on the use of device special files (see the `mem(7)` reference page). Access is done using two models, a device model and a memory map model.

### Access Using a Device Model

The device special file */dev/mem* represents physical memory. A process that can open this device can use **lseek()** and **read()** to copy physical memory into process virtual memory. If the process can open the device for output, it can use **write()** to patch physical memory.

The device special file */dev/kmem* represents kernel virtual memory (*kseg0* or *xkseg*). It can be opened, read and written similarly to */dev/mem*. Clearly both of these devices should have file permissions that restrict their use even for input.

### Access Using mmap()

The **mmap()** function allows a user process to map an open file into the process address space (see the `mmap(2)` reference page). When the file that is mapped is */dev/mem*, the process can map a specified segment of physical memory. The effect of **mmap()** is to set up a page table entry and TLB entry so that access to a range of virtual addresses in user space is redirected to the mapped physical addresses in cached physical memory (*kseg0* or the equivalent segment of *xkphys*).

The */dev/kmem* device, representing kernel virtual memory, cannot be used with `mmap()`. However, a third device special, */dev/mmem* (note the double “m”), represents access to only those addresses that are configured in the file */var/sysgen/master.d/mem*. As distributed, this file is configured to allow access to the free-running timer device and, in some systems, to graphics hardware.

For an example of mapped access to physical memory, see the example code in the `syssgi(2)` reference page related to the `SGL_QUERY_CYCLECNTR` option. In this operation, the address of the timer (a device register) is mapped into the process’s address space using a TLB entry. When the user process accesses the mapped address, the TLB entry converts it to an address in *kseg1/xkphys*, which then bypasses the cache.

### **Mapped Access Provided by a Device Driver**

A kernel-level device driver can provide mapped access to device registers or to memory allocated in kernel virtual space. An example of such a driver is shown in Part III, “Kernel-Level Drivers.”

### **Memory Use in Kernel-Level Drivers**

When you control a device from a kernel-level driver, your code executes in kernel virtual space. The allocation of memory for program text, local (stack) variables, and static global variables is handled automatically by the kernel. Besides designing data structures so they have a consistent size, you have to consider these special cases:

- dynamic memory allocation for data and for buffers
- transferring data between kernel space and user process space
- getting addresses of device registers to use for PIO

The kernel supplies utility functions to help you deal with each of these issues, all of which are discussed in Chapter 9, “Device Driver/Kernel Interface.”

### **Uncached Memory Access in the Challenge and Onyx Series**

Access to uncached memory is not supported. The Challenge and Onyx systems have coherent caches; cache coherency is maintained by the hardware, even under access from CPUs and concurrent DMA. There is never a need (and no approved way) to access uncached memory in these systems.

### Uncached Memory Access in the IP26 CPU

The IP26 CPU module is used in the Silicon Graphics Power Indigo<sup>2</sup> workstation and the Power Challenge M workstation. Both are desktide workstations using the R8000 processor chip.

Late in the design of these systems, the parity-based memory that had been planned for them was replaced with ECC memory (error-correcting code memory, which can correct for single-bit errors on the fly). ECC memory is also used in large multiprocessor systems from Silicon Graphics, where it has no effect on performance.

Owing to the hardware design of the IP26, ECC memory could be added with no impact on the performance of cached memory access, but uncached memory access can be permitted only when the CPU is placed in a special, “slow” access mode.

In some cases a kernel-level device driver must be sure that stored data has been written into main memory, rather than being held in the cache. There are two ways to ensure this:

- Store the data into cached memory, then use the `dkl_dcachc_wb()` function to force a range of cached addresses to be written to memory. This method works in all systems including the IP26; however, the function call is an expensive one when the amount of data is small.
- Write directly to uncached memory using addresses in `kseg1`. This works in all systems, but in the IP26 (only) it will fail unless the CPU is first put into “slow” mode.

In order to put the CPU into “slow” mode, call the function `ip26_enable_ucmem()`. As soon as the uncached store is complete, return the system to “fast” mode by calling `ip26_return_ucmem()`. (See the `ip26_ucmem(D3)` reference page.) While the CPU is in “slow” mode, several clock cycles are added to every memory access, so do not keep it in “slow” mode any longer than necessary.

These functions can be called in any system. They do nothing unless the CPU is an IP26. Alternatively, you could save the current CPU type using a function like the one shown in Example 2-2 on page 33, and call the functions only when that function returns `INV_IP26BOARD`.



---

## Device Configuration

This chapter discusses how IRIX establishes the inventory of available hardware, and how devices are represented to software.

This information is essential when your work involves attaching a new device or a new class of devices to IRIX. The information is helpful background material when you intend to control a device from a user-level process.

The following primary topics are covered in this chapter.

- “Hardware Inventory” on page 31 describes the hardware inventory table displayed by the *hinv* command and how the inventory is initialized.
- “Device Special Files” on page 34 describes the system of filenames in */dev* and how they are created.
- “Configuration Files” on page 39 summarizes the files used for system generation and kernel configuration.

### Hardware Inventory

In a conventional UNIX system, during bootstrap, each device driver probes the hardware attachments for which it is responsible, and adds information to a hardware inventory table. This is the case with IRIX through IRIX 6.3 for O2. The kernel maintains a hardware inventory table in kernel virtual memory. It is available to users and to programs.

**Note:** In the release of IRIX immediately following IRIX 6.3 for O2, the architecture of the hardware inventory changes radically. However, the functions described in this section continue to be supported for compatibility.

## Using the Hardware Inventory

The hardware inventory is used by users, administrators, and programmers.

### Contents of the Inventory

Using database terminology, the hardware inventory consists of a single table with the following columns:

Class	A code for the class of device; for example, audio, disk, processor, or network.
Type	A code for the type of device within its class; for example, FPU and CPU types within the processor class.
Controller	When applicable, the number of the controller, board, or attachment.
Unit	When applicable, the logical unit or device within a Controller number.
State	A descriptive number, such as the CPU model number.

### Displaying the Inventory with *hinv*

The *hinv* command formats all or selected rows of the inventory table for display (see the *hinv(1)* reference page), translating the numbers to readable form. The user or system administrator can use command options to select a class of entries or certain specific device types by name. The class or type can be qualified with a unit number and a controller number. For example,

```
hinv -c disk -b 1 -u 4
```

displays information about disk 4 on controller 1.

You can use *hinv* to check the result of installing new hardware. The new hardware should show up in the report after the system is booted following installation, provided that the associated device driver was called and was written correctly.

A full inventory report (*hinv -v*) is almost mandatory documentation for a software problem report, either submitted by your user to you, or by you to Silicon Graphics.

## Testing the Inventory In Software

Within a shell script, you can test the output of *hinv* most conveniently in the command exit status. The command sets exit status of 0 when it finds or reports any items. It sets status of 1 when it finds no items. The code in Example 2-1 could be used in a shell script to test the existence of a disk controller.

### Example 2-1 Testing the Hardware Inventory in a Shell Script

```
if hinv -s -c disk -b 1;
then ;
else echo No second disk controller;
fi ;
```

You can access the inventory table in a C program using the functions documented in the `getinvent(3)` reference page. The only access method supported is a sequential scan over the table, viewing all entries. Three functions permit access:

- **setinvent()** initializes or reinitializes the scan to the first row.
- **getinvent()** returns the next table row in sequence.
- **endinvent()** releases storage allocated by **setinvent()**.

These functions use static variables and should only be used by a single process within an address space. Reentrant forms of the same functions, which can safely be used in a multithreaded process, are also available (see `getinvent(3)`). Example 2-2 demonstrates the use of these functions.

The format of one inventory table row is declared as type *inventory\_t* in the *sys/invent.h* header file. This header file also supplies symbolic names for all the class and type numbers that can appear in the table, as well as containing commentary explaining the meanings of some of the numbers.

### Example 2-2 Function Returning Type Code for CPU Module

```
#include <stddef.h> /* for NULL */
#include <invent.h> /* includes sys/invent.h */
int getIPtypeCode()
{
    inv_state_t * pstate = NULL;
    inventory_t * work;
    int ret = 0;
    setinvent_r(&pstate);
    do {
```

```
        work = getinvent_r(pstate);
        if ( (INV_PROCESSOR == work->inv_class)
            && (INV_CPUBOARD == work->inv_type) )
            ret = work->inv_state;
    } while (!ret)
endinvent_r(pstate); /* releases pstate-> */
return ret;
}
```

## Creating an Inventory Entry

Device drivers supplied by Silicon Graphics add information to the hardware inventory table when they are called at their *pxinit()* or *pxedtinit()* entry points. One of these entry points is called by the IRIX kernel during bootstrap. (The small distinction between the two entry points is discussed in “Initialization Entry Points” on page 147.)

The function that adds a row to the inventory table is **add\_to\_inventory()**. Its prototype is declared in the include file *sys/invent.h*. The function takes arguments that are scalar values corresponding to the fields of the *inventory\_t* structure.

**Note:** The only valid inventory types and classes are those declared in *sys/invent.h*. Only those numbers can be decoded and displayed by the *hinvt* command, which prints an error message if it finds an unknown device class, and which prints nothing at all for an unknown device type within a known class. There is no provision for adding new device-class or device-type values for third-party devices.

## Device Special Files

Devices are represented in IRIX as in all conventional UNIX systems, as device special file nodes in the */dev* directory. These special file nodes are, in some cases, created automatically during the bootstrap process, and in some cases created manually by the system administrator. The device special file nodes contain the basic information that lets a user process connect to a device driver to use a device.

**Note:** The discussion in this section is a correct description of IRIX concepts through IRIX 6.3 for O2, and these concepts are referred to again and again in the rest of the book. However, be advised that in the next release of IRIX, these conventions are augmented by an entirely new facility, the hardware graph.

## Device Representation

The IRIX record of a file's existence is sometimes called an *inode*. The device special files consist of inodes only, with no associated data. The fields of the inode are used to encode the following critical information about a device:

Filename	Programs use the name of a device file to open the device using <b>open()</b> .
Permissions, Owner ID, Group ID	The file access permissions, owner ID, and group ID of a device file establish which users can read and which can write to that device.
Block or Character	A device file belongs to one of two classes, block or character, visible as the first letter of an <i>ls -l</i> display.
Major device number	A code for the device driver that controls this device.
Minor device number	A code specifying the unit or position of this device under its controller.

All this information is visible in a display produced by *ls -l*. The major and minor numbers are shown in the column used for file size for regular files. Examine the output of a command such as

```
ls -l /dev/* | more
```

A device special file can be used the same as a regular file in most IRIX commands; for example, a device file can be the target of a symbolic link, the destination of redirected input or output, and so on.

### Block Versus Character

IRIX supports two classes of device. A *block device* such as a disk drive transfers data in fixed size blocks between the device and memory, and usually has some ability to reposition the medium so as to read or write the same data again. The driver for a block device typically has to manage buffering, and it may schedule I/O operations in a different sequence than they are requested.

A *character device* such as a printer accepts or returns data as a stream of bytes, and usually acts as a sink or source of data—the medium cannot be repositioned and read again. The driver for a character device typically transfers data as soon as it is requested and completes one operation before accepting another request. Character devices are also called *raw* devices, because their input is not buffered.

### Major Device Number

The *major device number* recorded in the device special inode selects the device driver to service this device. When a device is opened, IRIX selects the driver to handle the device based on the major device number. Each device driver supports one or more specific major numbers. There are two unrelated ranges of major numbers, one for character device drivers and one for block device drivers.

The possible major numbers are declared and given names in the file *sys/major.h*. When you create a new kernel-level device driver you must choose a major number for it—a number not used by any other driver. Numbers 60-79 are not used by Silicon Graphics. (See “Selecting a Major Number” on page 228.)

In IRIX releases through 5.2 (and 6.0.x, which is based on 5.2), major numbers were limited to the range 0 through 254. Beginning with releases 5.3 and 6.1, the IRIX inode structure permits major numbers to have up to 14 bits of precision. However, major numbers are currently restricted to at most 9 bits to limit the size of kernel tables that are indexed by the major number.

In order to use this limit symbolically, use the name `L_MAXMAJ` defined in *sys/sysmacros.h*. When you declare a variable for a major device number in a program, use type *major\_t* declared in *sys/types.h*.

Normally a device driver services only one major number. However, it is possible to designate the same device driver to service more than one major number. In this case, the driver may need to discover the major number at execution time. The **getemajor()** function returns the number in use for a given request (see the `getemajor(D3)` reference page).

### Minor Device Number

The *minor device number* is passed to the device driver as an argument when the driver is called. (The major and minor numbers are passed together in a long integer called a *dev\_t*.) The minor device number is interpreted only by the device driver, so it can be a simple logical unit number, or it can contain multiple, encoded bit fields. For example:

- The IRIX tape device driver uses the minor device number to encode the options for rewind or no-rewind, byte-swap or nonswap, and fixed or variable blocking, along with the logical unit number.
- The IRIX disk device drivers encode the disk partition number into the minor device number along with a disk unit number. Both disk and tape devices encode the SCSI adapter number in the minor number.
- The IRIX generic SCSI driver encodes the adapter (bus) number, target (control unit) number, and logical unit number into the minor number (see “Generic SCSI Device Special Files” on page 82).

The IRIX inode structure permits minor numbers to have up to 18 bits of precision. In order to use this limit symbolically, use the name `L_MAXMIN` defined in *sys/sysmacros.h*. When you declare a variable for a minor device number in a program, use type *minor\_t* declared in *sys/types.h*.

With STREAMS drivers, the minor device number can be chosen arbitrarily during a CLONE open—see “Support for CLONE Drivers” on page 510.

### Defining Device Names

The device special files related to Silicon Graphics device drivers are created by execution of the script */dev/MAKEDEV*. Additional device special files can be created with administrator commands.

### IRIX Conventional Device Names

The device drivers distributed with IRIX depend on certain conventions for device names. These conventions are spelled out in the following reference pages: *intro(7)*, *dks(7)*, *dsreq(7)*, and *tps(7)*. For example, the components of a disk device name in */dev/dsk* include

- dks***c*            Constant prefix “dks” followed by bus adapter number *c*.
- d***u*                Constant letter “d” followed by disk SCSI ID number *u*.
- l***n*                Optionally, letter “l” (ell) and logical unit number *n* (used only when disk *u* controls multiple drives).
- sp** or **vh** or **vol**    Constant letter “s” and partition number *p*, or else “vh” for volume header, or “vol” for (entire) volume.

Programs throughout the system rely on the conventions for these device names. In addition, by convention the associated major and minor numbers agree with the names. For example, the logical unit and partition numbers that appear in a disk name are also encoded into the minor number.

### The Script MAKEDEV

The conventions for all the IRIX device special names are written into the script */dev/MAKEDEV*. This is a make file, but unlike most make files, it is not used to compile executable programs. It contains the logic to prepare device special names and their associated major and minor numbers and file permissions.

The MAKEDEV script is executed during IRIX startup from a script in */etc/rc2.d*. It is executed after all device drivers have been initialized, so it can use the output of the *hinvt* command to construct device names to suit the actual configuration.

The system administrator can invoke MAKEDEV to construct device special files. Administrator use of MAKEDEV is described in *IRIX Administration: System Configuration and Operation*.

### Making Device Files

You or a system administrator can create device special files explicitly using the commands *mknod* or *install*. Either command can be used in a make file such as you might create as part of the installation script for a product.

For details of these commands, see the `install(1)` and `mknod(1M)` reference pages, and *IRIX Administration: System Configuration and Operation*. The following is a hypothetical example of `install`:

```
# install -m 644 -u root -g sys -root /dev -chr 62,0
```

The `-chr` option specifies a character device, and `62,0` are the major and minor device numbers, respectively.

**Tip:** The `mknod` command is portable, being used in most UNIX systems. The `install` command is unique to IRIX, and has a number of features and uses beyond those of `mknod`. Examples of both can be found by reading `/dev/MAKEDEV`.

### Multiple Names for One Device

It is possible to point to the same device with more than one device special filename. This is done in the distributed IRIX system for several reasons:

- To supply default names for devices with specific names. For example, the default device `/dev/tapens` is a link to the first device file in `/dev/rmt/*`.
- To pass different parameters to the device driver. For example, the same tape device appears multiple times in `/dev/rmt/tps*`, with different combinations of `nr` (norewind), `ns` (nonswapped), and `v` (variable block) suffixes. The minor number for each name encodes these options for the same unit number.
- To supply both block and character drivers for the same device. For example, each disk device appears in `/dev/dsk/*` as a block device, and again in `/dev/rdsk/*` as a character device.

## Configuration Files

IRIX uses a number of configuration files to supplement its knowledge of devices and device drivers. This is a summary of the files. The use of each file for device driver purposes is described in more detail in other chapters. (The uses of these files for other system administration tasks is covered in *IRIX Administration: System Configuration and Operation*.)

Most configuration files used by the IRIX kernel are located in the directory */var/sysgen*. Files used by the X11 display system are generally in */usr/lib/X11*. With regard to device drivers, the important files are:

<i>/var/sysgen/master.d.*</i>	Descriptions of the attributes of kernel modules
<i>/var/sysgen/boot/*</i>	Kernel object modules
<i>/var/sysgen/system/*.*sm</i>	Device configuration information
<i>/var/sysgen/mtune/*</i>	Values and limits of tunable parameters
<i>/var/sysgen/stune</i>	New values for tunable parameters
<i>/usr/lib/X11/input/config/*</i>	Initialization commands for Xdm input modules

### Master Configuration Database

Every configurable module of the kernel (this includes kernel-level device drivers and some other service modules) is represented by a single file in the directory */var/sysgen/master.d*.

A file in *master.d* describes the attributes of a module of the kernel which is to be loaded at boot time. The general syntax of the file is documented in detail in the *master(4)* reference page. Only a subset of the syntax is used to describe a device driver module. In general, the *master.d* file specifies device driver attributes such as:

- the driver's *prefix*, a name that qualifies all its entry points
- whether it is a block, character, or STREAMS driver
- the major number serviced by the driver
- whether the driver can be loaded dynamically as needed
- whether the driver is multiprocessor-aware
- which of the possible driver entry points the driver supplies

For each module described in a *master.d* file there should be a corresponding object module in */var/sysgen/boot*. The creation of device driver modules and the syntax of *master.d* files is covered in detail in Chapter 10, "Building and Installing a Driver."

## System Configuration Files

The files */var/sysgen/system/\*.sm* direct the *lboot* command in loading the modules of the kernel at boot time. Although there are normally several files with the names of subsystems, all the files are treated as one single file. The contents of the files direct *lboot* in loading components that are described by files in */var/sysgen/master.d*, and in probing for devices to see if they exist.

The exact syntax of these files is documented in the *system(4)* reference page. The use of the VECTOR lines to probe for hardware is covered in this book in the context of each type of attachment..

## System Tuning Parameters

The IRIX kernel supports a variety of tunable parameters, some of which can be interrogated by device drivers. The current values of the parameters are recorded in files in */var/sysgen/mtune/\** (one file per major subsystem).

You or the system administrator can view the current settings using the *systune* command (see the *systune(1M)* reference page). The system administrator can use *systune* to request changes in parameters. Some changes take effect at once; others are recorded in a modified kernel that is loaded the next time the system boots.

To retrieve certain tuning parameters from within a kernel-level device driver, include the header file *sys/var.h*.

The use of *systune* and its related files is covered in *IRIX Administration: System Configuration and Operation*.

## X Display Manager Configuration

Most files related to the configuration of the X Display Manager *Xdm* are held in */var/X11*. These files are documented in reference pages such as *xm(1)* and in the programming manuals related to the X Windows System™.

One set of files, in */usr/lib/X11/input/config*, controls the initialization of nonstandard input devices. These devices use STREAMS modules, and their configuration is covered in Chapter 16, "STREAMS Drivers."



---

## Device Control Software

IRIX provides for two general methods of controlling devices, at the user level and at the kernel level. This chapter describes the architecture of these two software levels and points out the different abilities of each. This is important background material for understanding all types of device control. The chapter covers the following main topics:

- “User-Level Device Control” on this page summarizes five methods of device control for user-initiated processes.
- “Kernel-Level Device Control” on page 47 sets the concepts needed to understand kernel-level drivers.

### User-Level Device Control

In IRIX terminology, a *user-level* process is one that is initiated by a user (possibly the superuser). A user-level process runs in an address space of its own, with no access to the address space of other processes or to the kernel’s address space, except through explicit memory-sharing agreements.

In particular, a user-level process has no access to physical memory (which includes access to device registers) unless the kernel allows the process to share part of the kernel’s address space. (For more on physical memory, see Chapter 1, “Physical and Virtual Memory.”)

There are several ways in which a user-level process can control devices, which are summarized in the following topics:

- “EISA Mapping Support” on page 44 summarizes PIO access to the EISA bus.
- “VME Mapping Support” on page 44 summarizes PIO access to the VME bus.
- “PCI Mapping Support” on page 45 summarizes PIO access to the PCI bus.
- “User-Level DMA From the VME Bus” on page 45 summarizes DMA I/O managed from a user-level process.

- “User-Level Control of SCSI Devices” on page 45 summarizes DMA and command access to the SCSI bus.
- “Managing External Interrupts” on page 46 summarizes access to the external interrupt ports on Challenge and Onyx systems.
- “User-Level Interrupt Management” on page 46 summarizes the handling of some interrupts in a user-level process.

### **EISA Mapping Support**

In systems that support the EISA bus (Indigo<sup>2</sup> Maximum Impact and Indigo<sup>2</sup>, Challenge M, and their Power versions), IRIX contains a kernel-level device driver that supports memory-mapping EISA bus addresses into the address space of a user process (see “Overview of Memory Mapping” on page 53).

You can write a program that maps a portion of the EISA bus address space into the program address space. Then you can load and store from device registers directly.

For more details of PIO to the EISA bus, see Chapter 4, “User-Level Access to Devices.”

### **VME Mapping Support**

In systems that support the VME bus (Onyx, Challenge DM, Challenge L, Challenge XL, and their Power versions), IRIX contains a kernel-level device driver that supports mapping of VME bus addresses into the address space of a user process (see “Overview of Memory Mapping” on page 53).

You can write a program that maps a portion of the VME bus address space into the program address space. Then you can load and store from device registers directly.

For more details of PIO to the VME bus, see Chapter 4, “User-Level Access to Devices.”

## PCI Mapping Support

In systems that support the PCI bus (O2 and related workstations), a kernel-level device driver for a PCI device can provide support for the **mmap()** system function (see the **mmap(2)** reference page), and in this way can allow a user-level process to map some part of the I/O or memory space defined by a particular PCI device into the address space of the process (see “Overview of Memory Mapping” on page 53).

This must be done by a specific device driver; there can be no general-purpose bus mapping driver as there is for the VME bus. (This is because PCI devices are assigned bus address space dynamically, and there is no interface by which a general device driver could learn the bus addresses assigned.)

When a specific device driver supports PIO mapping, your program can load and store values directly to and from locations defined by the mapped device. For more details of PIO to the PCI bus, see Chapter 4, “User-Level Access to Devices.”

## User-Level DMA From the VME Bus

The Challenge L, Challenge XL, and Onyx systems and their Power versions contain a DMA engine that manages DMA transfers from VME devices, including VME slave devices that normally cannot do DMA.

The DMA engine in these systems can be programmed directly from code in a user-level process. Software support for this facility is contained in the *udmalib* package.

For more details of user DMA, see Chapter 4, “User-Level Access to Devices” and the *udmalib(3)* reference page.

## User-Level Control of SCSI Devices

IRIX contains a special kernel-level device driver whose purpose is to give user-level processes the ability to issue commands and read and write data on the SCSI bus. By using **ioctl()** calls to this driver, a user-level process can interrogate and program devices, and can initiate DMA transfers between memory buffers and devices.

The low-level programming used with the *dsreq* device driver is eased by the use of a library of utility functions documented in the *dslib(3)* reference page.

For more details on user-level SCSI access, see Chapter 5, “User-Level Access to SCSI Devices.”

## Managing External Interrupts

The Challenge L, Challenge XL, and Onyx systems and their Power versions have four external-interrupt output jacks and four external-interrupt input jacks on their back panels. In these systems, the device special file `/dev/ei` represents a device driver that manages access to these external interrupt ports.

Using `ioctl()` calls to this device (see “Overview of Device Control” on page 50), your program can

- enable and disable the detection of incoming external interrupts
- set the strobe length of outgoing signals
- strobe, or set a fixed level, on any of the four output ports

In addition, library calls are provided that allow very low-latency detection of an incoming signal.

For more information on external interrupt management, see Chapter 6, “Control of External Interrupts” and the `ei(7)` reference page.

## User-Level Interrupt Management

A facility introduced in IRIX 6.2 allows you to receive and handle certain device interrupts in a user-level program you write.

Your program calls a library function to register the interrupt-handling function. When the device generates an interrupt, the kernel branches directly into your handler. Because this handler runs as a subroutine of the kernel, it can use only a very limited set of system and library functions. However, it can refer to variables in the process address space, and it can wake up a process that is blocked, waiting for the interrupt to occur.

Combined with PIO, user-level interrupts allow you to test most of the logic of a device driver for a new device in user-level code.

In IRIX 6.3 for O2, support for user-level interrupts is limited to VME devices and to external interrupts in the Challenge L, Challenge XL, and Onyx systems and their POWER versions. In a future release, user-level interrupts will be supported for PCI devices as well.

For more details on user-level interrupts, see Chapter 7, “User-Level Interrupts” and the `uli(3)` reference page.

### **Memory-Mapped Access to Serial Ports**

The Audio/Serial Option (ASO) board for the Challenge and Onyx series provides six high-performance serial ports, each of which can be set to run at speeds as high as 115,200 bits per second. The features and administration of the Audio/Serial Option board are described in the *Audio/Serial Option User's Guide* (document 007-2645-001).

The serial ports of the ASO board can be accessed in the usual way, by opening a file to a device in the `/dev/tty*` group of names. However, for the minimum of latency and overhead, a program can open a device in the `/dev/aso_mmap` directory. These device files are managed by a device driver that permits the input and output ring buffers for the port to be mapped directly into the user process address space. The user-level program can spin on the input ring buffer pointers and detect the arrival of a byte of data in microseconds after the device driver stores it.

The details of the memory-mapping driver for ASO ports are spelled out in the `asoserns(7)` reference page (available only when the ASO feature has been installed).

## **Kernel-Level Device Control**

IRIX supports the conventional UNIX architecture in which a user process uses a kernel service to request a data transfer, and the kernel calls on a device driver to perform the transfer.

## Kinds of Kernel-Level Drivers

There are three distinct kinds of kernel-level drivers:

- A *character device driver* transfers data as a stream of bytes of arbitrary length. A character device driver is invoked when a user process issuing a system function call such as `read()` or `ioctl()`.
- A *block device driver* transfers data in blocks of fixed size. Normally a block driver is not called directly to support a user process. User reads and writes are directed to files, and the filesystem code calls the block driver to read or write whole disk blocks. Block drivers are also called for paging operations.
- A STREAMS driver is not a device driver, but rather can be dynamically installed to operate on the flow of data to and from any character device driver.

Overviews of the operation of STREAMS drivers are found in Chapter 16, “STREAMS Drivers.” The rest of this discussion is on character and block device drivers.

## Typical Driver Operations

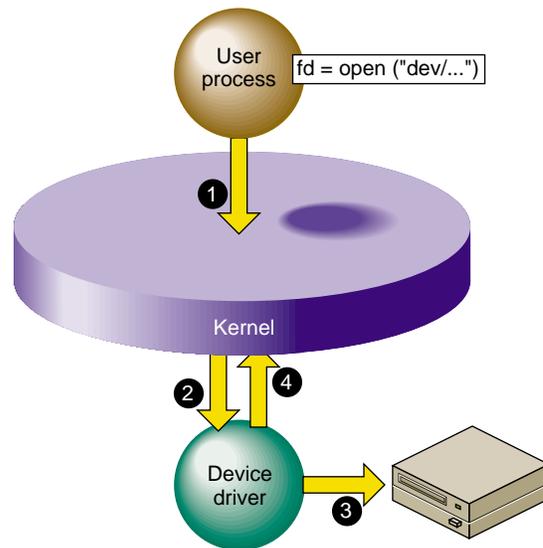
There are five different kinds of operations that a device driver can support:

- The open interaction is supported by all drivers; it initializes the connection between a process and a device.
- The control operation is supported by character drivers; it allows the user process to modify the connection to the device or to control the device.
- A character driver transfers data directly between the device and a buffer in the user process address space. This is typically done with programmed I/O (PIO) to transfer small quantities of data synchronously.
- Memory mapping enables the user process to perform PIO for itself.
- A block driver transfers one or more fixed-size blocks of data between the device and a buffer owned by a filesystem or the memory paging system. This is typically done with Direct memory access (DMA) to transfer larger quantities of data asynchronously under device control.

The following topics present a conceptual overview of the relationship between the user process, the kernel, and the kernel-level device driver. The software architecture that supports these interactions is documented in detail in Part III, “Kernel-Level Drivers,” especially Chapter 8, “Structure of a Kernel-Level Driver.”

### Overview of Device Open

Before a user process can use a kernel-controlled device, the process must open the device as a file. A high-level overview of this process, as it applies to a character device driver, is shown in Figure 3-1.



**Figure 3-1** Overview of Device Open

The steps illustrated in Figure 3-1 are:

1. The user process calls the **open()** kernel function, passing the name of a device special file (see “Device Special Files” on page 34 and the `open(2)` reference page).
2. The kernel notes the device major and minor numbers from the inode of the device special file (see “Device Representation” on page 35). The kernel uses the major device number to select the device driver, and calls the driver’s open entry point, passing the minor number and other data.
3. The device driver verifies that the device is operable, and prepares whatever is needed to operate it.
4. The device driver returns a return code to the kernel, which returns either an error code or a *file descriptor* to the process.

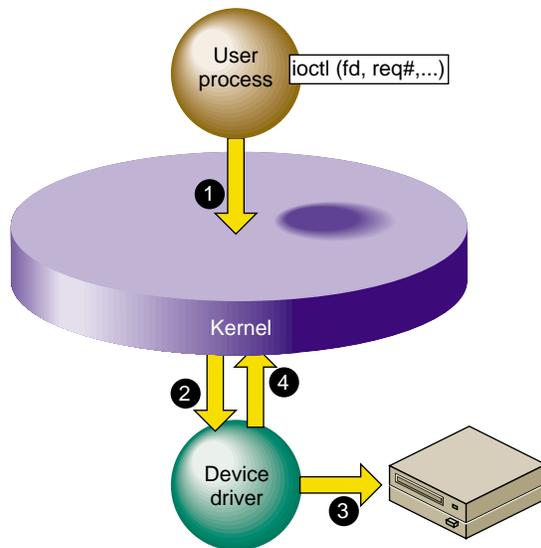
It is up to the device driver whether the device can be used by only one process at a time, or by more than one process. If the device can support only one user, and is already in use, the driver returns the EBUSY error code.

The **open()** interaction on a block device is similar, except that the operation is initiated from the filesystem code responding to a **mount()** request, rather than coming from a user process **open()** request (see the mount(1) reference page).

There is also a **close()** interaction so a process can terminate its connection to a device.

### Overview of Device Control

After the user process has successfully opened a character device, it can request control operations. Figure 3-2 shows an overview of this operation.



**Figure 3-2** Overview of Device Control

The steps illustrated in Figure 3-2 are:

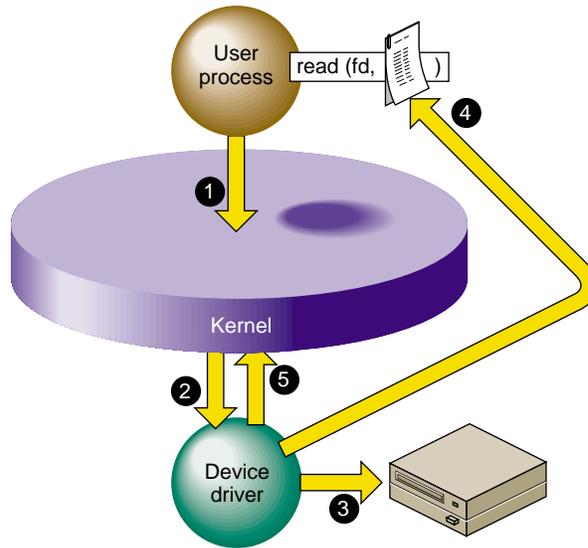
1. The user process calls the **ioctl()** kernel function, passing the file descriptor from **open** and one or more other parameters (see the **ioctl(2)** reference page).
2. The kernel uses the major device number to select the device driver, and calls the device driver, passing the minor device number, the request number, and an optional third parameter from **ioctl()**.
3. The device driver interprets the request number and other parameter, notes changes in its own data structures, and possibly issues commands to the device.
4. The device driver returns an exit code to the kernel, and the kernel (then or later) redispaches the user process.

Block device drivers are not asked to provide a control interaction. The user process is not allowed to issue **ioctl()** for a block device.

The interpretation of **ioctl** request codes and parameters is entirely up to the device driver. For examples of the range of **ioctl** functions, you might review some reference pages in volume 7, for example, **termio(7)**, **ei(7)**, and **arp(7P)**.

### **Overview of Character Device I/O**

Figure 3-3 shows a high-level overview of data transfer for a character device driver that uses programmed I/O.



**Figure 3-3** Overview of Programmed Kernel I/O

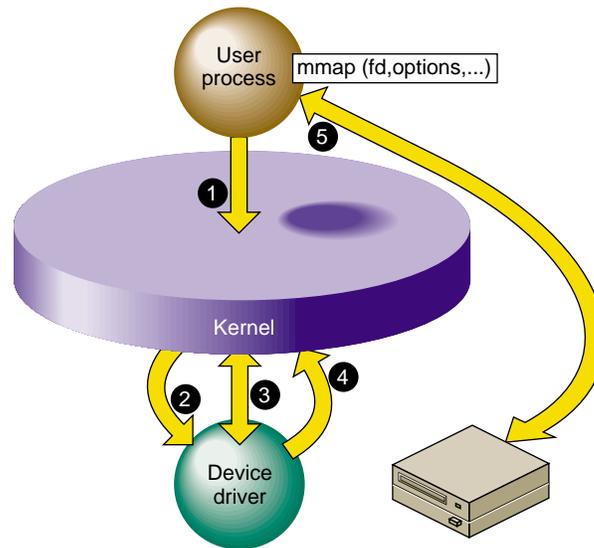
The steps illustrated in Figure 3-3 are:

1. The user process invokes the **read()** kernel function for the file descriptor returned by **open()** (see the `read(2)` and `write(2)` reference pages).
2. The kernel uses the major device number to select the device driver, and calls the device driver, passing the minor device number and other information.
3. The device driver directs the device to operate by storing into its registers in physical memory.
4. The device driver retrieves data from the device registers and uses a kernel function to store the data into the buffer in the address space of the user process.
5. The device driver returns to the kernel, which (then or later) dispatches the user process.

The operation of **write()** is similar. A kernel-level driver that uses programmed I/O is conceptually simple since it is basically a subroutine of the kernel.

### Overview of Memory Mapping

It is possible to allow the user process to perform I/O directly, by mapping the physical addresses of device registers into the address space of the user process. Figure 3-4 shows a high-level overview of this interaction.



**Figure 3-4** Overview of Memory Mapping

The steps illustrated in Figure 3-4 are:

1. The user process calls the **mmap()** kernel function, passing the file descriptor from `open` and various other parameters (see the `mmap(2)` reference page).
2. The kernel uses the major device number to select the device driver, and calls the device driver, passing the minor device number and certain other parameters from **mmap()**.
3. The device driver validates the request and uses a kernel function to map the necessary range of physical addresses into the address space of the user process.
4. The device driver returns an exit code to the kernel, and the kernel (then or later) redispaches the user process.
5. The user process accesses data in device registers by accessing the virtual address returned to it from the **mmap()** call.

Memory mapping can be supported only by a character device driver. When a user process applies `mmap()` to an ordinary disk file, the filesystem maps the file into memory. The filesystem may call a block driver to transfer pages of the file in and out of memory, but to the driver this is no different from any other read or write call.

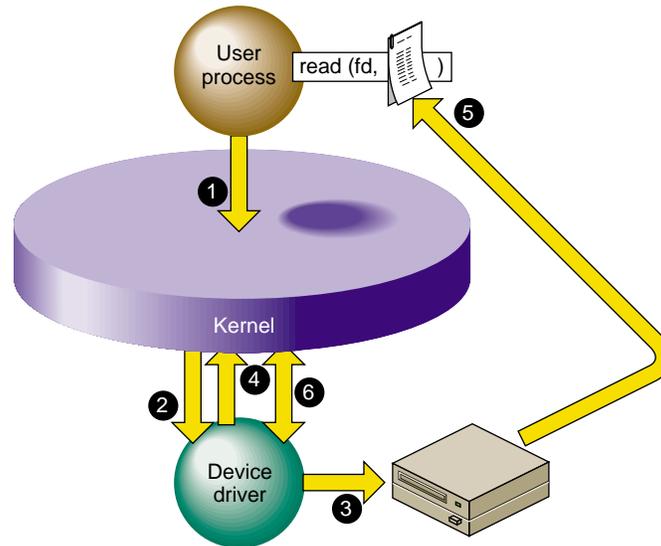
Memory mapping by a character device driver has the purpose of making device registers directly accessible to the process as memory addresses. A memory-mapping character device driver is very simple; it needs to support only `open()`, `mmap()`, and `close()` interactions. Data throughput can be higher when PIO is performed in the user process, since the overhead of the `read()` and `write()` system calls is avoided.

Silicon Graphics device drivers for the VME and EISA buses support memory mapping. This enables user-level processes to perform PIO to devices on these buses, as described under “EISA Mapping Support” on page 44 and “VME Mapping Support” on page 44. Character drivers for the PCI bus are allowed to support memory mapping.

It is possible to write a kernel-level driver that only maps memory, and controls no device at all. Such drivers are called *pseudo-device* drivers. For examples of pseudo-device drivers, see the `prf(7)` and `imon(7)` reference pages.

### Overview of Block I/O

Block devices and block device drivers normally use DMA (see “Direct Memory Access” on page 11). With DMA, the driver can avoid the time-consuming process of transferring data between memory and device registers. Figure 3-5 shows a high-level overview of a DMA transfer.



**Figure 3-5** Overview of DMA I/O

The steps illustrated in Figure 3-5 are:

1. The user process invokes the `read()` kernel function for a normal file descriptor (not necessarily a device special file). The filesystem (not shown) asks for a block of data.
2. The kernel uses the major device number to select the device driver, and calls the device driver, passing the minor device number and other information.
3. The device driver uses kernel functions to create a DMA map that describes the buffer in physical memory; then programs the device with target addresses by storing into its registers.
4. The device driver returns to the kernel after telling it to put to sleep the user process that called the driver.
5. The device itself stores the data to the physical memory locations that represent the buffer in the user process address space. While this is going on, the kernel may dispatch other processes.
6. When the device presents a hardware interrupt, the kernel invokes the device driver. The driver notifies the kernel that the user process can now resume execution. It resumes in the filesystem code, which moves the requested data into the user process buffer.

DMA is fundamentally asynchronous. There is no necessary timing relation between the operation of the device performing its operation and the operation of the various user processes. A DMA device driver has a more complex structure because it must deal with such issues as

- making a DMA map and programming a device to store into a buffer in physical memory
- blocking a user process, and waking it up when the operation is complete
- handling interrupts from the device
- the possibility that requests from other processes can occur while the device is operating
- the possibility that a device interrupt can occur while the driver is handling a request

The reward for the extra complexity of DMA is the possibility of much higher performance. The device can store or read data from memory at its maximum rated speed, while other processes can execute in parallel.

A DMA driver must be able to cope with the possibility that it can receive several requests from different processes while the device is busy handling one operation. This implies that the driver must implement some method of queuing requests until they can be serviced in turn.

The mapping between physical memory and process address space can be complicated. For example, the buffer can span multiple pages, and the pages need not be in contiguous locations in physical memory. If the device does not support *scatter/gather* operations, the device driver has to program a separate DMA operation for each page or part of a page—or else has to obtain a contiguous buffer in the kernel address space, do the I/O from that buffer, and copy the data from that buffer to the process buffer. When the device supports *scatter/gather*, it can be programmed with the starting addresses and lengths of each page in the buffer, and read and write into them in turn before presenting a single interrupt.

### Upper and Lower Halves

When a device can produce hardware interrupts, its kernel-level device driver has two distinct logical parts, called the “upper half” and the “lower half” (although the upper “half” is usually much more than half the code).

### Driver Upper Half

The upper half of a driver comprises all the parts that are invoked as a result of user process calls: the driver entry points that execute in response to **open()**, **close()**, **ioctl()**, **mmap()**, **read()** and **write()**.

These parts of the driver are called on behalf of a specific process. This is referred to as “having user context,” which means that they are executed under the identity of a specific process.

As a result, code in the upper half of the driver is allowed to request kernel services that can be delayed, or “sleep.” For example, code in the upper half of a driver can call **kmem\_alloc()** to request memory in kernel space, and can specify that if memory is not available, the driver can sleep until memory is available. Also, code in the upper half can wait on a semaphore until some event occurs, or can seize a lock knowing that it may have to sleep until the lock is released.

In each case, the entire kernel does not “sleep.” The driver upper half sleeps under the identity of the user process; but the kernel dispatches other processes to run. When the blocking condition is removed—when memory is available, the semaphore is posted, or the lock is released—the driver is scheduled for execution and resumes.

### Driver Lower Half

The lower half of a driver comprises the code that is called to respond to a hardware interrupt. An interrupt can occur at almost any time, including large parts of the time when the kernel is executing other services, including driver upper halves, and even driver lower halves for devices with lower-priority interrupts.

The kernel is not in a known state when executing a driver lower half, and there is no process context. Several things follow from this fact:

- It is very important that the interrupt be handled in the absolute minimum of time, since it may be delaying a kernel service or even the handling of a lower-priority interrupt.
- The lower-half code may not use any kernel service that can sleep (because there is no dispatchable process to be blocked and dispatched again later). Every authorized kernel service is documented as to whether it can sleep or not.

### Relationship Between Halves

Each half has its proper kind of work. In general terms, the upper half performs all validation and preparation, including allocating and deallocating memory and copying data between address spaces. It initiates the first device operation of a series and queues other operations. Then it waits on a semaphore.

The lower half verifies the correct completion of an operation. If another operation is queued, it initiates that operation. Then it posts the semaphore to awaken the upper half, and exits.

### Layered Drivers

IRIX allows for “layered” device drivers, in which one driver operates the actual hardware and the driver at the higher layer presents the programming interface. This approach is implemented for SCSI devices: actual management of the SCSI bus is delegated to a set of Host Adapter drivers. Drivers for particular kinds of SCSI devices call the Host Adapter driver through an indirect table to execute SCSI commands. SCSI drivers and Host Adapter drivers are discussed in detail in Chapter 13, “SCSI Device Drivers.”

### Combined Block and Character Drivers

A block device driver is called indirectly, from the filesystem, and it is not allowed to support the `ioctl()` entry point. In some cases, block devices can also be thought of as character devices. For example, a block device might return a string of diagnostic information, or it might be sensitive to dynamic control settings.

It is possible to support *both* block and character access to a device: block access to support filesystem operations, and character access in order to allow a user process (typically one started by a system administrator) to read, write, or control the device directly.

For example, the Silicon Graphics disk device drivers support both block and character access to disk devices. This is why you can find every disk device represented as a block device in the `/dev/dsk` directory and again as a character device in `/dev/rdisk` (“r” for “raw,” meaning character devices).

## Drivers for Multiprocessors

Many Silicon Graphics computers have multiple CPUs that execute concurrently. The CPUs share access to the single main memory, including a single copy of the kernel address space. In principle, all CPUs can execute in the kernel code simultaneously. In principle, the upper half of a device driver could be entered simultaneously by as many different processes as there are CPUs in the system (up to 36 in a Challenge or Onyx system).

A device driver written for a uniprocessor system cannot tolerate concurrent execution by multiple CPUs. For example, a uniprocessor driver has scalar variables whose values would be destroyed if two or more processes updated them concurrently.

In order to make uniprocessor drivers work in multiprocessors, IRIX by default uses only CPU 0 to execute calls to upper-half code of character and STREAMS drivers. This ensures that at most one process executes in any upper half at one time. (Network and block device drivers do not receive this service.)

It is not difficult to design a kernel-level driver to execute safely in any CPU of a multiprocessor. Each critical data object must be protected by a lock or semaphore, and particular techniques must be used to coordinate between the upper and lower halves. These techniques are discussed in "Planning for Multiprocessor Use" on page 174.

When you have made a driver multiprocessor-safe, you compile it with a particular flag value that IRIX recognizes. From then on, the driver upper half is executed on any CPU of a multiprocessor. This can improve performance, since processes that use the driver are not required to wait for CPU 0 to be available.

## Loadable Drivers

Some drivers are needed whenever the system is running, but others are needed only occasionally. IRIX allows you to create a kernel-level device driver or STREAMS driver that is not loaded at boot time, but only later when it is needed.

A loadable driver has the same purposes as a nonloadable one, and uses the same interfaces to do its work. A loadable driver can be configured for automatic loading when its device is opened. Alternatively it can be loaded on command using the *ml* program (see the *ml(1)* and *mload(4)* reference pages).

A loadable driver remains in memory until its device is no longer in use, or until the administrator uses *ml* to unload it. A loadable driver remains in memory indefinitely, and cannot be unloaded, unless it provides a *pxunload()* entry point (see “Entry Point *unload()*” on page 170).

There are some small differences in the way a loadable driver is compiled and configured (see “Configuring a Loadable Driver” on page 239).

One operational difference is that a loadable driver is not available in the miniroot, the standalone system administration environment used for emergency maintenance. If a driver might be required in the miniroot, it can be made nonloadable, or it can be configured for “autoregistration” (see “Registration” on page 241).

## PART TWO

# Device Control From Process Space

### **Chapter 4:** User-Level Access to Devices

How a user-level process can access and control devices on various buses.

### **Chapter 5:** User-Level Access to SCSI Devices

How a user-level process can execute commands and transfer data to a SCSI device.

### **Chapter 6:** Control of External Interrupts

How a user-level process creates or responds to external interrupt signals in the Challenge and Power Challenge systems.

### **Chapter 7:** User-Level Interrupts

How a user-level process can trap and respond to device interrupts with low latency and the fewest context switches.



---

## User-Level Access to Devices

Programmed I/O (PIO) refers to loading and storing data between program variables and device registers. This is done by setting up a memory mapping of a device into the process address space, so that the program can treat device registers as if they were volatile memory locations. This chapter discusses the methods of setting up this mapping, and the performance that can be obtained. The main topics are as follows:

- “VME Programmed I/O” on page 64 discusses PIO mapping of VME devices.
- “EISA Programmed I/O” on page 68 discusses PIO mapping of EISA devices.
- “VME User-Level DMA” on page 72 discusses the use of the DMA engine in a Challenge or Onyx system.
- “PCI Programmed I/O” on page 78 discusses PIO mapping of PCI devices.

**Note:** Of these topics, only “PCI Programmed I/O” on page 78 is applicable to O2 workstations. The other topics all require the use of different hardware—Challenge or Onyx systems for VME, and Indigo<sup>2</sup> systems for EISA. These topics are included here in case you are using an O2 workstation to develop applications for one of those systems which is running IRIX 5.3 or 6.2.

## VME Programmed I/O

The VME bus is available on Silicon Graphics systems such as the Challenge and Onyx. If you are writing or maintaining a VME-based user-level application to execute on a system with a VME bus, the information in this section is of interest.

### Mapping a VME Device Into Process Address Space

As discussed in “Physical Device Addresses” on page 4, an I/O device is represented as an address, or range of addresses, in the address space of its bus. A kernel-level device driver has the ability to set up a mapping between the bus address of a device register and a location in the address space of a user-level process. When this has been done, the device register appears to be a variable in memory. The program can assign values to it, or refer to it in expressions.

### Learning VME Device Addresses

In order to map a VME device for PIO, you must know the following points:

- the VME bus number on which the device resides  
Challenge and Onyx systems support as many as five VME buses. The first is number 0. Use the *hinv* command on the Challenge or Onyx system to display the numbers of others.
- the VME address space in which the device resides  
This will be either A16, A24, or A32—the A64 space is not supported for PIO.
- the VME address space modifier the device uses—either supervisory (s) or nonprivileged (n)
- the VME bus addresses associated with the device  
This must be a sequential range of VME bus addresses that spans all the device registers you need to map.

This information is normally supplied by the manufacturer of a third-party VME device. You can find these values for Silicon Graphics equipment by examining the

*/var/sysgen/system/irix.sm* file, in which each configured VME device is specified by a VECTOR line. When you examine a VECTOR line, note the following parameter values:

<i>bustype</i>	Specified as <i>VME</i> for VME devices. The VECTOR statement can be used for other types of buses as well.
<i>adapter</i>	The number of the VME bus where the device is attached.
<i>iospace</i> , <i>iospace2</i> , <i>iospace3</i>	Each <i>iospace</i> group specifies the VME address space and modifier, the starting bus address, and the size of a segment of VME address space used by this device.

Within each *iospace* parameter group you find keywords and numbers for the address space, modifier, and addresses for a device. The following is an example of a VECTOR line:

```
VECTOR: bustype=VME module=cdsio ipl=5 ctlr=0 adapter=0
iospace=(A24S,0xF00000,0x10000) probe_space=(A24S,0xF0FFFF,1)
```

This example specifies a VME device (*bustype=VME*) on bus 0 (*adapter=0*). The device resides in the A24 address space in supervisory mode (*iospace=(A24S...)*). Its first VME bus address is 0xF0 0000 and it covers a span of 0x01 0000 (64K) addresses—in other words, 0xF0 0000 through 0xF0 FFFF.

For third-party VME devices, look for a VECTOR line supplied by the manufacturer, usually stored in some file in */var/sysgen/system*.

### Opening a Device Special File

When you know the device addresses, you can open a device special file that represents the correct range of addresses. The device special files for VME mapping are found in */dev/vme*.

The naming convention for these files is documented in the *usrvme(7)* reference page. Briefly, each file is named **vmeBaSM**, where

<i>B</i>	is one or two digits for the bus number, for example 0 or 53
<i>S</i>	is two digits for the address space, 16, 24, or 32
<i>M</i>	is the modifier, either <i>s</i> for supervisory or <i>n</i> for nonprivileged

The device special file for the device described by the example VECTOR line in the preceding section would be */dev/vme/vme0a24s*.

In order to map a device on a particular bus and address space, you must open the corresponding file in */dev/vme*.

### Using the `mmap()` Function

When you have successfully opened the device special file, you use the file descriptor as the primary input parameter in a call to the `mmap()` system function.

This function has many different uses, all of which are documented in the `mmap(2)` reference page. For purposes of mapping a VME device into memory, the parameters should be as follows (using the names from the reference page):

<i>addr</i>	Should be NULL to permit the kernel to choose the address in user process space.
<i>len</i>	The length of the span of VME addresses, as documented in the <i>iospace</i> group in the VECTOR line.
<i>prot</i>	PROT_READ for input, PROT_WRITE for output, or the logical sum of those names when the device will be used for both input and output.
<i>flags</i>	MAP_SHARED. Add MAP_PRIVATE if this mapping is not to be visible to child processes created with the <code>sproc()</code> function (see the <code>sproc(2)</code> reference page).
<i>fd</i>	The file descriptor returned from opening the device special file in <i>/dev/vme</i> .
<i>off</i>	The starting VME bus address, as documented in the <i>iospace</i> group in the VECTOR line.

The value returned by `mmap()` is the virtual address that corresponds to the starting VME bus address. When the process accesses that address, the access is implemented by data transfer to the VME bus.

### Map Size Limits

There are limits to the amount and location of VME bus address space that can be mapped for PIO. The system architecture can restrict the span of mappable addresses, and kernel resource constraints can impose limits.

In all systems that support the VME bus it is possible to map all of A16 space.

In the Silicon Graphics Challenge and Onyx systems, all of A24 and A32 space can be used for PIO mappings, but there is a limit on the size of each map. Each bus mapping uses a hardware register that can span as much as 8 MB of contiguous VME bus addresses—so a single `mmap(0)` call can map at most 8 MB. There are as many as 12 mapping registers available for user mapping on each bus, so by making successive `mmap(0)` calls for adjacent 8 MB blocks of VME space you can map up to 96 MB of VME space into user process space from a single bus.

## VME PIO Access

Once a VME device has been mapped into memory, your program reads from the device by referencing the mapped address, and writes to the device by storing into the mapped address. Example 4-1 displays a sketch of a hypothetical function that maps a device and copies one register into another.

### Example 4-1 Opening and Using a Hypothetical VME Device

```
#define SPECFILE "/dev/vme/vmela16n"
typedef unsigned short int busdata; /* device data item */
typedef volatile busdata busreg;   /* device register */
#define MAPSIZE 8*sizeof(vmereg)
#define BUSADDR 0xff00
int vmefunc()
{
    busreg *mapped;
    busdata sample;
    int specfd = open(SPECFILE,O_RDWR);
    if (-1 == specfd) return error;
    mapped = mmap(NULL,          /* kernel pick address */
                  REGSIZE,      /* size of mapped area */
                  PROT_READ|PROT_WRITE, /* protection flags */
                  MAP_SHARED,    /* mapping flags */
                  specfd,        /* special file */
                  BUSADDR)      /* file offset */
    if (!mapped) return error;
    sample = busreg[0];          /* read A16N at 0xff00 */
    busreg[1] = sample;          /* write A16N at 0xff02 */
}
```

A PIO read is synchronous at the hardware level. The CPU executes a register-load instruction that does not complete until data has been returned from the device, up the system bus, to the CPU (see “CPU Access to Device Registers” on page 10). This can take 1 or 2 microseconds in a Challenge system.

A PIO write is not necessarily synchronous at the hardware level. The CPU executes a register-store instruction that is complete as soon as the physical address and data have been placed on the system bus. The actual VME write operation on the VME bus can take 1 or more microseconds to complete. During that time the CPU can execute dozens or even hundreds more instructions from cache memory.

### VME PIO Bandwidth

On a Challenge L or Onyx system, the maximum rate of PIO output is approximately 750K writes per second. The maximum rate of PIO input is approximately 250K reads per second. The corresponding data rate depends on the number of bytes transferred on each operation, as summarized in Table 4-1.

**Table 4-1** VME Bus PIO Bandwidth

Data Unit Size	Read	Write
D8	0.25 MB/second	0.75 MB/second
D16	0.5 MB/second	1.5 MB/second
D32	1 MB/second	3 MB/second

**Note:** The numbers in Table 4-1 were obtained by doing continuous reads, or continuous writes, to a device in the Challenge chassis. When reads and writes alternate, add approximately 1 microsecond for each change of direction. The use of a repeater to extend to an external card cage would add 200 nanoseconds or more to each transfer.

### EISA Programmed I/O

The EISA bus is supported by Silicon Graphics Indigo<sup>2</sup> workstation. If you are writing or maintaining a user-level EISA application to be executed on an Indigo<sup>2</sup>, the following information is of interest.

## Mapping an EISA Device Into Memory

As discussed in “Physical Device Addresses” on page 4, an I/O device is represented as an address or range of addresses in the address space of its bus. A kernel-level device driver has the ability to set up a mapping between the bus address of a device register and an arbitrary location in the address space of a user-level process. When this has been done, the device register appears to be a variable in memory—the program can assign values to it, or refer to it in expressions.

### Learning EISA Device Addresses

In order to map an EISA device for PIO, you must know the following points:

- which EISA bus adapter the device is on  
In all current Silicon Graphics systems, there is only one EISA bus adapter, and its number is 0.
- whether you need access to the EISA bus memory or I/O address space
- the address and length of the desired registers within the address space

You can find all these values by examining files in the */var/sysgen/system* directory, especially the */var/sysgen/system/irix.sm* file, in which each configured EISA device is specified by a VECTOR line. When you examine a VECTOR line, note the following parameter values:

<i>bustype</i>	Specified as <i>EISA</i> for EISA devices. The VECTOR statement can be used for other types of buses as well.
<i>adapter</i>	The number of the bus where the device is attached (0).
<i>iospace</i> , <i>iospace2</i> , <i>iospace3</i>	Each <i>iospace</i> group specifies the address space, starting bus address, and the size of a segment of bus address space used by this device.

Within each *iospace* parameter group you find keywords and numbers for the address space and addresses for a device. The following is an example of a VECTOR line (which must be a single physical line in the system file):

```
VECTOR: bustype=EISA module=if_ec3 ctlr=1
iospace=(EISAIO,0x1000,0x1000)
exprobe_space=(r,EISAIO, 0x1c80,4,0x6010d425,0xffffffff)
```

This example specifies a device that resides in the I/O space at offset 0x1000 (the slot-1 I/O space) for the usual length of 0x1000 bytes. The *exprobe\_space* parameter suggests that a key device register is at 0x1c80.

### Opening a Device Special File

When you know the device addresses, you can open a device special file that represents the correct range of addresses. The device special files for EISA mapping are found in */dev/eisa*.

The naming convention for these files is as follows: Each file is named **eisaBaM**, where

*B* is a digit for the bus number (0)

*M* is the modifier, either *io* or *mem*

The device special file for the device described by the example VECTOR line in the preceding section would be */dev/vme/eisa0aio*.

In order to map a device on a particular bus and address space, you must open the corresponding file in */dev/eisa*.

### Using the `mmap()` Function

When you have successfully opened the device special file, you use the file descriptor as the primary input parameter in a call to the **mmap()** system function.

This function is documented for all its many uses in the `mmap(2)` reference page. For purposes of mapping EISA devices, the parameters should be as follows (using the names from the reference page):

- addr* Should be NULL to permit the kernel to choose an address in user process space.
- len* The length of the span of bus addresses, as documented in the *iospace* group in the VECTOR line.
- prot* PROT\_READ, or PROT\_WRITE, or the logical sum of those names when the device is used for both input and output.
- flags* MAP\_SHARED, with the addition of MAP\_PRIVATE if this mapping is not to be visible to child processes created with the **sproc()** function (see the `sproc(2)` reference page).

- fd*            The file descriptor from opening the device special file in */dev/eisa*.
- off*           The starting bus address, as documented in the *iospace* group in the VECTOR line.

The value returned by **mmap()** is the virtual memory address that corresponds to the starting bus address. When the process accesses that address, the access is implemented by data transfer to the EISA bus.

**Note:** When programming EISA PIO, you must always be aware that EISA devices generally store 16-bit and 32-bit values in “small-endian” order, with the least-significant byte at the lowest address. This is opposite to the order used by the MIPS CPU under IRIX. If you simply assign to a C unsigned integer from a 32-bit EISA register, the value will appear to be byte-inverted.

### EISA PIO Bandwidth

The EISA bus adapter is a device on the GIO bus. The GIO bus runs at either 25 MHz or 33 MHz, depending on the system model. Each EISA device access takes multiple GIO cycles, as follows:

- The base time to do a native GIO read (of up to 64 bits) is 1 microsecond.
- A 32-bit EISA slave read adds 15 GIO cycles to the base GIO read time; hence one EISA access takes 19 GIO cycles, best case.
- A 4-byte access to a 16-bit EISA device requires 10 more GIO cycles to transfer the second 2-byte group; hence a 4-byte read to a 16-bit EISA slave requires 25 GIO cycles.
- Each wait state inserted by the EISA device adds four GIO cycles.

Table 4-2 summarizes best-case (no EISA wait states) data rates for reading and writing a 32-bit EISA device, based on these considerations.

**Table 4-2** EISA Bus PIO Bandwidth (32-Bit Slave, 33-MHz GIO Clock)

<b>Data Unit Size</b>	<b>Read</b>	<b>Write</b>
1 byte	0.68 MB/sec	1.75 MB/sec
2 byte	1.38 MB/sec	3.51 MB/sec
4 bytes	2.76 MB/sec	7.02 MB/sec

Table 4-3 summarizes the best-case (no wait state) data rates for reading and writing a 16-bit EISA device.

**Table 4-3** EISA Bus PIO Bandwidth (16-Bit Slave, 33-MHz GIO Clock)

<b>Data Unit Size</b>	<b>Read</b>	<b>Write</b>
1 byte	0.68 MB/sec	1.75 MB/sec
2 byte	1.38 MB/sec	3.51 MB/sec
4 bytes	2.29 MB/sec	4.59 MB/sec

## VME User-Level DMA

A DMA engine is included as part of each VME bus adapter in a Challenge or Onyx system. The DMA engine is unique to the Challenge architecture. It performs efficient, block-mode, DMA transfers between system memory and VME bus slave cards—cards that would normally be capable of only PIO transfers.

The DMA engine greatly increases the rate of data transfer compared to PIO, provided that you transfer at least 32 contiguous bytes at a time. The DMA engine can perform D8, D16, D32, D32 Block, and D64 Block data transfers in the A16, A24, and A32 bus address spaces.

## Using the `udmalib` Functions

All DMA engine transfers are initiated by a special device driver. However, you do not access this driver through `open/read/write` system functions. Instead, you program it through a library of functions. The functions are documented in the `udmalib(3)` reference page. They are used in the following sequence:

1. Call `dma_open()` to initialize action to a particular VME card.
2. Call `dma_allocbuf()` to allocate storage to use for DMA buffers.
3. Call `dma_mkparms()` to create a descriptor for an operation, including the buffer, the length, and the direction of transfer.
4. Call `dma_start()` to execute each transfer. This function does not return until the transfer is complete.

The `dma_start()` function takes the VME bus address of the particular slave device register that will provide or accept a series of data items. Before starting a DMA transfer, and possibly between transfers, you may need to program the VME controller with other commands. You would do this using PIO (see “VME Programmed I/O” on page 64).

**Tip:** The `dma_start()` function operates synchronously, polling the VME adapter hardware to find out when the DMA transfer is complete. In order to get parallel execution, consider calling `dma_start()` from a separate process.

### Buffer Allocation for User DMA

A buffer allocated by `dma_allocbuf()` is rounded up to a multiple of the memory page size, and is locked in memory to avoid page-faults during the DMA transfer. There is some overhead in creating a buffer, so for best performance the program should allocate the required buffers of the necessary size during initialization. However, if the total size of the buffers is a significant fraction of the available real memory, the large number of locked pages can hurt system performance.

You can only use the allocated buffer for DMA; it is not possible to provide your own buffer (for example, a buffer in a shared memory arena) for use by the DMA engine. When the data is produced by one process and written by another, this design can mean that the data has to be copied from an application buffer to the DMA buffer.

**Tip:** One way to avoid copying is to call `dma_allocbuf()` early, during program setup, before creating subprocesses using `sproc()`. Processes made with `sproc()` share their parent process address space, including buffer space created by `dma_allocbuf()`, so you

can have a process generating or consuming data in one part of the allocated buffer space while a different process executes **dma\_start()** to write or read data in a different part. It is of course essential to synchronize the use of the different buffer segments so that each area is used by only one process at a time.

### Allocation of Descriptors

Each call to **dma\_mkparms()** allocates a small block of memory and fills it with constants that describe a particular transfer operation. The primary input to **dma\_mkparms()** is a *vme\_parms\_t* object (declared in *udmalib.h*) containing the following important fields:

<i>vp_block</i>	1 for a block-mode transfer; 0 for a normal transfer.
<i>vp_datumsz</i>	The width of transfer units: VME_DS_BYTE (8-bit transfers) VME_DS_HALFWORD (16-bit transfers) VME_DS_WORD (32-bit transfers), or VME_DS_DBLWORD (64-bit transfers).
<i>vp_dir</i>	The direction of transfer: either VME_READ (from the VME bus) or VME_WRITE (to the VME bus).
<i>vp_throt</i>	For a block-mode transfer, the number of bytes to transfer in one burst without yielding the bus. Set to either VME_THROT_256 (the usual size, supported by most VME slaves that allow block transfer) or VME_THROT_2048 to allow up to 2,048 bytes per burst.
<i>vp_release</i>	The bus arbitration mode: VME_REL_RWD to release the bus as soon as the transfer is over, or VME_REL_ROR to release only when another bus master wants the bus. Use VME_REL_ROR for best speed when no bus masters are on the bus. Use VME_REL_RWD when other bus masters may be present and active.
<i>vp_addrmod</i>	The address modifier value that selects the target VME address space. Names of the form VME_*AMOD are declared for these values in <i>sys/vmereg.h</i> , for example VME_A16NPAMOD for the nonprivileged A16 space.

During initialization you call **dma\_mkparms()** to create a descriptor for each unique combination of parameters that your program will use. In each call to **dma\_start()**, you pass the descriptor that contains the appropriate set of parameters. A descriptor can be used in multiple **dma\_start()** calls.

## Advantages of User DMA

The `dma_start()` function and its related functions operate in user space; they do not make system calls to the kernel. This has two important effects. First, overhead is reduced, since there are no mode switches between user and kernel, as there are for `read()` and `write()`. This is important, because the DMA engine is often used for frequent, small inputs and outputs.

Second, `dma_start()` does not block the calling process, in the sense of suspending it and possibly allowing another process to use the CPU. It waits in a test loop until the operation is complete. As you can infer from Table 4-4, typical transfer times range from 50 to 250 microseconds. You can calculate the approximate duration of a call to `dma_start()` based on the amount of data and the operational mode.

You can use the `udmalib` functions to access a VME Bus Master device, if the device can respond in slave mode. However, this would normally be less efficient than using the Master device's own DMA circuitry.

While you can initiate only one DMA engine transfer per bus, it is possible to program a DMA engine transfer from each bus in the system, concurrently.

## DMA Engine Bandwidth

The maximum performance of the DMA engine for D32 transfers is summarized in Table 4-4. Performance with D64 Block transfers is somewhat less than twice the rate shown in Table 4-4. Transfers for larger sizes are faster because the setup time is amortized over a greater number of bytes.

**Table 4-4** VME Bus Bandwidth, DMA Engine, D32 Transfer

Transfer Size	Read	Write	Block Read	Block Write
32	2.8 MB/sec	2.6 MB/sec	2.7 MB/sec	2.7 MB/sec
64	3.8 MB/sec	3.8 MB/sec	4.0 MB/sec	3.9 MB/sec
128	5.0 MB/sec	5.3 MB/sec	5.6 MB/sec	5.8 MB/sec
256	6.0 MB/sec	6.7 MB/sec	6.4 MB/sec	7.3 MB/sec
512	6.4 MB/sec	7.7 MB/sec	7.0 MB/sec	8.0 MB/sec

**Table 4-4 (continued)** VME Bus Bandwidth, DMA Engine, D32 Transfer

Transfer Size	Read	Write	Block Read	Block Write
1024	6.8 MB/sec	8.0 MB/sec	7.5 MB/sec	8.8 MB/sec
2048	7.0 MB/sec	8.4 MB/sec	7.8 MB/sec	9.2 MB/sec
4096	7.1 MB/sec	8.7 MB/sec	7.9 MB/sec	9.4 MB/sec

**Note:** The throughput that can be achieved in VME DMA is very sensitive to several factors:

- The other activity on the VME bus.
- The blocksize (larger is better).
- Other overhead in the loop requesting DMA operations.

The loop used to generate the figures in Table 4-4 contained no activity except calls to `dma_start()`.

- the response time of the target VME board to a read or write request, in particular the time from when the VME adapter raises Data Strobe (DS) and the time the slave device raises Data Acknowledge (DTACK).

For example, if the slave device takes 500 ns to raise DTACK, there will always be fewer than 2 M data transfers per second.

### Example User DMA Function

The hypothetical function displayed in Example 4-2 is called to perform a series of DMA transfers from a specified device address. The code does not reflect any device initialization; all setup is presumably done by the caller. The code does not show the processing of the data, which is done in the hypothetical function `processOneBuffer()`.

**Example 4-2** User-Level DMA Access to VME

```

/*
|| This function assumes that any device programming needed to
|| prepare the device for input has been done (using PIO) before
|| the function is called. The bus number and device address
|| are function parameters.
*/

```

```

#define BLOCK_SIZE_TO_USE 4096
#include <udmalib.h>
#include <sys/vmereg.h>

extern void processOneBuffer(void *pBuffer);

int
readBlocks(int iBusNum, __uint32_t uiDevAddress)
{
    udmaid_t    *hEngine;    /* handle returned by dma_open */
    void        *pDMAbuffer; /* pointer from dma_allocbuf */
    vmeparms_t  sParms;     /* operation parms for dma_mkparms */
    udmaprmt_t *hParms;     /* handle returned by dma_mkparms */
    int         iStartCode; /* return code of dma_start */

    /*
    || Open the DMA engine. Terminate if it won't open.
    */
    hEngine = dma_open(DMA_VMEBUS, iBusNum);
    if (!hEngine)
    {
        perror("dma_open");
        return(-1);
    }

    /*
    || Allocate a special buffer for I/O. If that fails,
    || release the engine and terminate.
    */
    hDMAbuffer = dma_allocbuf(hEngine, BLOCK_SIZE_TO_USE);
    if (!hDMAbuffer)
    {
        perror("dma_allocbuf");
        dma_close(hEngine);
        return(-2);
    }

    /*
    || Set up the VME parameters and "make" them. A different set
    || of parameters is needed for each combination of vmeparms_t
    || values, buffer, and size. This example uses only one set.
    */
    sParms.vp_block = 0;                /* this device does not do blocks */
    sParms.vp_datumsz = VME_DS_WORD;    /* this is a 32-bit device */
    sParms.vp_dir = VME_READ;           /* input operation */
    sParms.vp_throt = VME_THROT_256;    /* smaller burst size */
    sParms.vp_release = VME_REL_ROR;    /* release on request */
    sParms.vp_addrmod = VME_A32NPAMOD; /* address modifier */

```

```
hParms = dma_mkparms(hEngine, &sParms, pDMAbuffer, BLOCK_SIZE_TO_USE);
if (!hParms)
{
    perror("dma_mkparms");
    dma_freebuf(pDMAbuffer);
    dma_close(hEngine);
    return(-3);
}
/*
|| Read and process blocks until error.
*/
for(iStartCode=0;iStartCode==0;)
{
    iStartCode = dms_start(hEngine, uiDevAddress, hParms);
    if (!iStartCode)
        processOneBuffer(pDMAbuffer);
}
/*
|| Clean up and exit.
*/
dma_freeparms(hEngine, hParms);
dma_freebuf(pDMAbuffer);
dma_close(hEngine);
return 0;
}
```

## PCI Programmed I/O

For an overview of the PCI bus and its hardware implementation in Silicon Graphics systems, see Chapter 15, "PCI Device Drivers."

### Mapping a PCI Device Into Process Address Space

As discussed in "Physical Device Addresses" on page 4, an I/O device is represented as an address, or range of addresses, in the address space of its bus. A kernel-level device driver has the ability to set up a mapping between an address on an I/O bus and an arbitrary location in the address space of a user-level process. When this has been done, the bus location appears to be a variable in memory. The program can assign values to it, or refer to it in expressions.

The PCI bus addresses managed by a device are not wired or jumpered into the board; they are established dynamically at the time the system attaches the device. The assigned bus addresses can vary from one day to the next, as devices are added to or removed from that PCI bus adapter. For this reason, you cannot know the bus addresses of a PCI device in advance, so as to write them into a configuration file (like the VECTOR statements that document the bus addresses of VME and EISA devices) or into the source code of a program.

In order to map bus addresses for a particular device, you must open the device special file that represents that device. You pass the file descriptor for the opened device to the **mmap()** function. If the device driver for the device supports memory mapping—mapping is an optional feature of a PCI device driver—the mapping is set up.

The PCI bus defines three address spaces: configuration space, I/O space, and memory space. It is up to the device driver which of the spaces it allows you to map. Some device drivers may set up a convention allowing you to map in different spaces.

### Opening a Device Special File

The device special files for PCI devices are established by the system administrator. You need to know the pathname of the file in */dev* in order to open the device. This pathname is passed to the **open()** system function, along with flags representing the type of access (see the **open(2)** reference page).

### Using the **mmap()** Function

When you have successfully opened the device special file, you use the file descriptor as the primary input parameter in a call to the **mmap()** system function.

This function is documented for all its many uses in the **mmap(2)** reference page. For purposes of mapping a PCI device into memory, the parameters should be as follows (using the names from the reference page):

<i>addr</i>	Should be NULL to permit the kernel to choose an address in user process space.
<i>len</i>	The length of the span of PCI addresses to map.
<i>prot</i>	PROT_READ for input, PROT_WRITE for output, or the logical sum of those names when the device will be used for both input and output.

<i>flags</i>	MAP_SHARED. Add MAP_PRIVATE if this mapping is not to be visible to child processes created with the <b>sproc()</b> function (see the sproc(2) reference page).
<i>fd</i>	The file descriptor returned from opening the device special file in <i>/dev/vme</i> .
<i>off</i>	The offset into the device address space. The treatment of this value is up to the device driver, which can interpret it different ways. Commonly, the driver will treat this as an offset into the memory address space defined by the device.

The value returned by **mmap()** is the virtual address that corresponds to the starting VME bus address. When the process accesses that address, the access is implemented by data transfer to the VME bus.

### Map Size Limits

There are limits to the amount and location of PCI bus address space that can be mapped for PIO. The system architecture can restrict the span of mappable addresses, and kernel resource constraints can impose limits. In order to create the map, the PCI device driver has to create a software object called a PIO map. In some systems, only a limited number of PIO maps can be active at one time. However, in the O2 workstation, the number of PIO maps is limited only by kernel memory constraints.

---

## User-Level Access to SCSI Devices

IRIX contains a programming library, called *dslib*, that allows you to control SCSI devices from a user-level process. This chapter documents the functions in *dslib*, including the following topics:

- “Overview of the dsreq Driver” on page 82 gives a summary of the features and use of the generic SCSI device driver.
- “Generic SCSI Device Special Files” on page 82 documents the format of the names and major and minor numbers of generic SCSI files.
- “The dsreq Structure” on page 85 gives details of the request structure that is the primary input to the generic SCSI driver.
- “Testing the Driver Configuration” on page 92 documents the use of the `DS_CONF ioctl()` operation.
- “Using dslib Functions” on page 94 describes the functions that make it simpler to use the generic SCSI driver.
- “Using the Special DS\_RESET and DS\_ABORT Calls” on page 93 describes two special functions of the generic SCSI driver.

You must understand the SCSI interface in order to command a SCSI device. For several SCSI information resources, see “Other Sources of Information” on page xxvii.

If you are specifically interested in using audio data from a CDROM or DAT drive, you should use the special-purpose libraries for CDROM and DAT that are included in the IRIS Digital Media Development Environment. These libraries are built upon the generic SCSI driver, but provide convenient, audio-oriented functions. For more information on these libraries, see the *IRIS Digital Media Programming Guide*, document number 008-1799-040.

If your interest is in controlling SCSI devices at the kernel level, see Part IV, “SCSI Device Drivers.”

## Overview of the *dsreq* Driver

IRIX includes a generic SCSI device driver, the *dsreq* driver, through which a user-level program can issue SCSI commands to SCSI devices. This is a character device driver that supports only **open()**, **close()** and **ioctl()** operations (see “Kinds of Kernel-Level Drivers” on page 48, and also the *open(2)*, *close(2)* and *ioctl(2)* reference pages).

The formal documentation of the *dsreq* driver is found in the *ds(7)* reference page. In order to invoke its services, you prepare a *dsreq* data structure describing the operation and pass it to the device driver using an **ioctl()** call. The device driver issues the SCSI command you specify, and sleeps until it has completed. Then it returns the status in the *dsreq* structure.

You can request operations for input and output as well as issuing control and diagnostic commands. The *dsreq* structure for input and output operations specifies a buffer in memory for data transfer. The *dsreq* driver handles the task of locking the buffer into memory (if necessary) and managing a DMA transfer of data.

The programming interface supported by the generic SCSI driver is quite primitive. A library of higher-level functions makes it easier to use. This library is formally documented in the *dslib(3)* reference page, and is described under “Using *dslib* Functions” on page 94.

## Generic SCSI Device Special Files

The creation and use of device special files is discussed under “Device Special Files” on page 34. A device special file represents a device, and is the mechanism for associating a device with a kernel-level device driver.

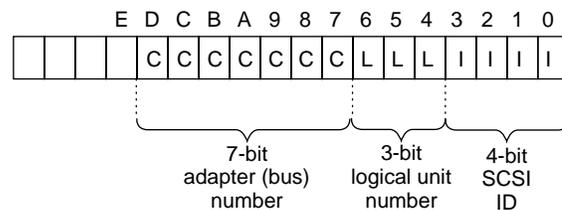
The device special files in the */dev/scsi* directory are all associated with the *dsreq* driver. A basic set of these names is created automatically by the */dev/MAKEDEV* script (see “The Script *MAKEDEV*” on page 38). You have to create additional device special files if you need to control logical units other than logical unit 0.

## Major and Minor Device Numbers in /dev/scsi

Device special files in */dev/scsi* have one of the following major device numbers:

- 195 for devices on a SCSI bus (files */dev/scsi/sc\**).
- 196 for devices on a *jag* (VME) SCSI bridge (files */dev/scsi/jag\**).

The minor number of these files encodes the adapter number, the SCSI ID, and the LUN, using the bit assignments shown in Figure 5-1.



**Figure 5-1** Bit Assignments in SCSI Device Minor Numbers

## Form of Filenames in /dev/scsi

Each device special filename in the */dev/scsi* directory reflects the values of the device's adapter (bus) number, SCSI ID, and logical unit number (LUN).

**Tip:** The character between the SCSI ID and the LUN in these names is the letter "l." When reading or copying these device names, take care not to write a digit 1 instead. This is a frequent error.

### Names of SCSI Devices on a SCSI Bus

Devices attached directly to a SCSI bus have names in this form:

<b>sc</b>	Prefix "sc" for SCSI attachment.
0 to 137	Number of the SCSI adapter, typically 0 or 1.
<b>d</b>	Constant letter "d" for device.
0 to 7 (to 15 for wide SCSI)	SCSI ID of the target device or control unit, as set by switches on the device itself.

- l (letter ell)      Constant letter "l" for logical unit.
- 0 to 7              Logical unit number (LUN) of this device, typically 0.

A typical device name would be `/dev/scsi/sc1d3l0` meaning a SCSI device configured as ID 3 on SCSI bus 1. Either this device has no logical units, or this is the first logical unit on device 3.

### Names of SCSI Devices on the Jag (VME Bus) Controller

Machines in the Challenge and Onyx systems can optionally have SCSI devices attached to the VME bus through a bridge using the `jag` device driver. These devices are also represented in `/dev/scsi` with names of the following form:

- jag**              Prefix "jag" for VME/SCSI attachment.
- 0 to 4              Number of the VME adapter, typically 0 or 1.
- d**                Constant letter "d" for device.
- 0 to 7 (to 15 for wide SCSI)      SCSI ID of the target device or control unit, as set by switches on the device itself.
- l (letter ell)      Constant letter "l" for logical unit.
- 0 to 7              Logical unit number (LUN) of this device, typically 0.

A typical device name would be `/dev/scsi/jag1d3l0` meaning a SCSI device configured as ID 3 on VME bus 1. Either the device has no logical units, or this is the first logical unit on device 3.

### Creating Additional Names in `/dev/scsi`

The script `/dev/MAKEDEV`, which runs each time the system boots, creates 16 files for each existing SCSI or jag bus. These files represent the possible SCSI ID numbers 0-15 on each bus, with a logical number of 0. If you want to control a device with LUN 0, the device special file exists.

In order to control a device with a LUN of 1-7, you must create an additional device special file, using the *mknod* or *install* command (see the *install(1)* reference page). For example, before you can operate logical unit 2 of device 5 on SCSI bus 1, you must create */dev/scsi/sc1d5l2* using a command such as

```
install -F /dev/scsi -m 600 -u root -g sys \  
-chr 195,165 sc1d5l2
```

### Relationship to Other Device Special Files

The files in */dev/scsi* describe many of the same devices that are described by files in */dev/dsk*, */dev/tape*, and other directories. There is a security exposure in that a user-level program could use a */dev/scsi* file to do almost anything to a disk or tape, including total erasure.

The *dsreq* device driver forces exclusivity with itself; that is, a given */dev/scsi* file can be opened only by one process at a time. However, a device could be open through the *dsreq* driver at the same time it is open by another process, or by a filesystem, through a different device special file and device driver. For example, a disk volume could be simultaneously open through the name */dev/scsi/sc0d0l0* and through */dev/rdisk/dks0d1s0*.

The process that opens a generic SCSI device can request exclusivity using the `O_EXCL` option to `open()`. In that case, the open is rejected when the device is already open through another driver; and no other driver can open the device until the generic device file is closed.

## The dsreq Structure

The primary input to most *dsreq* `ioctl()` calls, as well as the primary input to most `dslib` functions, is the *dsreq* structure. This structure is declared in */usr/include/sys/dsreq.h*, a header file that rewards careful study.

The important fields of the *dsreq* structure are shown in Table 5-1. Some of the field values are expanded in the following topics. The *sys/dsreq.h* file declares macros for access to many of the fields. Use these macros (listed in Table 5-1) in both expressions and assignments in order to insulate your code against future changes.

**Table 5-1** Fields of the *dsreq* Structure

Field Name	Macro	Purpose
<i>ds_flags</i>	FLAGS(dp)	Bits used to determine device driver actions. See “Values for <i>ds_flags</i> ” on page 87.
<i>ds_time</i>	TIME(dp)	Timeout value in milliseconds. If the command does not complete, it is ended with an error code. The driver sets a default of 5000 (5 seconds) when this is set to zero. <b>dsopen()</b> initializes it to 10000.
<i>ds_private</i>	PRIVATE(dp)	Field for use by the calling program. <b>dsopen()</b> uses this field to point to its “context” data (see “Using <i>dsopen()</i> and <i>dsclose()</i> ” on page 95).
<i>ds_cmdbuf</i>	CMDBUF(dp)	Address of SCSI command string to be sent.
<i>ds_cmdlen</i>	CMDLEN(dp)	Length of the SCSI command string.
<i>ds_databuf</i>	DATABUF(dp)	Address of a single data buffer. See “Data Transfer Options” on page 89.
<i>ds_dataalen</i>	DATALEN(dp)	Length of data buffer.
<i>ds_sensebuf</i>	SENSEBUF(dp)	Address to receive sense data after an error.
<i>ds_senselen</i>	SENSELEN(dp)	Length of sense buffer in bytes.
<i>ds_iovbuf</i>	IOVBUF(dp)	Address of an <i>iov_t</i> structure. See “Data Transfer Options” on page 89.
<i>ds_iovlen</i>	IOVLEN(dp)	Length of data described by <i>ds_iovbuf</i> .
<i>ds_link</i>		This field is not supported, and should be zero-filled.
<i>ds_synch</i>		This field is not supported, and should be zero-filled.
<i>ds_revcode</i>		Intended for the version code of the <i>dsreq</i> driver, not currently set to a useful value.
<i>ds_ret</i>	RET(dp)	Return code for the requested operation. See Table 5-3 on page 89.

**Table 5-1 (continued)** Fields of the dsreq Structure

Field Name	Macro	Purpose
<i>ds_status</i>	STATUS(dp)	SCSI status byte from the operation. See Table 5-4 on page 91.
<i>ds_msg</i>	MSG(dp)	The first byte of a message returned by the target. See Table 5-5 on page 91.
<i>ds_cmdsent</i>	CMDSENT(dp)	Length of command string actually sent (same as <i>ds_cmdlen</i> , unless an error occurs).
<i>ds_datasent</i>	DATASENT(dp)	Length of data transferred.
<i>ds_sensesent</i>	SENSESENT(dp)	Length of sense data received.

The dslib library contains functions to simplify the preparation and execution of a *dsreq* request; see “Using dslib Functions” on page 94.

## Values for ds\_flags

The possible flag values in the *ds\_flags* field are listed in Table 5-2. The flag values are designed for the most flexible, capable type of bus, device, and device driver. Not all values are supported, and different host adapters can support different combinations.

**Table 5-2** Flag Values for ds\_flags

Constant Name	Supported by Any Driver?	Meaning When Set to 1
DSRQ_ASYNC	Yes	Return at once, do not sleep until the operation is complete.
DSRQ_SENSE	Yes	Get sense data following an error on the requested command.
DSRQ_TARGET	No	Act as the SCSI target, not the SCSI initiator.
DSRQ_SELATN	Yes	Select with ATN.
DSRQ_DISC	Yes	Allow identify disconnect.
DSRQ_SYNXFR	Yes	Negotiate a synchronous transfer if possible. Needed only to switch into synchronous mode. Repeated negotiation is wasteful.

**Table 5-2 (continued)** Flag Values for `ds_flags`

Constant Name	Supported by Any Driver?	Meaning When Set to 1
<code>DSRQ_ASYNXFR</code>	Yes	Negotiate an asynchronous transfer. Needed only to return to asynch after a synchronous transfer. Repeated negotiation is wasteful.
<code>DSRQ_SELMSG</code>	No	A specific select is coded in the message. This feature is not supported.
<code>DSRQ_IOV</code>	Yes	Use the <code>iov_t</code> from <code>ds_iovbuf</code> , not the single buffer from <code>ds_databuf</code> (see “Data Transfer Options” on page 89).
<code>DSRQ_READ</code>	Yes	This is a data input command (as opposed to an immediate command or an output).
<code>DSRQ_WRITE</code>	Yes	This is a data output command (as opposed to an immediate command or an input).
<code>DSRQ_MIXRDWR</code>	No	This command can both read and write.
<code>DSRQ_BUF</code>	No	Buffer the input and copy to the supplied buffer, instead of direct input to the buffer.
<code>DSRQ_CALL</code>	No	Notify completion (with <code>DSRQ_ASYNC</code> ).
<code>DSRQ_ACKH</code>	No	Hold ACK asserted.
<code>DSRQ_ATNH</code>	No	Hold ATN asserted.
<code>DSRQ_ABORT</code>	No	Send ABORT messages until the bus is clear. Useful only with SCSI commands that have the immediate bit set.
<code>DSRQ_TRACE</code>	Yes	Trace this request (accepted but has no effect).
<code>DSRQ_PRINT</code>	Yes	Print this request (accepted but has no effect).
<code>DSRQ_CTRL1</code>	Yes	Request with host control bit 1.
<code>DSRQ_CTRL2</code>	Yes	Request with host control bit 2.

In order to find out which flags are supported by a particular driver, use the `DS_CONF` operation (see “Testing the Driver Configuration” on page 92).

## Data Transfer Options

When reading or writing data, you have two design options:

- You can transfer a single segment of data directly between the device and a buffer you supply (set neither `DSRQ_BUF` nor `DSRQ_IOV`).
- You can transfer segments of data between the device and a series of one or more memory locations based on an `iov_t` object (set `DSRQ_IOV`).

All read/write requests are done using DMA. The “scatter/gather” support of `DSRQ_IOV` is presently restricted to only one memory segment, so it is not greatly different from single-buffer I/O. If you elect to use it, the `iov_t` structure is declared in `sys/iov.h` (see also the part of the `read(2)` reference page that deals with the `readv()` function).

During a direct transfer using either a single buffer or scatter/gather, the data buffer spaces are locked in memory.

The maximum amount of data you can transfer in one operation is set by the host adapter driver for the bus, and can be retrieved with an `ioctl()` (see “Testing the Driver Configuration” on page 92). The maximum length for a buffered transfer is returned by the same `ioctl()`. It can be less than the direct-transfer size because there may be a limit on the size of kernel memory that can be allocated.

## Return Codes and Status Values

A zero return code in the `ds_ret` field signifies success. The possible nonzero return codes are summarized in Table 5-3 and are declared in `sys/dsreq.h`. Not all return codes are possible with every driver.

**Table 5-3** Return Codes From SCSI Operations

Constant Name	Meaning
<code>DSRT_DEVSCSI</code>	General failure from SCSI driver.
<code>DSRT_MULT</code>	General software failure, typically a SCSI-bus request.
<code>DSRT_CANCEL</code>	Operation cancelled in host adapter driver.
<code>DSRT_REVCODE</code>	Software level mismatch, recompile application.

**Table 5-3 (continued)** Return Codes From SCSI Operations

Constant Name	Meaning
DSRT_Again	Try again, recoverable SCSI-bus error.
DSRT_HOST	Failure reported by host adapter driver for the bus in use.
DSRT_NOSEL	No unit responded to select.
DSRT_SHORT	Incomplete transfer (not an error). See <i>ds_datasent</i> .
DSRT_OK	Not returned at this time.
DSRT_SENSE	Command returned with status; sense data successfully retrieved from SCSI host (see <i>ds_sensesent</i> ).
DSRT_NOSENSE	Command with status, error occurred while trying to get sense data from SCSI host.
DSRT_TIMEOUT	Command did not complete in the time allowed by <i>ds_timeout</i> .
DSRT_LONG	Data transfer overran bounds ( <i>ds_dataalen</i> ).
DSRT_PROTO	Miscellaneous protocol failure.
DSRT_EBSY	Busy dropped unexpectedly; protocol error.
DSRT_REJECT	Message rejected; protocol error.
DSRT_PARITY	Parity error on SCSI bus; protocol error.
DSRT_MEMORY	Memory error in system memory.
DSRT_CMDO	Protocol error during command phase.
DSRT_STAI	Protocol error during status phase.
DSRT_UNIMPL	Command not implemented; protocol error.

The possible SCSI status value in the *ds\_status* field are summarized in Table 5-4.

**Table 5-4** SCSI Status Codes

Constant Name	Meaning
STA_GOOD	The target has successfully completed the SCSI command.
STA_CHECK	An error or exception was detected. Sense was attempted if DSRQ_SENSE was specified.
STA_BUSY	Command not attempted; addressed unit is busy.
STA_IGOOD	Linked SCSI command completed.
STA_RESERV	Command aborted because it tried to access a logical unit or an extent within a logical unit that reserves that type of access to another SCSI device.

The possible SCSI message byte values in the *ds\_msg* field are summarized in Table 5-5.

**Table 5-5** SCSI Message Byte Values

Constant Name	Meaning
MSG_COMPL	Command complete.
MSG_XMSG	Extended message (only byte returned).
MSG_SAVEP	Initiator should save data pointers.
MSG_RESTP	Initiator restore data pointers.
MSG_DISC	Disconnect.
MSG_IERR	Initiator detected error.
MSG_ABORT	Abort.
MSG_REJECT	Optional message rejected, not supported.
MSG_NOOP	Empty message.
MSG_MPARITY	Parity error during Message In phase.
MSG_LINK	Linked command complete.
MSG_LINKF	Linked command complete with flag.

**Table 5-5 (continued)** SCSI Message Byte Values

Constant Name	Meaning
MSG_BRESET	Bus device reset.
MSG_IDENT	Value 0x80, first of the 0x80-0xFF identifier messages.

## Testing the Driver Configuration

Different buses have different host adapter drivers that can have different features. The *dsreq* device driver supports an **ioctl()** call that retrieves the configuration of the driver for the bus where the device resides. This call fills in the fields of a structure of type *dsconf* (declared in *sys/dsreq.h*) listed in Table 5-6.

**Table 5-6** Fields of the *dsconf* Structure

Field Name	Contents
<i>dsc_flags</i>	DSRQ flags honored by this driver (see Table 5-2 on page 87)
<i>dsc_preset</i>	DSRQ preset values (defaults) that are merged with the input <i>ds_flags</i> using logical OR in any request.
<i>dsc_bus</i>	Number of this SCSI bus, as encoded in the device minor number.
<i>dsc_imax</i>	Maximum target ID for this bus (7 for SCSI, 15 for wide SCSI).
<i>dsc_lmax</i>	Maximum number LUN values per ID on this bus.
<i>dsc_iomax</i>	Maximum length of a single I/O transfer.
<i>dsc_biomax</i>	Maximum length of a buffered I/O transfer.

The code in Example 5-1 shows a function that tests if a particular flag is supported by a particular bus. The input arguments are a file descriptor for an open device special file, and a flag value (or values) from *sys/dsreq.h*.

**Example 5-1** Testing the Generic SCSI Configuration

```

uint
test_dsreq_flags(int dev_fd, uint flag)
{
    dsconf_t config;
    int ret;
    ret = ioctl(dev_fd, DS_CONF, &config);
    if (!ret) { /* no problem in ioctl */
        return (flag & config.dsc_flags);
    } else { /* ioctl failure */
        return 0; /* not supported, it seems */
    }
}

```

A program could use the function in Example 5-1 to find out if a particular feature is supported. For example, a test of support for the DSRQ\_SYNXFER feature could be coded as follows:

```

if (test_dsreq_flags(the_dev, DSRQ_SYNXFER)) {
    /* synchronous negotiation is supported */...
}

```

**Using the Special DS\_RESET and DS\_ABORT Calls**

Two special functions of the generic SCSI driver are available only as `ioctl()` calls, not through `dslib` functions.

**Using DS\_ABORT**

The `DS_ABORT ioctl()` sends a SCSI ABORT message to the bus, target, and LUN defined by the file descriptor. The resulting status is returned in the `dsreq` that is also specified. The host adapter driver waits until no commands are pending on that bus, so there is no point in using this function to cancel anything but an immediate command such as a rewind. An example of this call is as follows:

```

ioctl(dev_fd, DS_ABORT, &some_dsreq);

```

## Using DS\_RESET

The DS\_RESET `ioctl()` function causes a reset of the SCSI bus specified by the file descriptor. The resulting status is returned in the `dsreq` that is also specified. This powerful operation should be used with great care, because it terminates all pending activity on the bus.

## Using dslib Functions

The functions in the dslib library are built upon calls to the `dsreq` device driver, and simplify the process of allocating a `dsreq` structure, setting values in it, and executing commands. The formal documentation of the library is found in `dslib(3)`. The source code is distributed with the system in the `/usr/share/src/irix/examples/scsi` directory so that you can read and extend it. (This directory installs as part of the `irix_dev` software component, and the `examples` directory does not install by default.)

### dslib Functions

In order to use the functions in the library, you include `/usr/include/dslib.h` in your code, and link with the `-lds` option so as to link `/usr/lib/libds.so`. Then the functions summarized in Table 5-7 are available.

**Table 5-7** dslib Function Summary

Function Name	Purpose	Page
<code>ds_ctostr</code>	Look up a string in a table using an integer key.	page 99
<code>ds_vtostr</code>	Look up a string in a table using an integer key.	page 99
<code>dsopen</code>	Open a device special file and allocate a <code>dsreq</code> for use with it.	page 95
<code>dsfclose</code>	Free the <code>dsreq</code> structure and close the device.	page 95
<code>doscsireq</code>	Perform an operation on a device as specified in a <code>dsreq</code> .	page 97
<code>filldsreq</code>	Set values in fields of a <code>dsreq</code> structure.	page 97
<code>fillg0cmd</code>	Set up the <code>dsreq</code> structure for a group 0 SCSI command.	page 98
<code>fillg1cmd</code>	Set up the <code>dsreq</code> structure for a group 1 SCSI command.	page 98

**Table 5-7 (continued)** dslib Function Summary

Function Name	Purpose	Page
<b>inquiry12</b>	Issue an Inquiry command and retrieve information from the device concerning such things as its type.	page 100
<b>modeselect15</b>	Issue a group 0 Mode Select command to a SCSI device.	page 100
<b>modesense1a</b>	Send a group 0 Mode Sense command to a device to retrieve a parameter page from the device.	page 101
<b>read08</b>	Issue a group 0 Read command in disk-drive form.	page 102
<b>readextended28</b>	Issue a group 1 Read command in disk-drive form.	page 102
<b>readcapacity25</b>	Issue a Read Capacity command.	page 103
<b>requestsense03</b>	Issue a Request Sense command and test or probe for the device.	page 104
<b>reserveunit16</b>	Issue a Reserve Unit command.	page 104
<b>releaseunit17</b>	Issue a Release Unit command.	page 104
<b>senddiagnostic1d</b>	Issue a Send Diagnostic command to test if the device or the SCSI bus is online, or run a self-test on the device.	page 105
<b>testunitready00</b>	Issue a Test Unit Ready command to the SCSI device.	page 106
<b>write0a</b>	Issue a group 0 Write command to the SCSI device.	page 106
<b>writeextended2a</b>	Issue an extended Write command to the SCSI device.	page 106

### Using dsopen() and dsfclose()

The **dsopen()** function opens a device special file for a generic SCSI device, and allocates a *dsreq* structure initialized for use with that device. The function prototype is

```
struct dsreq* dsopen(char *opath, int oflags);
```

The arguments are

- opath*            The name of the device special file as a character string, for example “/dev/scsi/jag0d710” (see “Form of Filenames in /dev/scsi” on page 83).
- oflags*            The *oflag* value expected by **open()** when opening this device special file. **O\_EXCL** has special meaning; see “Relationship to Other Device Special Files” on page 85.

If the **open()** call fails or memory cannot be allocated, the function returns **NULL**. Otherwise it allocates a *dsreq* structure as well as generous buffers for command and sense strings. The following fields of the *dsreq* are initialized:

- ds\_time*            Set to 10000 (10 second timeout).
- ds\_private*        Set to the address of the context that contains the *dsreq* as well as the command and sense buffers.
- ds\_cmdbuf*        Set to the address of the command buffer.
- ds\_cmdlen*        Set to the length of the allocated command buffer.
- ds\_sensebuf*      Set to the address of the allocated sense buffer.
- ds\_senselen*      Set to the length of the sense buffer.

Other fields of the *dsreq* are cleared to zero.

**Note:** Other functions in *dslib* assume that a *dsreq* has been initialized by **dsopen()**. In particular they assume the *ds\_private* value points to a context block. You should not attempt to use any *dsreq* structure with a *dslib* function except one returned by **dsopen()**; and you should not use a *dsreq* opened for one file with another file.

The **dsfclose()** function releases the *dsreq* structure and close the device. Its prototype is `void dsfclose(struct dsreq *dsp);`

The only argument is the *dsreq* created by **dsopen()**.

## Issuing a Request With `doscsireq()`

The `doscsireq()` function issues a SCSI request by passing a *dsreq* to the SCSI device driver using an `ioctl()` call. The *dsreq* must have been prepared completely beforehand. The function prototype is

```
int doscsireq(int fd, struct dsreq *dsp);
```

The arguments are as follows:

*fd*                    The file descriptor for the open device file.  
*dsp*                    The address of the *dsreq* prepared by `dsopen()`.

Normally the returned value is the SCSI status byte. When the requested operation ends with Busy or Reserve Conflict status, the function sleeps 2 seconds and tries the operation up to four times. The returned value is -1 when the device driver rejects the `ioctl()` or the third retry ends in failure.

## SCSI Utility Functions

The functions `filldsreq()`, `fillg0cmd()`, `fillg1cmd()`, `ds_vtostr()`, and `ds_ctostr()` are not oriented toward particular SCSI operations, but are used to construct your own task-oriented SCSI functions.

### Using `filldsreq()`

The `filldsreq()` function is used to set the *ds\_flags*, *ds\_databuf*, and *ds\_datalen* members of a *dsreq* structure. Its prototype is

```
void filldsreq(struct dsreq *dsp, uchar_t *data, long datalen, long flags)
```

The arguments are as follows:

*dsp*                    The address of a *dsreq* prepared by `dsopen()`.  
*data*                    The address of a buffer area.  
*datalen*                The length of the buffer area.  
*flags*                    Flag values for *ds\_flags* (see “Values for *ds\_flags*” on page 87).

The bits in *flags* are added to *ds\_flags* with an OR; they do not replace the contents of the field.

**Note:** Besides the specified values, the function also sets 10000 in *ds\_timeout* and clears *ds\_link*, *ds\_synch*, and *ds\_ret* to zero.

### Using `fillg0cmd()` and `fillg1cmd()`

The `fillg0cmd()` function stores a group 0 (6-byte) SCSI command in a command buffer. The `fillg1cmd()` stores a group 1 (10-byte) SCSI command in the buffer. Both functions set the *ds\_cmdbuf* and *ds\_cmdlen* fields of a *dsreq*. The function prototypes are:

```
void fillg0cmd(struct dsreq *dsp, uchar_t *cmdbuf, b0, ..., b5)
void fillg1cmd(struct dsreq *dsp, uchar_t *cmdbuf, b0, ..., b9)
```

The arguments are as follows:

<i>dsp</i>	The address of any <i>dsreq</i> .
<i>cmdbuf</i>	The address of a buffer to receive the command string.
<i>b0, b1,...</i>	Expressions for the successive bytes of a SCSI command.

In typical use, the arguments are as follows:

<i>dsp</i>	The address of a <i>dsreq</i> initialized by <code>dsopen()</code> .
<i>cmdbuf</i>	The command buffer allocated by <code>dsopen()</code> , whose address is stored in the <i>ds_cmdbuf</i> field of the <i>dsreq</i> .
<i>b0</i>	A SCSI command verb expressed as one of the constants declared in <code>dslib.h</code> , for example <code>G0_INQU</code> .

A typical call resembles the following:

```
fillg0cmd(dsp, (uchar_t *)CMDBUF(dsp), G0_INQU, 1, inq_page, 0,
B1(datalen), 0);
```

The macros `B1()`, `B2()`, and `B4()` defined in `sys/dsreq.h` are useful for expressing halfword and word values as byte sequences.

### Using `ds_vtostr()` and `ds_ctostr()`

The dslib library module contains six static tables that can be used to convert between numeric values and character strings for message display. The tables are summarized in Table 5-8. The table definitions are in the source file *dstab.c*.

**Table 5-8** Lookup Tables in dslib

External Name	Type	Table Contents
cmdnametab	vtab	Names for SCSI command bytes, for example "Test Unit."
cmdstatustab	vtab	Names for SCSI status byte codes, for example "BUSY."
dsrqnametab	vtab	Descriptions of flag values from <i>ds_flags</i> , for example "select with (without) attn" for DSRQ_SELATN.
dsrtnametab	vtab	Descriptions of return values in <i>ds_ret</i> , for example "parity error on SCSI bus" for DSRT_PARITY.
msgnametab	vtab	Descriptions of SCSI message bytes, for example "Save Pointers."
sensekeytab	ctab	Descriptions of SCSI sense byte values, for example "Illegal Request."

The `ds_vtostr()` function searches any of the five vtab tables for the string matching an integer key. The `ds_ctostr()` function searches a ctab (currently, only sensekeytab is a ctab) for the string matching a key. The function prototypes are

```
char * ds_vtostr(unsigned long v, struct vtab *table);
char * ds_ctostr(unsigned long v, struct ctab *table);
```

Each function searches the specified table for a row containing the numeric value *v*, and returns address of the corresponding string. If there is no such row, the functions return the address of a zero-length string.

## Using Command-Building Functions

The remaining functions in `dslib` each construct and execute a specific type of common SCSI command. Each function follows this general pattern:

1. Use `fillg0cmd()` or `fillg1cmd()` to set up the command string, based on the function's arguments.
2. Use `filldsreq()` to set up the remaining fields of the `dsreq` structure.
3. Execute the command using `doscsireq()`.
4. Return the value returned by `doscsireq()`.

You can construct similar, additional functions using the utility functions in this same way. In particular you are likely to need to construct your own function to issue Read commands.

### **inquiry12()**—Issue an Inquiry Command

The `inquiry12()` function prepares and issues an Inquiry command to retrieve device-specific information. The function prototype is

```
int inquiry12(struct dsreq *dsp, caddr_t data, long datalen, int vu);
```

The arguments are as follows:

<i>dsp</i>	The address of a <code>dsreq</code> structure prepared by <code>dsopen()</code> .
<i>data</i>	The address of a buffer to receive the inquiry response.
<i>datalen</i>	The length of the buffer, at least 36 and typically 64.
<i>vu</i>	The least-significant two bits are used to set the vendor-specific bits in the Control byte in the command.

### **modeselect15()**—Issue a Group 0 Mode Select Command

The `modeselect15()` function prepares and issues a group 0 Mode Select command. This command is used to control a variety of standard and vendor-specific device parameters. Typically, `modesense1A0` is first used to retrieve the current parameters. The function prototype is

```
int modeselect15(struct dsreq *dsp, caddr_t data, long datalen,
                int save, int vu);
```

The arguments are as follows:

<i>dsp</i>	The address of a <i>dsreq</i> structure prepared by <b>dsopen()</b> .
<i>data</i>	The address of a mode data page to send.
<i>datalen</i>	The length of the data.
<i>save</i>	The least significant bit sets the SP bit in the command.
<i>vu</i>	The least-significant two bits are used to set the vendor-specific bits in the Control byte in the command.

### **modesense1a()**—Send a Group 0 Mode Sense Command

The **modesense1a()** function prepares and issues a group 0 Mode Sense command to a SCSI device to retrieve a page of device-dependent information. The function prototype:

```
int modesense1a(struct dsreq *dsp, caddr_t data, long datalen,
               int pagectrl, int pagecode, int vu);
```

The arguments are as follows:

<i>dsp</i>	The address of a <i>dsreq</i> structure prepared by <b>dsopen()</b> .
<i>data</i>	The address of a buffer to receive the page of data.
<i>datalen</i>	The length of the buffer.
<i>pagectrl</i>	The least significant 2 bits are set as the PCF bits in the command.
<i>pagecode</i>	The least significant 6 bits are set as the page number.
<i>vu</i>	The least-significant two bits are used to set the vendor-specific bits in the Control byte in the command.

For reference, the PCF codes are as follows:

0	Current values.
1	Changeable values.
2	Default values.
3	Saved values.

For reference, some page numbers are as follows:

- 0 Vendor unique.
- 1 Read/write error recovery.
- 2 Disconnect/reconnect.
- 3 Direct access device format; parallel interface; measurement units.
- 4 Rigid disk geometry; serial interface.
- 5 Flexible disk; printer options.
- 6 Optical memory.
- 7 Verification error.
- 8 Caching.
- 9 Peripheral device.
- 63 (0x3f) Return all pages supported.

#### **read08() and readextended28()—Issue a Read Command**

The **read08()** and **readextended28()** functions prepare and issue particular forms of SCSI Read commands. The Read and extended Read commands have so many variations that it is unlikely that either of these functions will work with your device. However, you can use them as models to build additional variations on Read. Do not preempt the function names.

The function prototypes are

```
int
read08(struct dsreq *dsp, caddr_t data, long datalen,
        long lba, int vu);
int
readextended28(struct dsreq *dsp, caddr_t data, long datalen,
               long lba, int vu);
```

The arguments are as follows:

- dsp* The address of a *dsreq* structure prepared by **dsopen()**.
- data* The address of a buffer to receive the data.
- datalen* The length of the buffer (not exceeding 255 for **read08**)

<i>lba</i>	The logical block address for the start of the read (not exceeding 16 bits for <b>read08</b> )
<i>vu</i>	The least-significant two bits are used to set the vendor-specific bits in the Control byte in the command.

The functions set the transfer length in the command to the number of bytes given by *datalen*. This is often incorrect; many devices want a number of blocks of some size. Function **read08()** sets only 16 bits from *lba* as the logical block number, although the SCSI command format permits another 5 bits to be encoded in the command. For these and other reasons you are likely to need to create customized Read functions of your own.

### **readcapacity25()—Issue a Read Capacity Command**

The **readcapacity25()** function prepares and issues a Read Capacity command to a SCSI device. The function prototype is

```
int
readcapacity25(struct dsreq *dsp, caddr_t data, long datalen,
              long lba, int pmi, int vu);
```

The arguments are as follows:

<i>dsp</i>	The address of a <i>dsreq</i> structure prepared by <b>dsopen()</b> .
<i>data</i>	The address of a buffer to receive the capacity data.
<i>datalen</i>	The length of the buffer, typically 8.
<i>lba</i>	Last block address, 0 unless <i>pmi</i> is nonzero.
<i>pmi</i>	The least-significant bit is used to set the partial medium indicator (PMI) bit of the command.
<i>vu</i>	The least-significant two bits are used to set the vendor-specific bits in the Control byte in the command.

When *pmi* is 0, *lba* should be given as 0 and the command returns the device capacity. When *pmi* is 1, the command returns the last block following block *lba* before which a delay (seek) will occur.

**requestsense03()—Issue a Request Sense Command**

The **requestsense03()** function prepares and issues a Request Sense command. If you include `DSRQ_SENSE` in the *flag* argument to **doscsireq0()**, a Request Sense is sent automatically after an error in a command. The function prototype is

```
int
requestsense03(struct dsreq *dsp, caddr_t data,
               long datalen, int vu);
```

The arguments are:

<i>dsp</i>	The address of a <i>dsreq</i> structure prepared by <b>dsopen0()</b> .
<i>data</i>	The address of a buffer to receive the sense data.
<i>datalen</i>	The length of the buffer.
<i>vu</i>	The least-significant two bits are used to set the vendor-specific bits in the Control byte in the command.

**reserveunit16() and releaseunit17()—Control Logical Units**

The **reserveunit16()** function prepares and issues a Reserve Unit command to reserve a logical unit, causing it to return Reservation Conflict status to requests from other initiators. The **releaseunit17()** function prepares and issues a Release Unit command to release a reserved unit. The function prototypes are

```
int
reservunit16(struct dsreq *dsp, caddr_t data, long datalen,
             int tpr, int tpdid, int extent, int res_id, int vu);
int
releaseunit17(struct dsreq *dsp,
              int tpr, int tpdid, int extent, int res_id, int vu);
```

The arguments are as follows:

<i>dsp</i>	The address of a <i>dsreq</i> structure prepared by <b>dsopen0()</b> .
<i>data</i>	The address of data to send with the Reserve Unit. (This may be NULL for <b>reservunit16()</b> which does not normally transfer data.)
<i>datalen</i>	The length of the data (typically 0).
<i>tpr</i>	The least-significant bit is used to set the Third-Party Reservation bit in the command: 1 means the reservation is on behalf of another initiator.

<i>tpdid</i>	The device ID for the device to hold the reservation: 0 unless <i>tpr</i> is 1.
<i>extent</i>	The least-significant bit sets the least-significant bit of byte 1 of the command string.
<i>res_id</i>	Passed as byte 2 of the command string.
<i>vu</i>	The least-significant two bits are used to set the vendor-specific bits in the Control byte in the command.

### **senddiagnostic1d()**—Issue a Send Diagnostic Command

The **senddiagnostic1d()** function prepares and issues a Send Diagnostic command. The function prototype is

```
int
senddiagnostic1d(struct dsreq *dsp, caddr_t data, long datalen,
                 int self, int dofl, int uofl, int vu);
```

The arguments are as follows:

<i>dsp</i>	The address of a <i>dsreq</i> structure prepared by <b>dsopen()</b> .
<i>data</i>	The address of a page or pages of diagnostic parameter data to be sent.
<i>datalen</i>	The length of the data (0 if none).
<i>self</i>	The least-significant bit sets the Self Test (ST) bit in the command: 1 means return status from the self-test; 0 means hold the results.
<i>dofl</i>	The least-significant bit sets the Device Offline bit of the command: 1 authorizes tests that can change the status of other logical units.
<i>uofl</i>	The least-significant bit sets the Unit Offline bit of the command: 1 authorizes tests that can change the status of the logical unit.
<i>vu</i>	The least-significant two bits are used to set the vendor-specific bits in the Control byte in the command.

When *self* is 1, the status reflects the success of the self-test. You should either set the DSRQ\_SENSE flag in the *dsreq* so that if the self-test fails, a Sense command will be issued, or be prepared to call **requestsense03()**. When *self* is 0, you can use a Read Diagnostic command to return detailed results of the test (however, dslib does not contain a predefined function for Read Diagnostic).

**testunitready00—Issue a Test Unit Ready Command**

The **testunitready00** function prepares and issues a Test Unit Ready command to a SCSI device. The function prototype is

```
int
testunitready00(struct dsreq *dsp);
```

This function is reproduced here in Example 5-2 as an example of how other command-oriented functions can be created.

**Example 5-2** Code of the testunitread00() Function

```
int
testunitready00(struct dsreq *dsp)
{
    fillg0cmd(dsp, CMDBUF(dsp), G0_TEST, 0, 0, 0, 0, 0);
    filldsreq(dsp, 0, 0, DSRQ_READ|DSRQ_SENSE);
    return(doscsireq(getfd(dsp), dsp));
}
```

**write0a() and writeextended2a()—Issue a Write Command**

The **write0a()** function prepares and issues a group 0 Write command. The **writeextended2a()** function prepares and issues an extended (10-byte) Write command. As with Read commands (see “read08() and readextended28()—Issue a Read Command” on page 102), Write commands have many device-specific features, and you will very likely have to create your own customized version of these functions.

The function prototypes are

```
int
write0a(struct dsreq *dsp, caddr_t data, long datalen,
        long lba, int vu);
int
writeextended2a(struct dsreq *dsp, caddr_t data, long datalen,
               long lba, int vu);
```

The arguments are as follows:

<i>dsp</i>	The address of a <i>dsreq</i> structure prepared by <b>dsopen()</b> .
<i>data</i>	The address of the data to be sent.
<i>datalen</i>	The length of the data (at most 255 for <b>write0a</b> )

<i>lba</i>	The logical block address (at most 16 bits for <b>write0a</b> )
<i>vu</i>	The least-significant two bits are used to set the vendor-specific bits in the Control byte in the command.

## Example dslib Program

The program in Example 5-3 illustrates the use of the dslib functions. This is an edited version of a program that can be obtained in full from Dave Olson's home page, <http://reality.sgi.com/employees/olson/Olson/index.html>.

### Example 5-3 Program That Uses dslib Functions

```
#ident "scsicontrol.c: $Revision $"

#include <sys/types.h>
#include <stddef.h>
#include <stdlib.h>
#include <unistd.h>
#include <ctype.h>
#include <errno.h>
#include <fcntl.h>
#include <stdio.h>
#include <string.h>
#include <dslib.h>

typedef struct
{
    uchar  pqt:3; /* peripheral qual type */
    uchar  pdt:5; /* peripheral device type */
    uchar  rmb:1, /* removable media bit */
          dtq:7; /* device type qualifier */
    uchar  iso:2, /* ISO version */
          ecma:3, /* ECMA version */
          ansi:3; /* ANSI version */
    uchar  aenc:1, /* async event notification supported */
          trmiop:1, /* device supports 'terminate io process msg */
          res0:2, /* reserved */
          respfmt:3; /* SCSI 1, CCS, SCSI 2 inq data format */
    uchar  ailen; /* additional inquiry length */
    uchar  res1; /* reserved */
    uchar  res2; /* reserved */
    uchar  reladr:1, /* supports relative addressing (linked cmds) */

```

```
        wide32:1, /* supports 32 bit wide SCSI bus */
        wide16:1, /* supports 16 bit wide SCSI bus */
        synch:1, /* supports synch mode */
        link:1, /* supports linked commands */
        res3:1, /* reserved */
        cmdq:1, /* supports cmd queuing */
        softre:1; /* supports soft reset */
    uchar vid[8]; /* vendor ID */
    uchar pid[16]; /* product ID */
    uchar prl[4]; /* product revision level*/
    uchar vendsp[20]; /* vendor specific; typically firmware info */
    uchar res4[40]; /* reserved for scsi 3, etc. */
    /* more vendor specific information may follow */
} inqdata;

struct msel {
    unsigned char rsv, mtype, vendspec, blkdesclen; /* header */
    unsigned char dens, nblks[3], rsv1, bsize[3]; /* block desc */
    unsigned char pgnum, pglen; /* modesel page num and length */
    unsigned char data[240]; /* some drives get upset if no data requested
    on sense*/
};

#define hex(x) "0123456789ABCDEF" [ (x) & 0xF ]

/* only looks OK if nperline a multiple of 4, but that's OK.
 * value of space must be 0 <= space <= 3;
 */
void
hprint(unsigned char *s, int n, int nperline, int space)
{
    int i, x, startl;

    for(startl=i=0;i<n;i++) {
        x = s[i];
        printf("%c%c", hex(x>>4), hex(x));
        if(space)
            printf("%. *s", ((i%4)==3)+space, " ");
        if ( i%nperline == (nperline - 1) ) {
            putchar('\t');
            while(startl < i) {
                if(isprint(s[startl]))
                    putchar(s[startl]);
                else
                    putchar('.');
            }
        }
    }
}
```

```

        startl++;
    }
    putchar('\n');
}
}
if(space && (i%ncperline))
    putchar('\n');
}

/* aenc, trmiop, reladr, wbus*, synch, linkq, softre are only valid if
 * if respfmt has the value 2 (or possibly larger values for future
 * versions of the SCSI standard). */

static char pdt_types[][16] = {
    "Disk", "Tape", "Printer", "Processor", "WORM", "CD-ROM",
    "Scanner", "Optical", "Jukebox", "Comm", "Unknown"
};
#define NPDT (sizeof pdt_types / sizeof pdt_types[0])

void
printing(struct dsreq *dsp, inqdata *inq, int allinq)
{
    if(DATASENT(dsp) < 1) {
        printf("No inquiry data returned\n");
        return;
    }
    printf("%-10s", pdt_types[(inq->pdt < NPDT) ? inq->pdt : NPDT-1]);
    if (DATASENT(dsp) > 8)
        printf("%12.8s", inq->vid);
    if (DATASENT(dsp) > 16)
        printf("%16s", inq->pid);
    if (DATASENT(dsp) > 32)
        printf("%4s", inq->prl);
    printf("\n");
    if(DATASENT(dsp) > 1)
        printf("ANSI vers %d, ISO ver: %d, ECMA ver: %d; ",
            inq->ansi, inq->iso, inq->ecma);
    if(DATASENT(dsp) > 2) {
        unchar special = *(inq->vid-1);
        if(inq->respfmt >= 2 || special) {
            if(inq->respfmt < 2)
                printf("\nResponse format type %d, but has "
                    "SCSI-2 capability bits set\n", inq->respfmt);
        }
    }
}

```

```
printf("supports: ");
if(inq->aenc)
    printf(" AENC");
if(inq->trmiop)
    printf(" termiop");
if(inq->reladr)
    printf(" reladdr");
if(inq->wide32)
    printf(" 32bit");
if(inq->wide16)
    printf(" 16bit");
if(inq->synch)
    printf(" synch");
if(inq->synch)
    printf(" linkedcmds");
if(inq->cmdq)
    printf(" cmdqueing");
if(inq->softre)
    printf(" softreset");
}
if(inq->respfmt < 2) {
    if(special)
        printf(". ");
    printf("inquiry format is %s",
        inq->respfmt ? "SCSI 1" : "CCS");
}
}
putchar('\n');
printf("Device is ");
/* do test unit ready only if inquiry successful, since many
   devices, such as tapes, return inquiry info, even if
   not ready (i.e., no tape in a tape drive). */
if(testunitready00(dsp) != 0)
    printf("%s\n",
        (RET(dsp)==DSRT_NOSEL) ? "not responding" : "not ready");
else
    printf("ready");
printf("\n");
}

/* inquiry cmd that does vital product data as spec'ed in SCSI2 */
int
vpinquiry12( struct dsreq *dsp, caddr_t data, long datalen, char vu,
            int page)
{
```

```

    fillg0cmd(dsp, (uchar_t *)CMDBUF(dsp), G0_INQU, 1, page, 0, B1(datalen),
        B1(vu<<6));
    filldsreq(dsp, (uchar_t *)data, datalen, DSRQ_READ|DSRQ_SENSE);
    return(doscsireq(getfd(dsp), dsp));
}

int
startunit1b(struct dsreq *dsp, int startstop, int vu)
{
    fillg0cmd(dsp, (uchar_t *)CMDBUF(dsp), 0x1b, 0, 0, 0, (uchar_t)startstop, B1(vu<<6));
    filldsreq(dsp, NULL, 0, DSRQ_READ|DSRQ_SENSE);
    dsp->ds_time = 1000 * 90; /* 90 seconds */
    return(doscsireq(getfd(dsp), dsp));
}

int
myinquiry12(struct dsreq *dsp, uchar_t *data, long datalen, int vu, int neg)
{
    fillg0cmd(dsp, (uchar_t *)CMDBUF(dsp), G0_INQU, 0, 0, 0, B1(datalen), B1(vu<<6));
    filldsreq(dsp, data, datalen, DSRQ_READ|DSRQ_SENSE|neg);
    dsp->ds_time = 1000 * 30; /* 90 seconds */
    return(doscsireq(getfd(dsp), dsp));
}

int
dsreset(struct dsreq *dsp)
{
    return ioctl(getfd(dsp), DS_RESET, dsp);
}

void
usage(char *prog)
{
    fprintf(stderr,
        "Usage: %s [-i (inquiry)] [-e (exclusive)] [-s (sync) | -a (async)]\n"
        "\t[-l (long inq)] [-v (vital proddata)] [-r (reset)] [-D (diagseltest)]\n"
        "\t[-H (halt/stop)] [-b blksize]\n"
        "\t[-g (get host flags)] [-d (debug)] [-q (quiet)] scsidevice [...]\n",
        prog);
    exit(1);
}

main(int argc, char **argv)
{

```

```
struct dsreq *dsp;
char *fn;
/* int because they must be word aligned. */
int errs = 0, c;
int vital=0, doreset=0, exclusive=0, dosync=0;
int dostart = 0, dostop = 0, dosenddiag = 0;
int doing = 0, printname = 1;
unsigned bsize = 0;
extern char *optarg;
extern int optind, opterr;

opterr = 0; /* handle errors ourselves. */
while ((c = getopt(argc, argv, "b:HDSaserdvlqCiq")) != -1)
switch(c) {
case 'i':
    doing = 1; /* do inquiry */
    break;
case 'D':
    dosenddiag = 1;
    break;
case 'r':
    doreset = 1; /* do a scsi bus reset */
    break;
case 'e':
    exclusive = O_EXCL;
    break;
case 'd':
    dsdebug++; /* enable debug info */
    break;
case 'q':
    printname = 0; /* print devicename only if error */
    break;
case 'v':
    vital = 1; /* set evpd bit for scsi 2 vital product data */
    break;
case 'H':
    dostop = 1; /* send a stop (Halt) command */
    break;
case 'S':
    dostart = 1; /* send a startunit/spinup command */
    break;
case 's':
    dosync = DSRQ_SYNXFR; /* attempt to negotiate sync scsi */
    break;
case 'a':
```

```
        dosync = DSRQ_ASYNXFR; /* attempt to negotiate async scsi */
        break;
default:
    usage(argv[0]);
}

if(optind >= argc || optind == 1) /* need at 1 arg and one option */
    usage(argv[0]);

while (optind < argc) { /* loop over each filename */
    fn = argv[optind++];
    if(printname) printf("%s: ", fn);
    if((dsp = dsopen(fn, O_RDONLY|exclusive)) == NULL) {
        /* if open fails, try pre-pending /dev/scsi */
        char buf[256];
        strcpy(buf, "/dev/scsi/");
        if((strlen(buf) + strlen(fn)) < sizeof(buf)) {
            strcat(buf, fn);
            dsp = dsopen(buf, O_RDONLY|exclusive);
        }
        if(!dsp) {
            if(!printname) printf("%s: ", fn);
            fflush(stdout);
            perror("cannot open");
            errs++;
            continue;
        }
    }
}

/* try to order for reasonableness; reset first in case
 * hung, then inquiry, etc. */

if(doreset) {
    if(dsreset(dsp) != 0) {
        if(!printname) printf("%s: ", fn);
        printf("reset failed: %s\n", strerror(errno));
        errs++;
    }
}

if(doingq) {
    int inqbuf[sizeof(inqdata)/sizeof(int)];
    if(myinquiry12(dsp, (uchar_t *)inqbuf, sizeof inqbuf, 0, dosync)) {
        if(!printname) printf("%s: ", fn);
        printf("inquiry failure\n");
    }
}
```

```
        errs++;
    }
    else
        printing(dsp, (inqdata *)inqbuf, 0);
}

if(vital) {
    unsigned char *vpinq;
    int vpingbuf[sizeof(inqdata)/sizeof(int)];
    int vpingbuf0[sizeof(inqdata)/sizeof(int)];
    int i, serial = 0, asciidef = 0;
    if(vpinquiry12(dsp, (char *)vpingbuf0,
        sizeof(vpingbuf)-1, 0, 0)) {
        if(!printname) printf("%s: ", fn);
        printf("inquiry (vital data) failure\n");
        errs++;
        continue;
    }
    if(DATASENT(dsp) <4) {
        printf("vital data inquiry OK, but says no"
            "pages supported (page 0)\n");
        continue;
    }
    vpinq = (unsigned char *)vpingbuf0;
    printf("Supported vital product pages: ");
    for(i = vping[3]+3; i>3; i--) {
        if(vpinq[i] == 0x80)
            serial = 1;
        if(vpinq[i] == 0x82)
            asciidef = 1;
        printf("%2x ", vpinq[i]);
    }
    printf("\n");
    vpinq = (unsigned char *)vpingbuf;
    if(serial) {
        if(vpinquiry12(dsp, (char *)vpingbuf,
            sizeof(vpingbuf)-1, 0, 0x80) != 0) {
            if(!printname) printf("%s: ", fn);
            printf("inquiry (serial #) failure\n");
            errs++;
        }
        else if(DATASENT(dsp)>3) {
            printf("Serial #: ");
            fflush(stdout);
            /* use write, because there may well be
```

```
        *nulls; don't bother to strip them out */
        write(1, vpinq+4, vpinq[3]);
        printf("\n");
    }
}
if(asciiodef) {
    if(vpinquiry12(dsp, (char *)vpinqbuf,
        sizeof(vpinqbuf)-1, 0, 0x82) != 0) {
        if(!printname) printf("%s: ", fn);
        printf("inquiry (ascii definition) failure\n");
        errs++;
    }
    else if(DATASENT(dsp)>3) {
        printf("Ascii definition: ");
        fflush(stdout);
        /* use write, because there may well be
        *nulls; don't bother to strip them out */
        write(1, vpinq+4, vpinq[3]);
        printf("\n");
    }
}
}

if(dostop && startunitlb(dsp, 0, 0)) {
    if(!printname) printf("%s: ", fn);
    printf("stopunit fails\n");
    errs++;
}

if(dostart && startunitlb(dsp, 1, 0)) {
    if(!printname) printf("%s: ", fn);
    printf("startunit fails\n");
    errs++;
}

if(dosenddiag && senddiagnostic1d(dsp, NULL, 0, 1, 0, 0, 0)) {
    if(!printname) printf("%s: ", fn);
    printf("self test fails\n");
    errs++;
}
}
dsclose(dsp);
}
return(errs);
}
```



---

## Control of External Interrupts

Some Silicon Graphics computer systems can generate and receive external interrupt signals. These are simple, two-state signal lines that cause an interrupt in the receiving system.

The external interrupt hardware is managed by a kernel-level device driver that is distributed with IRIX and automatically configured when the system supports external interrupts. The driver provides two abilities to user-level processes:

- the ability to change the state of an outgoing interrupt line, so as to interrupt the system to which the line is connected)
- the ability to capture an incoming interrupt signal with low latency

**Note:** Some software for external interrupt support is closely tied to the hardware of the system. The features described in this chapter are hardware-dependent and in many cases cannot be ported from one system type to another without making software changes. There is no external interrupt support in the O2 workstation.

There is a hardware-independent way to capture incoming external interrupts: the user-level interrupt facility (ULI). To learn about ULI, see Chapter 7, “User-Level Interrupts,” especially “Registering an External Interrupt Handler” on page 131.

## External Interrupts in Challenge and Onyx Systems

The hardware architecture of the Challenge series supports external interrupt signals as follows:

- Four jacks for outgoing signals are available on the master IO4 board. A user-level program can change the level of these lines individually.
- Two jacks for incoming interrupt signals are also provided. The input lines are combined with logical OR and presented as a single interrupt; a program cannot distinguish one input line from another.

The electrical interface to the external interrupt lines is documented at the end of the ei(7) reference page.

A program controls the outgoing signals by interacting with the external interrupt device driver. This driver is associated with the device special file */dev/ei*, and is documented in the ei(7) reference page.

### Generating Outgoing Signals

A program can generate an outgoing signal on any one of the four external interrupt lines. To do so, it first opens */dev/ei*. Then it can apply **ioctl()** on the file descriptor to switch the outgoing lines. The principal functions are summarized in Table 6-1.

**Table 6-1** Functions for Outgoing External Signals

Operation	Typical ioctl() Call
Set pulse width to <i>N</i> microseconds.	<code>ioctl(eifd, EIIOCSETOPW, N)</code>
Return current output pulse width.	<code>ioctl(eifd, EIIOCGETOPW, &amp;var)</code>
Send a pulse on some lines <i>M</i> . <sup>a</sup>	<code>ioctl(eifd, EIIOCSTROBE, M)</code>
Set a high (active, asserted) level on lines <i>M</i> .	<code>ioctl(eifd, EIIOCSETHI, M)</code>
Set a low (inactive, deasserted) level on lines <i>M</i> .	<code>ioctl(eifd, EIIOCSETLO, M)</code>

a. *M* is an unsigned integer whose bits 0, 1, 2, and 3 correspond to the external interrupt lines 0, 1, 2, and 3. At least one bit must be set.

In the Challenge and Onyx series, the level on an outgoing external interrupt line is set directly from a CPU. The device driver generates a pulse (function EIIOCSTROBE) by asserting the line, then spinning in a disabled loop until the specified pulse time has elapsed, and finally deasserting the line. Clearly, if the pulse width is set to much more than the default of 5 microseconds, pulse generation could interfere with the handling of other interrupts.

The calls to assert and deassert the outgoing lines (functions EIIOCSETHI and EIIOCSETLO) do not tie up the kernel. However, direct assertion of the outgoing signal should be used only when the desired signal frequency and pulse duration are measured in milliseconds or seconds. No user-level program can hope to generate pulse durations measured in microseconds by calling these functions. For one thing, the minimum guaranteed interrupt service time is 200 microseconds. An interrupt occurring between the call to assert the signal and the call to deassert it will stretch the intended pulse width by at least 200 microseconds.

## Receiving Incoming External Interrupts

An important feature of the Challenge and Onyx external input line is that interrupts are triggered by the level of the signal, not by the transition from deasserted to asserted. This means that, whenever external interrupts are enabled and any of the input lines are in the asserted state, an external interrupt occurs. The interface between your program and the external interrupt device driver is affected by this hardware design. The functions for incoming signals are summarized in Table 6-2.

**Table 6-2** Functions for Incoming External Interrupts

Operation	Typical ioctl() Call
Enable receipt of external interrupts.	<code>ioctl(eifd, EIIOCENABLE) eicinit();</code>
Disable receipt of external interrupts.	<code>ioctl(eifd, EIIOCDISABLE)</code>
Block in the driver until an interrupt occurs.	<code>ioctl(eifd, EIIOCRECV, &amp;eiargs)</code>
Request a signal when an interrupt occurs.	<code>ioctl(eifd, EIIOCSTSIG, <i>signumber</i>)</code>
Wait in an enabled loop for an interrupt.	<code>eicbusywait(1)</code>
Set expected pulse width of incoming signal.	<code>ioctl(eifd, EIIOCSETIPW, <i>microsec</i>)</code>

**Table 6-2 (continued)** Functions for Incoming External Interrupts

Operation	Typical ioctl() Call
Set expected time between incoming signals.	<code>ioctl(eifd, EIIOCSETSPW, <i>microsec</i>)</code>
Return current expected time values.	<code>ioctl(eifd, EIIOCGETIPW, &amp;<i>var</i>)</code> <code>ioctl(eifd, EIIOCGETSPW, &amp;<i>var</i>)</code>

**Detecting Invalid External Interrupts**

The external interrupt handler maintains two important numbers:

- the expected input pulse duration in microseconds
- the minimum pulse-to-pulse interval, called the “stuck” pulse width because it is used to detect when an input line is “stuck” in the asserted state

When the external interrupt device driver is entered to handle an interrupt, it waits with interrupts disabled until time equal to the expected input pulse duration has passed since the interrupt occurred. The default pulse duration is 5 microseconds, and it typically takes longer than this to recognize and process the interrupt, so no time is wasted in the usual case. However, if a long expected pulse duration is set, the interrupt handler might have to waste some cycles waiting for the end of the pulse.

At the end of the expected pulse duration, the interrupt handler counts one external interrupt and returns to the kernel, which enables interrupts and returns to the interrupted process.

Normally the input line is deasserted within the expected duration. However, if the input line is still asserted when the time expires, another external interrupt occurs immediately. The external interrupt handler notes that it has been reentered within the “stuck” pulse time since the last interrupt. It assumes that this is still the same input pulse as before. In order to prevent the stuck pulse from saturating the CPU with interrupts, the interrupt handler disables the external interrupt signal.

External interrupts remain disabled for one timer tick (10 milliseconds). Then the device driver reenables external interrupts. If an interrupt occurs immediately, the input line is still asserted. The handler disables external interrupts for another, longer delay. It continues to delay and to test the input signal in this manner until it finds the signal deasserted.

### Setting the Expected Pulse Width

You can set the expected input pulse width and the minimum pulse-to-pulse time using `ioctl()`. For example, you could set the expected pulse width using a function like the one shown in Example 6-1.

**Example 6-1**     Function to Test and Set External Interrupt Pulse Width

```
int setEIPulseWidth(int eifd, int newWidth)
{
    int oldWidth;
    if ( (0==ioctl(eifd, EIIOCGETIPW, &oldWidth))
        && (0==ioctl(eifd, EIIOCSETIPW, newWidth)) )
        return oldWidth;
    perror("setEIPulseWidth");
    return 0;
}
```

The function retrieves the original pulse width and returns it. If either `ioctl()` call fails, it returns 0.

The default pulse width is 5 microseconds. Pulse widths shorter than 4 microseconds are not recommended.

Since the interrupt handler keeps interrupts disabled for the duration of the expected width, you want to specify as short an expected width as possible. However, it is also important that all legitimate input pulses terminate within the expected time. When a pulse persists past the expected time, the interrupt handler is likely to detect a “stuck” pulse, and disable external interrupts for several milliseconds.

Set the expected pulse width to the duration of the longest valid pulse. It is not necessary to set the expected width longer than the longest valid pulse. A few microseconds are spent just reaching the external interrupt handler, which provides a small margin for error.

### Setting the Stuck Pulse Width

You can set the minimum pulse-to-pulse width using code like that in Example 6-1, using constants `EIIOCGETSPW` and `EIIOCSETSPW`.

The default stuck-pulse time is 500 microseconds. Set this time to the nominal pulse-to-pulse interval, minus the largest amount of “jitter” that you anticipate in the signal. In the event that external signals are not produced by a regular oscillator, set this value to the expected pulse width plus the duration of the shortest expected “off” time, with a minimum of twice the expected pulse width.

For example, suppose you expect the input signal to be a 10 microsecond pulse at 1000 Hz, both numbers plus or minus 10%. Set the expected pulse width to 10 microseconds to ensure that all pulses are seen to complete. Set the stuck pulse width to 900 microseconds, so as to permit a legitimate pulse to arrive 10% early.

### Receiving Interrupts

The external interrupt device driver offers you four different methods of receiving notification of an interrupt. You can

- have a signal of your choice delivered to your process
- test for interrupt-received using either an `ioctl()` call or a library function
- sleep until an interrupt arrives or a specified time expires
- spin-loop until an interrupt arrives

You would use a signal when interrupts are infrequent and irregular, and when it is not important to know the precise arrival time. Signal latency can be milliseconds long in some extreme cases. For this reason, it would not be wise to use signals to handle a high rate of interrupts, nor to expect to time interrupts closely. Use a signal when, for example, the external interrupt represents a human-operated switch or some kind of out-of-range alarm condition.

The `ioctl(EIIOCRCV)` call tests for an interrupt, or suspends the caller until an interrupt arrives or a timeout expires (see the `ei(7)` reference page for details). Use this method when interrupts arrive frequently enough that it is worthwhile devoting a process to handling them.

The **ioctl()** call is a fairly costly method of polling, since it entails entry to and exit from the kernel. This is not significant if the polling is infrequent—for example, if one poll call is made every 60th of a second.

When the **ioctl()** call is used to wait for an interrupt, an unknown amount of time can pass between the moment when the interrupt handler unblocks the process and the moment when the kernel dispatches the process. This makes it impossible to timestamp the interrupt at the microsecond level.

In order to detect an incoming interrupt with minimum latency, use the library function **eicbusywait()** (see the ei(7) reference page). This function does not switch into kernel mode, so it is a very fast method of polling for an interrupt. However, if you ask it to wait until an interrupt occurs, it waits by spinning on a repeated test for an interrupt. This monopolizes the CPU, so this form of waiting can be used reliably only by a process running in an isolated CPU. (If there are other processes to run, or interrupts to handle, the polling loop in **eicbusywait()** shares the CPU and can be preempted for long periods.)

The benefit of **eicbusywait()** is that, in an isolated, nonpreemptive CPU, control returns to the calling process in negligible time after the interrupt handler detects the interrupt, so the interrupt can be handled quickly and timed precisely.



---

## User-Level Interrupts

The user-level interrupt (ULI) facility complements the ability to perform programmed I/O (PIO) from user space. You can use PIO to initiate a device action that leads to a device interrupt, and you can intercept and handle the interrupt in your program. For function prototypes and other details, see the `uli(3)` reference page.

**Note:** In IRIX 6.3 for O2, ULI support applies to interrupts from the VME bus and to external interrupts only. Neither source of interrupts is available in the O2 workstation, but this chapter is included in case you are developing software for other systems.

### Overview of ULI

In the past, PIO could only be synchronous: the program wrote to a device register, then polled the device until the operation was complete. With ULI, the program can manage a device that causes interrupts on the VME bus. You set up a handler function within your program. The handler is called whenever the device causes an interrupt.

In IRIX 6.3 for O2, user-level interrupts are supported for VME-bus devices, and for external interrupts on the Challenge and Onyx systems.

When using ULI with a VME device, you use VME PIO to initiate device actions and to transfer data to and from device registers (see “VME Programmed I/O” on page 64). When using ULI to trap external interrupts, you enable the interrupts with `ioctl()` calls to the external interrupt handler (see Chapter 6, “Control of External Interrupts”).

## The User Level Interrupt Handler

The ULI handler is a function within your program. It is entered asynchronously from the IRIX kernel's interrupt-handling code. The kernel transfers from the kernel address space into the user process address space, and makes the call in user (not privileged kernel) execution mode. Despite this more complicated linkage, you can think of the ULI handler as a subroutine of the kernel's interrupt handler. As such, the performance of the ULI handler has a direct bearing on the system's interrupt response time.

Like the kernel's interrupt handler, the ULI handler can be entered at almost any time, regardless of what code is being executed by the CPU—a process of your program or a process of another program, executing in user space or in a system function. In fact, the ULI handler can be entered from one CPU while the your program executes concurrently in another CPU. Your normal code and your ULI function can execute in true concurrency, accessing the same global variables.

## Restrictions on the ULI Handler

Because the ULI handler is called in a special context of the kernel's interrupt handler, it is severely restricted in the system facilities it can use. The list of features the ULI handler may not use includes the following:

- Any use of floating-point calculations. The kernel does not take time to save floating-point registers during an interrupt trap. The floating-point coprocessor is turned off and an attempt to use it in the ULI handler causes a SIGILL (illegal instruction) exception.
- Any use of IRIX system functions. Because most of the IRIX kernel runs with interrupts enabled, the ULI handler could be entered while a system function was already in progress. System functions do not support reentrant calls. In addition, many system functions can sleep, which an interrupt handler may not do.
- Any storage reference that causes a page fault. The kernel cannot suspend the ULI handler for page I/O. Reference to an unmapped page causes a SIGSEGV (memory fault) exception.
- Any calls to C library functions that might violate the preceding restrictions.

There are very few library functions that you can be sure will use no floating point and make no system calls. Unfortunately, library functions such as **sprintf()**, often used in debugging, must be avoided.

In essence, the ULI handler should only do such things as

- store data in program variables to record the interrupt event
  - A ring buffer is a data structure that is suitable for concurrent access.
- program the device as required to clear the interrupt or acknowledge it
  - The ULI handler has access to the whole program address space, including any mapped-in devices, so it can perform PIO loads and stores.
- post a semaphore to wake up the main process
  - This must be done using a ULI function.

## Planning for Concurrency

Since the ULI handler can interrupt the program at any point, or run concurrently with it, the program must be prepared for concurrent execution. There are two areas to consider: global variables, and library routines.

### Debugging With Interrupts

The asynchronous, possibly concurrent entry to the ULI handler can confuse a debugging monitor such as *dbx*. Some strategies for dealing with this are covered in the `uli(3)` reference page.

### Declaring Global Variables

When variables can be modified by both the main process and the ULI handler, you must take special care to avoid race conditions.

An important step is to specify `-D_SGI_REENTRANT_FUNCTIONS` to the compiler, so as to get the reentrant versions of the C library functions. This ensures that, if the main process and the ULI handler both enter the C library, there will be no collision over global variables.

You can declare the global variables that are shared with the ULI handler with the keyword “volatile,” so that the compiler will generate code to load the variables from memory on each reference. However, the compiler never holds global values in registers over a function call, and you almost always have a function call (such as `ULI_block_intr()`) preceding a test of a shared global variable.

## Using Multiple Devices

The ULI feature allows a program to open more than one interrupting device. You register a handler for each device. However, the program can only wait for a specific interrupt to occur; that is, the `ULI_sleep()` function specifies the handle of one particular ULI handler. This does not mean that the main program must sleep until that particular interrupt handler is entered, however. Any ULI handler can waken the main program, as discussed under “Interacting With the Handler” on page 131.

## Setting Up

A program initializes for ULI in the following major steps:

1. Open the device special file for the device.
2. For a VME device, map the device addresses into process memory (see “Mapping a VME Device Into Process Address Space” on page 64).
3. Lock the program address space in memory.
4. Initialize any data structures used by the interrupt handler.
5. Register the interrupt handler.
6. Interact with the device and the interrupt handler.

Any time after the handler has been registered, an interrupt can occur, causing entry to the ULI handler.

## Opening the Device Special File

Devices are represented by device special files (see “Device Special Files” on page 34). In order to gain access to a device, you open the device special file that represents it.

The file that represents the external interrupt lines on a Challenge or Onyx system is `/dev/ei`. It can be opened by more than one process at a time, under rules that are spelled out in the `ei(7)` reference page. The methods of opening `/dev/ei`, configuring pulse widths, and generating output pulses are covered in Chapter 6, “Control of External Interrupts,” and remain the same.

The files that represent VME control units are */dev/vme/vme\**. The rules for opening one of these files and mapping a device into memory are covered under “VME Programmed I/O” on page 64, and remain the same. The difference is that, with ULI, you can map in a device that can cause interrupts.

The program should open the device and (for a VME device) verify that the device exists and is active before proceeding.

## Locking the Program Address Space

The ULI handler must not reference a page of program text or data that is not present in memory. You prevent this by locking the pages of the program address space in memory. The simplest way to do this is to call the **plock()** system function:

```
if (plock(PROCLOCK))
{ perror("plock"); exit(); }
```

The **plock()** function has two possible difficulties. One is that the calling process must have superuser privilege (see the **plock(2)** reference page). This may not pose a problem if the program needs superuser privilege in any case, for example in order to open a device special file. The second is that it locks all text and data pages. In a very large program this may be so much memory that system performance is harmed.

The **mpin()** function can be used by unprivileged programs to lock a limited number of pages. The limit is set by the tunable system parameter *maxlkmem*. (Check its value—typically 2000—in */var/sysgen/mtune/kernel*. See the **systune(1)** reference page for instructions on changing a tunable parameter.)

In order to use **mpin()**, you must specify the exact address ranges to be locked. Provided that the ULI handler refers only to global data and its own code, it is relatively simple to derive address ranges that encompass the needed pages. If the ULI handler calls any library functions, the library DSO needs to be locked as well. The smaller the scope of the ULI handler, the easier it is to use **mpin()**.

## Registering the Interrupt Handler

When the program is ready to start operations, it registers its ULI handler. The ULI handler is a function that matches the prototype

```
void function_name(void *arg);
```

The registration function takes arguments with the following purposes:

- the file descriptor of the device special file
- four arguments related to the ULI handler:
  - the address of the handler function
  - an argument value to be passed to the handler on each interrupt—typically a pointer to a work area that is unique to the interrupting device, supposing the program is using more than one device
  - the size, and an optional address, of memory to be used as stack space when calling the handler
- a count of semaphores to be allocated for use with this interrupt

You can ask the ULI support to allocate a stack space by passing a null pointer for the stack argument. When the ULI handler is as simple a function as it normally is, the default stack size of 1024 bytes is ample.

The semaphores are allocated and maintained by the ULI support. They are used to coordinate between the program process and the interrupt handler, as discussed under “Interacting With the Handler” on page 131. You should specify one semaphore for each independent process that can wait for interrupts from this handler. Normally one semaphore is sufficient.

The returned value is a handle that is used to identify this interrupt in other functions. Once registered, the ULI handler remains registered until the program terminates (there is no function for un-registration).

## Registering an External Interrupt Handler

The `ULI_register_ei()` function takes the arguments described in the preceding topic. Once it has successfully registered your handler, all external interrupts are directed to that handler.

It is important to realize that, so long as a ULI handler is registered, none of the other interrupt-reporting features supported by the external interrupt device driver (see Chapter 6, “Control of External Interrupts” and the `ei(7)` reference page) operate any more. These restrictions include the facts that:

- The per-process external interrupt queues are not updated.
- Signals requested by `ioctl(EIIOCSETSIG)` are not sent.
- Calls to `ioctl(EIIOCRCV)` sleep until they are interrupted by a timeout, a signal, or because the program using ULI terminated and an interrupt arrived.
- Calls to the library function `eicbusywait()` do not terminate.

Clearly you should not use ULI for external interrupts when there are other programs running that also use them.

## Registering a VME Interrupt Handler

The `ULI_register_vme()` function takes two additional arguments:

- the interrupt level that the device uses.
- a word that contains, or receives, an interrupt vector number

The interrupt level used by a device is normally set by hardware and documented in the VECTOR line that defines the device (see “Learning VME Device Addresses” on page 64).

Some VME devices have a fixed interrupt vector number; others are programmable. You pass a fixed vector number to the function. If the number is programmable, you pass 0, and the function allocates a number. You must then use PIO to program the vector number into the device.

## Interacting With the Handler

The program process and the ULI handler synchronize their actions using two functions.

When the program cannot proceed without an interrupt, it calls **ULI\_sleep()**, specifying

- the handle of the interrupt for which to wait
- the number of the semaphore to use for waiting

Typically only one process ever calls **ULI\_sleep()** and it specifies waiting on semaphore 0. However, it is possible to have two or more processes that wait. For example, if the device can produce two distinct kinds of interrupts—normal and high-priority, perhaps—you could set up an independent process for each interrupt type. One would sleep on semaphore 0, the other on semaphore 1.

When an ULI handler is entered, it wakes up a program process by calling **ULI\_wakeup()**, specifying the semaphore number to be posted. The handler must know which semaphore to post, based on the values it can read from the device or from program variables.

The **ULI\_sleep()** call can terminate early, for example if a signal is sent to the process. The process that calls **ULI\_sleep()** must test to find the reason the call returned—it is not necessarily because of an interrupt.

The **ULI\_wakeup()** function can be called from normal code as well as from a ULI handler function. It could be used within any type of asynchronous callback function to wake up the program process.

The **ULI\_wakeup()** call also specifies the handle of the interrupt. When you have multiple interrupting devices, you have the following design choices:

- You can have one child process waiting on the handler for each device. In this case, each ULI handler specifies its own handle to **ULI\_wakeup()**.
- You can have a single process that waits on any interrupt. In this case, the main program specifies the handle of one particular interrupt to **ULI\_sleep()**, and every ULI handler specifies that same handle to **ULI\_wakeup()**.

### **Achieving Mutual Exclusion**

The program can gain exclusive use of global variables with a call to **ULI\_block\_intr()**. This function does not block receipt of the hardware interrupt, but does block the call to the ULI handler. Until the program process calls **ULI\_unblock\_intr()**, it can test and update global variables without danger of a race condition. This period of time should be as short as possible, because it extends the interrupt latency time. If more than one

hardware interrupt occurs while the ULI handler is blocked, it will be called for only the last-received interrupt.

There are other techniques for safe handling of shared global variables besides blocking interrupts. One important, and little-known, set of tools is the **test\_and\_set()** group of functions documented in the test\_and\_set(3) reference page. These instructions use the Load Linked and Store Conditional instructions of the MIPS instruction set to safely update global variables in various ways.

## Sample Program

The program listed in Example 7-1 is a hypothetical example of how user-level interrupts can be used to handle external interrupts in a Challenge and Onyx system.

### Example 7-1 Hypothetical ULI Program

```

/* This program demonstrates use of the External Interrupt source
 * to drive a User Level Interrupt.
 *
 * The program requires the presence of an external interrupt cable looped
 * back between output number 0 and one of the inputs on the machine on
 * which the program is run.
 */
#include <sys/ei.h>
#include <sys/uli.h>
#include <sys/lock.h>
#include <unistd.h>
#include <stdlib.h>
#include <stdio.h>
#include <fcntl.h>
/* The external interrupt device file is used to access the EI hardware */
#define EIDEV "/dev/ei"
static int eifd;
/* The user level interrupt id. This is returned by the ULI registration
 * routine and is used thereafter to refer to that instance of ULI
 */
static void *ULIid;
/* Variables which are shared between the main process thread and the ULI
 * thread may have to be declared as volatile in some situations. For
 * example, if this program were modified to wait for an interrupt with
 * an empty while() statement, e.g.
 *     while(!intr);

```

```
* the value of intr would be loaded on the first pass and if intr is
* false, the while loop will continue forever since only the register
* value, which never changes, is being examined. Declaring the variable
* intr as volatile causes it to be reloaded from memory on each iteration.
* In this code however, the volatile declaration is not necessary since
* the while() loop contains a function call, e.g.
*   while(!intr)
*       ULI_sleep(ULIid, 0);
* The function call forces the variable intr to be reloaded from memory
* since the compiler cannot determine if the function modified the value
* of intr. Thus the volatile declaration is not necessary in this case.
* When in doubt, declare your globals as volatile.
*/
static int intr;
/* This is the actual interrupt service routine. It runs
* asynchronously with respect to the remainder of this program, possibly
* simultaneously, on an MP machine. This function must obey the ULI mode
* restrictions, meaning that it may not use floating point or make
* any system calls. (Try doing so and see what happens.) Also, this
* function should be written to execute as quickly as possible, since it
* runs at interrupt level with lower priority interrupts masked.
* The system imposes a 1-second time limit on this function to prevent
* the cpu from freezing if an infinite loop is inadvertently programmed
* in. Try inserting an infinite loop to see what happens.
*/
static void
intrfunc(void *arg)
{
    /* Set the global flag indicating to the main thread that an
    * interrupt has occurred, and wake it up
    */
    intr = 1;
    ULI_wakeup(ULIid, 0);
}
/* This function creates a new process and from it, generates a
* periodic external interrupt.
*/
static void
signaler(void)
{
    int pid;
    if ((pid = fork()) < 0) {
        perror("fork");
        exit(1);
    }
}
```

```
if (pid == 0) {
    while(1) {
        if (ioctl(eifd, EIIOCSTROBE, 1) < 0) {
            perror("EIIOCSTROBE");
            exit(1);
        }
        sleep(1);
    }
}
/* The main routine sets everything up, then sleeps waiting for the
 * interrupt to wake it up.
 */
int
main()
{
    /* open the external interrupt device */
    if ((eifd = open(EIDEV, O_RDONLY)) < 0) {
        perror(EIDEV);
        exit(1);
    }
    /* Set the target cpu to which the external interrupt will be
     * directed. This is the cpu on which the ULI handler function above
     * will be called. Note that this is entirely optional, but if
     * you do set the interrupt cpu, it must be done before the
     * registration call below. Once a ULI is registered, it is illegal
     * to modify the target cpu for the external interrupt.
     */
    if (ioctl(eifd, EIIOCSETINTRCPU, 1) < 0) {
        perror("EIIOCSETINTRCPU");
        exit(1);
    }
    /* Lock the process image into memory. Any text or data accessed
     * by the ULI handler function must be pinned into memory since
     * the ULI handler cannot sleep waiting for paging from secondary
     * storage. This must be done before the first time the ULI handler
     * is called. In the case of this program, that means before the
     * first EIIOCSTROBE is done to generate the interrupt, but in
     * general it is a good idea to do this before ULI registration
     * since with some devices an interrupt may occur at any time
     * once registration is complete
     */
    if (plock(PROCLOCK) < 0) {
        perror("plock");
        exit(1);
    }
}
```

```
    }
    /* Register the external interrupt as a ULI source. */
    ULId = ULI_register_ei( eifd, /* the external interrupt device */
                          intrfunc, /* the handler function pointer */
                          0, /* the argument to the handler */
                          1, /* the number of semaphores needed */
                          NULL, /* the stack to use (supply one) */
                          0); /* the stack size to use (default) */

    if (ULId == 0) {
        perror("register ei");
        exit(1);
    }
    /* Enable the external interrupt. */
    if (ioctl(eifd, EIIOCENABLE) < 0) {
        perror("EIIOCENABLE");
        exit(1);
    }
    /* Start creating incoming interrupts. */
    signaler();
    /* Wait for the incoming interrupts and report them. Continue
     * until the program is terminated by ^C or kill.
     */
    while (1) {
        intr = 0;
        while(!intr) {
            if (ULI_sleep(ULId, 0) < 0) {
                perror("ULI_sleep");
                exit(1);
            }
        }
        printf("sleeper woke up\n");
    }
}
```

## PART THREE

# Kernel-Level Drivers

### **Chapter 8:** Structure of a Kernel-Level Driver

The software structure of a block or character device driver: the entry points it provides for kernel use, and how it communicates with user-level processes

### **Chapter 9:** Device Driver/Kernel Interface

A topical survey of the facilities the IRIX kernel provides to device drivers.

### **Chapter 10:** Building and Installing a Driver

How a kernel-level driver is compiled, loaded, and linked with the IRIX kernel.

### **Chapter 11:** Testing and Debugging a Driver

How a kernel-level driver is tested and debugged using *symmon* and other facilities.

### **Chapter 12:** Driver Example

Annotated code of a simple memory-mapping device driver.



---

## Structure of a Kernel-Level Driver

A kernel-level device driver consists of a module of subroutines that supply services to the kernel. The subroutines are public entry points in the driver. When events occur, the kernel calls these entry points. The driver takes action and returns a result code.

This chapter discusses when the driver entry points are called, what parameters they receive, and what actions they are expected to take. For a conceptual overview of the kernel and drivers, see “Kernel-Level Device Control” on page 47. For details on how a driver is compiled, linked, and added to IRIX, see Chapter 10, “Building and Installing a Driver.”

**Note:** This chapter discusses device drivers. The entry point conventions for STREAMS drivers are covered in Chapter 16, “STREAMS Drivers.” Additional entry points supported only for PCI drivers are covered in Chapter 15, “PCI Device Drivers.”

The primary topics covered in this chapter are:

- “Summary of Driver Structure” on page 140 summarizes the entry points and how they are made known to the kernel.
- “Driver Flag Constant” on page 145 documents a public constant the driver must supply.
- “Initialization Entry Points” on page 147 documents the entry points that are called at boot time and when a loadable driver is loaded.
- “Open and Close Entry Points” on page 150 documents the entry points called by the **open()** and **close()** kernel functions.
- “Control Entry Point” on page 154 documents the entry point called by the **ioctl()** kernel function.
- “Data Transfer Entry Points” on page 155 documents the entry points called by the **read()** and **write()** kernel functions.
- “Poll Entry Point” on page 158 documents the entry point called by the **poll()** kernel function.

- “Memory Map Entry Points” on page 161 documents the entry points called by the **mmap()** kernel function.
- “Interrupt Entry Point” on page 167 documents the entry point called to handle a device interrupt.
- “Support Entry Points” on page 170 documents the entry points that support kernel operations and system administration.
- “Planning for Multiprocessor Use” on page 174 points out methods for writing a multiprocessor-aware driver.

## Summary of Driver Structure

A driver consists of a binary object module in ELF format stored in the */var/sysgen/boot* directory. As a program, the driver consists of a set of functional entry points that supply services to the IRIX kernel. There is a large set of entry points to cover different situations, but no single driver supports all possible entry points.

The entry points that a driver supports must be named according to a specified convention. The *lboot* command uses entry point names to build tables used by the kernel.

### Entry Point Naming and lboot

The device driver makes known which entry points it supports by giving them public names in its object module. The *lboot* command links together the object modules of drivers and other kernel modules to make a bootable kernel. *lboot* recognizes the entry points by the form of their names.

### Driver Name Prefix

A device driver must be described by a file in the */var/sysgen/master.d* directory (see “Master Configuration Database” on page 40). One of the items in that configuration file specifies the driver *prefix*, a string of 1 to 14 characters that is unique to that driver. For example, the prefix of the SCSI driver is *scsi\_*.

The prefix string is defined in the */var/sysgen/master.d* file only. The string does not have to appear as a constant in the driver, and the name of the driver object file does not have to correspond to the prefix (although the object module typically has a related name).

The *lboot* command recognizes driver entry points by searching the driver object module for public names that begin with the prefix string. For example, the entry point for the **open()** operation must have a name that consists of the prefix string followed by the letters “open.”

In this book, entry point names are written as follows: *pfxopen*, where *pfx* stands for the driver’s prefix string.

### Kernel Switch Tables

The IRIX kernel maintains tables that allow it to dispatch calls to device drivers quickly. These tables are built by *lboot* based on the device major numbers and the names of the driver entry points. The tables are named as follows:

<i>bdevsw</i>	Table of block device drivers
<i>cdevsw</i>	Table of character device drivers
<i>fmodsw</i>	Table of STREAMS drivers
<i>vfssw</i>	Table of filesystem modules (not related to device drivers)

The tables for block and character drivers have one row for each major device number, and one column for each possible driver entry point. As *lboot* loads a driver, it fills in that driver’s row of a switch table with the addresses of the driver’s entry points. Where an entry point is not defined, *lboot* leaves the address of a null routine that returns the ENODEV error code.

The sizes of the switch tables are fixed at boot time in order to minimize kernel data space. The table sizes are tunable parameters that can be set with *sysstune* (see the *sysstune(1)* reference page).

When a driver is loaded dynamically (see “Configuring a Loadable Driver” on page 239), the associated row of the switch table is not filled at link time but rather is filled when the driver is loaded. When you add new, loadable drivers, you might need to specify a larger switch table. The *IRIX Administration: System Configuration and Operation* book documents these tunable parameters.

## Entry Point Summary

The names of all possible driver entry points and their purposes are summarized in Table 8-1. STREAMS drivers are covered in Chapter 16.

**Table 8-1** Entry Points in Alphabetic Order

Entry Point	Purpose	Discussion	Reference Page
<i>pfxattach</i>	PCI device attach entry point.	page 386	
<i>pfxclose</i>	Note the device is not in use.	page 153	close(D3)
<i>pfxdevflag</i>	Constant flag bits for driver features.	page 145	devflag(D1)
<i>pfxdetach</i>	PCI device detach entry point.	page 401	
<i>pfxedtinit</i>	Initialize driver from VECTOR statement.	page 148	edtinit(D2)
<i>pfxhalt</i>	Prepare for system shutdown.	page 171	halt(D2)
<i>pfxinit</i>	Initialize driver at load or boot time.	page 148	init(D2)
<i>pfxintr</i>	Handle device interrupt (not used).	page 167	intr(D2)
<i>pfxiocctl</i>	Implement control operations.	page 154	iocctl(D2)
<i>pfxmap</i>	Implement memory-mapping (IRIX).	page 163	map(D2)
<i>pfxmmap</i>	Implement memory-mapping (SVR4).	page 165	mmap(D2)
<i>pfxopen</i>	Connect a process to a device. Connect a stream module.	page 150 page 501	open(D2)
<i>pfxpoll</i>	Implement device event test.	page 160	poll(D2)
<i>pfxprint</i>	Display diagnostic about block device.	page 172	print(D2)
<i>pfxread</i>	Implement device input.	page 155	read(D2)
<i>pfxrput</i>	STREAMS message on read queue.	page 502	put(D2)
<i>pfxsize</i>	Return logical size of block device.	page 172	size(D2)
<i>pfxsrv</i>	STREAMS service queued messages.	page 503	srv(D2)
<i>pfxstart</i>	Initialize driver at load or boot time.	page 149	start(D2)
<i>pfxstrategy</i>	Input/output for a block device.	page 157	strategy(D2)

**Table 8-1 (continued)** Entry Points in Alphabetic Order

Entry Point	Purpose	Discussion	Reference Page
<i>pxunload</i>	Prepare loadable module for unloading.	page 170	unload(D2)
<i>pxunmap</i>	Note the end of a memory mapping.	page 166	unmap(D2)
<i>pxwput</i>	STREAMS message on write queue.	page 502	put(D2)
<i>pxwrite</i>	Implement device output.	page 155	write(D2)

The use of entry points in different types of drivers is summarized in Table 8-2. The columns of Table 8-2 show the different types of drivers. The table cells show whether a given entry point is optional (O), required (R), or not allowed (N).

**Table 8-2** Use of Driver Entry Points

Entry Point	Character	Block	Pseudo	STREAMS
<i>pxattach</i>	R (PCI)	R (PCI)	N	N
<i>pxclose</i>	R	R	R	R
<i>pxdetach</i>	R (PCI)	R (PCI)	N	N
<i>pxdevflag</i>	O	O	O	O
<i>pxedtinit</i>	O	O	N	N
<i>pxhalt</i>	O	O	O	O
<i>pxinit</i>	O	O	O	O
<i>pxintr</i>	O	O	N	N
<i>pxioctl</i>	O	O	O	N
<i>pxmap</i>	O	O	O	N
<i>pxmmap</i>	O	O	O	N
<i>pxopen</i>	R	R	R	R
<i>pxpoll</i>	O	O	N	O
<i>pxprint</i>	N	O	N	N
<i>pxread</i>	O	N	O	N

**Table 8-2 (continued)** Use of Driver Entry Points

Entry Point	Character	Block	Pseudo	STREAMS
<i>pxrput</i>	N	N	N	R
<i>pxsize</i>	N	R	N	N
<i>pxsrv</i>	N	N	N	R
<i>pxstart</i>	O	O	O	O
<i>pxstrategy</i>	N	R	N	N
<i>pxunload</i>	O	O	O	O
<i>pxunmap</i>	O	O	O	N
<i>pxwput</i>	N	N	N	R
<i>pxwrite</i>	O	N	O	N

As can be seen from Table 8-2, no driver supports all entry points.

- A minimal driver for a character device supports *pxinit()*, *pxopen()*, *pxread()*, *pxwrite()*, and *pxclose()*. The *pxioctl()* and *pxpoll()* entry points are optional. (The *pxattach()* and *pxdetach()* entry points are also required for a PCI device.)
- A minimal pseudo-device driver supports *pxstart()*, *pxopen()*, *pxmap()*, *pxunmap()*, and *pxclose()* (the latter two as stubs).
- A minimal block device driver supports *pxedtinit()*, *pxopen()*, *pxsize()*, *pxstrategy()*, and *pxclose()*. (The *pxattach()* and *pxdetach()* entry points are also required for a PCI device.)

## Driver Flag Constant

Any device driver or STREAMS module should define a public name *pfxddevflag* as a static integer. This integer contains a bitmask with zero or more of the following flags, which are declared in *sys/conf.h*:

D_MP	The driver is prepared for multiprocessor systems.
D_WBACK	The driver handles its own cache-writeback operations.
D_MT	The driver is prepared for a multithreaded kernel.
D_OLD	The driver implements IRIX 4.x semantics.

The flag names are declared in the header file *sys/ddi.h*. A typical definition would resemble the following:

```
int testdrive_devflag = D_MP;
```

A STREAMS module should also provide this flag, but the only relevant bit value for a STREAMS driver is D\_MP (see “Driver Flag Constant” on page 500).

The flag value is saved in the kernel switch table with the driver’s entry points (see “Kernel Switch Tables” on page 141).

When a driver does not define a *pfxddevflag*, *lboot* saves a word containing D\_OLD by default. See the note regarding D\_OLD on page 147.

### Flag D\_MP

You specify D\_MP in *pfxddevflag* to tell *lboot* that your driver is designed to operate in a multiprocessor system. The top half of the driver is designed to cope with multiple concurrent entries in multiple CPUs. The top and bottom halves synchronize through the use of semaphores or locks and do not rely on interrupt masking for critical sections. These issues are discussed further under “Planning for Multiprocessor Use” on page 174.

When D\_MP is not present in *pfxddevflag*, IRIX ensures that the driver code, including the upper-half entry points and the interrupt handler, executes only on CPU 0 of a multiprocessor. This ensures behavior similar to a uniprocessor, but can cause a performance bottleneck when either the device or CPU 0 is heavily used.

**Note:** This flag is only tested when loading a character driver or STREAMS driver. There is no special handling for block drivers and network drivers in multiprocessors. Block and network drivers *must be* multiprocessor-aware.

### Flag D\_WBACK

You specify D\_WBACK in *pxdevflag* to tell *lboot* that a block driver performs any necessary cache write-back operations through explicit calls to **dki\_dcache\_wb()** and related functions (see the *dki\_dcache\_wb(D3)* reference page).

When D\_WBACK is not present in *pxdevflag*, the **physiock()** function ensures that all cached data related to *buf\_t* structures is written back to main memory before it enters the driver's strategy routine. (See the *physiock(D3)* reference page and "Entry Point strategy()" on page 157.)

### Flag D\_MT

This flag is defined in IRIX 6.2 but has no effect in that release. The next major release of IRIX will run driver interrupt routines as threads of control within the kernel address space. D\_MT indicates that this driver understands that it can be run as one or more cooperating threads, and uses kernel synchronization primitives to serialize access to driver common data structures.

### Flag D\_OLD

The D\_OLD flag exists only to retain compatibility with certain drivers written originally for IRIX 4.x. It changes two features of the kernel-to-driver interface:

- The first argument to the *pxopen()* entry is a *dev\_t* value instead of the pointer-to-*dev\_t* that is now standard.
- The driver sets its return code by storing it into a global, *u.u\_error*, instead of returning it as the result of the function call.

D\_OLD is incompatible with D\_MP.

When a driver has no *pxdevflag* constant, *lboot* assumes it is a D\_OLD driver.

**Note:** The D\_OLD flag value, and the incompatible interface that it implies, is supported for compatibility only. Support for this flag, and support for drivers that use it (or that have no *pxdevflag* constant), *will be withdrawn* in the release of IRIX after 6.2. It may be removed from certain platform-specific releases of IRIX 6.2. Silicon Graphics urges you to revise any driver that depends on D\_OLD to use current semantics for parameters and return codes.

## Initialization Entry Points

The kernel calls a driver to initialize itself at any of three different entry points, as follows:

<i>pxinit</i>	Initialize self-defining hardware or a pseudo-device.
<i>pxedtinit</i>	Initialize a hardware device based on VECTOR data.
<i>pxstart</i>	General initialization.

Each call has different abilities. A driver may define any combination of the three entry points. It is not uncommon to define both a *pxstart()* and one of *pxedtinit()* or *pxinit()*.

### When Initialization Is Performed

The initialization entry points of ordinary (nonloadable) drivers are called during system startup, after interrupts have been enabled and before the message “The system is coming up” is displayed. In all cases, interrupts are enabled and basic kernel services are available at this time. However, other loadable or optional kernel modules might not have been initialized, depending on the sequence of statements in the files in */var/sysgen/system*.

Whenever a driver is initialized, the entry points are called in the following sequence:

1. *pxinit()* is called first.
2. *pxedtinit()* is called once for each VECTOR statement in reverse order of the VECTOR statements found in */var/sysgen/system* files.
3. *pxstart()* is called last.

### Initialization of Loadable Drivers

A loadable driver (see “Loadable Drivers” on page 59) is initialized any time it is loaded. This can occur more than once, if the driver is loaded, unloaded, and reloaded. When a loadable driver is configured for autoregister, it is loaded with other drivers during system startup. (For more information on autoregister, see “Configuring a Loadable Driver” on page 239.) Such a driver is initialized at system startup time along with the nonloadable drivers.

### Entry Point `init()`

The `pxinit()` entry point is called once during system startup or when a loadable driver is loaded. It receives no input arguments; its prototype is simply:

```
void pxinit(void);
```

You can use this entry point to initialize a hardware device that is self-defining; that is, all the information the driver needs is either coded into the driver, or can be gotten by probing the device itself. You can also use `pxinit()` to initialize a pseudo-device driver; that is, a driver that does not have real hardware attached.

A driver that is brought into the system by a `USE` or `INCLUDE` line in a system configuration file (see “Configuring a Kernel” on page 238) typically initializes in the `pxinit()` entry point.

### Entry Point `edtinit()`

The `pxedtinit()` entry is designed to initialize devices that are configured using the `VECTOR` statement in the system configuration file (see “System Configuration Files” on page 41). The entry point name is a contraction of “early device table initialization.”

The `VECTOR` statement specifies hardware details about a device on the VME, GIO, or EISA bus (on systems that have one of those buses), including iospace addresses, interrupt level, and an integer parameter. The `VECTOR` statement can specify a “probe” parameter that lets the kernel test for the existence of the specified hardware.

When the kernel processes a `VECTOR` statement during bootstrap and the probe is successful (or no probe is specified), the kernel stores the `VECTOR` parameters in a structure of type `edt_t`. (This structure is declared in `sys/edt.h`.)

Each time the kernel loads a driver that is named in a VECTOR statement, the kernel calls the driver's `pfxedtinit()` entry one time for each VECTOR statement that named that driver and had a successful probe (or that had no probe). VECTOR statements are processed in reverse sequence to the order in which they are coded in `/var/sysgen/system` files.

The prototype of the `pfxedtinit()` entry is

```
void pfxedtinit(edt_t *e);
```

The `edt_t` contains at least the following fields (see the `system(4)` reference page for the corresponding VECTOR parameters):

`e_bus_type` Integer specifying the bus type; constant values are declared in `sys/edt.h`, for example `ADAP_VME`, `ADAP_GIO`, or `ADAP_EISA`.

`e_adap` Integer specifying the adapter (bus) number.

`e_ctlr` Value from the VECTOR `ctlr=` parameter; typically the device minor number.

`e_space` Array of up to three I/O space structures of type `iospace_t`.

The difference between `pfxininit()` and `pfxedtinit()` is that `pfxedtinit()` is parameterized with information from the VECTOR line, and is called once for each VECTOR line that is associated with real hardware.

A driver that uses `pfxedtinit()` needs to save the `edt_t` information in a data structure. If the driver supports multiple devices—that is, if it can be called for multiple VECTOR statements—it needs to allocate an array or chain of structures, and save new data on each entry.

### Entry Point `start()`

The `pfxstart()` entry point is called at system startup, and whenever a loadable driver is loaded. It is called after `pfxedtinit()` and `pfxininit()`, but before any other entry point such as `pfxopen()`. The `pfxstart()` entry point receives no arguments; its prototype is simply

```
void pfxstart(void);
```

The `pfxstart()` entry point is a suitable place to allocate a poll-head structure using `phalloc()`, as discussed in “Use and Operation of `poll(2)`” on page 159.

## Open and Close Entry Points

The `pxopen()` and `pxclose()` entries for block and character devices are called when a device comes into use and when use of it is finished. For a conceptual overview of the `open()` process, see “Overview of Device Open” on page 49.

### Entry Point `open()`

The kernel calls a device driver’s `pxopen()` entry when a process executes the `open()` system call on any device special file (see the `open(2)` reference page). It is also called when a process executes the `mount()` system call on a block device (see the `mount(2)` reference page). ( For the `pxopen()` entry point of a STREAMS driver, see “Entry Point `open()`” on page 501.)

The prototype of `pxopen()` is as follows:

```
int pxopen(dev_t *devp, int oflag, int otyp, cred_t *crp);
```

The argument values are

- |                    |   |
|--------------------|---|
| <code>*devp</code> | Pointer to a <code>dev_t</code> value from which you can extract both the major and minor device numbers.   |
| <code>otyp</code>  | An integer flag specifying the source of the call: a user process opening a character device or block device, or another driver.                              |
| <code>oflag</code> | Flag bits specifying user mode options on the <code>open()</code> call.   |
| <code>crp</code>   | A <code>cred_t</code> object—an opaque structure for use in authentication. Standard access privileges to the special device file have already been verified. |

**Note:** When the driver’s `pxdevflag` entry contains `D_OLD` or when `pxdevflag` is not defined, the first argument to `pxopen()` is a `dev_t` value, not a pointer to a `dev_t` value. See “Flag `D_OLD`” on page 146.

The `open(2)` reference page discusses the kind of work the `pxopen()` entry point can do. In general, the driver is expected to verify that this user process is permitted access in the way specified in `otyp` (reading, writing, or both) for the device specified in `*devp`. If access is not allowable, the driver returns a nonzero error code from `sys/errno.h`, for example `ENOMEM` or `EBUSY`.

### Use of the Device Number

When the driver supports a single device with no logical unit divisions, the device number is of little interest except for diagnostic displays. When the driver supports multiple devices, or a device with multiple logical units, the minor device number is the key to locating the device information. The device number can also encode device options, as discussed under “Minor Device Number” on page 37.

When the driver supports the *pfxedtinit()* entry, the driver needs a way to associate the different *edt\_t* structures passed to *pfxedtinit()* with the device numbers passed to *pfxopen()* and other routines. One solution is to require that the *ctrl=* value from the VECTOR statement—which is passed in the *e\_ctrl* field of *edt\_t*—must be the same as the device minor number.

### Use of the Open Type

The *otyp* flag distinguishes between the following possible sources of this call to *pfxopen()* (the constants are defined in *sys/open.h*).

- a call to open a character device (OTYP\_CHR)
- a call to open a block device (OTYP\_BLK)
- a call to mount a block device as a filesystem (OTYP\_MNT)
- a call to open a block device as swapping device (OTYP\_SWP)
- a call direct from a device driver at a higher level (OTYP\_LYR)

Typically a driver is written only to be a character driver or a block driver, and can be called only through the switch table for that type of device. When this is the case, the *otyp* value has little use.

It is possible to have the same driver treated as both block and character, in which case the driver needs to know whether the **open()** call addressed a block or character special device. It is possible for a block device to support different partitions with different uses, in which case the driver might need to record the fact that a device has been mounted, or opened as a swap device.

With all open types except OTYP\_LYR, *pfxopen()* is called for every open or mount operation, but *pfxclose()* is called only when the last close or unmount occurs. The OTYP\_LYR feature is used almost exclusively by drivers distributed with IRIX, like the

host adapter SCSI driver (see “Host Adapter Concepts” on page 303). For each open of this type, there is one call to *pxclose()*.

### Use of the Open Flag

The interpretation of the open mode flags is up to the designer of the driver. Four modes can be requested (declared in *sys/file.h*):

- FREAD           Input access wanted.
- FWRITE          Output access wanted (both FREAD and FWRITE may be set, corresponding to O\_RDWR mode).
- FNDELAY or     Return at once, do not sleep if the open cannot be done immediately.  
FNONBLOCK
- FEXCL           Request exclusive use of the device.

You decide which of the flags have meaning with respect to the abilities of this device. You can return an EINVAL error when an unsupported mode is requested.

A key decision is whether the device can be opened only by one process at a time, or by multiple processes. If multiple opens are supported, a process can still request exclusive access with the FEXCL mode.

When the device can be used by only one process, or when FEXCL access is supported, the driver must keep track of the fact that the device is open. When the device is busy, the driver can test the FNDELAY and FNONBLOCK flags; if either is set, it can return EBUSY. Otherwise, the driver should sleep until the device is free; this requires coordination with the *pxclose()* entry point.

### Use of the cred\_t Object

The *cred\_t* object passed to *pxopen()*, *pxclose()*, and *pxioctl()* can be used with the *drv\_priv()* function to find out if the effective calling user ID is privileged or not (see the *drv\_priv(D3)* reference page). Do not examine the object in detail, since its contents are subject to change from release to release.

### Saving the Size of a Block Device

In a block device driver, the *pxsize()* entry point will be called soon after *pxopen()* (see “Entry Point size()” on page 172). It is typically best to calculate or read the device capacity at open time, and save it to be reported from *pxsize()*.

### Saving the User ABI

If your driver is, or might be, compiled to the 64-bit model for use with a 64-bit IRIX kernel, and if it supports the *pxioctl()* or *pxpoll()* entry points, the driver should test and save the user process’s programming model during an open. For details, see “Handling 32-Bit and 64-Bit Execution Models” on page 173.

### Entry Point close()

The kernel calls the *pxclose()* entry when the last process calls *close()* or *umount()* for the device special file. It is important to know that when the device can be opened by multiple processes, *pxclose()* is not called for every *close()* function, but only when the last remaining process closes the device and no other processes have it open.

The function prototype and arguments of *pxclose()* are

```
int pxclose(dev_t dev, int flag, int otyp, cred_t *crp);
```

The arguments are the same as were passed to *pxopen()*. However, the flag argument is not necessarily the same as at any particular call to *open()*.

It is up to you to design the meaning of “close” for this type of device. The close(D2) reference page discusses some of the actions the driver can do. Some considerations are:

- If the device is opened and closed frequently, you may decide to retain dynamic data structures.
- If the device can perform an action such as “rewind” or “eject,” you decide whether that action should be done upon close. Possibly the choice of acting or not acting can be set by an *ioctl()* call; or possibly the choice can be encoded into the device minor number—for example, the no-rewind-on-close option is encoded in certain tape minor device numbers.
- If the *pxopen()* entry point supports exclusive access, and it can be waiting for the device to be free, *pxclose()* must release the wait.

The `pxclose()` entry can detect an error and report it with a return code. However, the file is closed or unmounted regardless.

## Control Entry Point

The `pxioctl()` entry point is called by the kernel when a user process executes the `ioctl()` system call (see the `ioctl(2)` reference page). This entry point is allowed in character drivers only. Block device drivers do not support it, and STREAMS drivers pass control information as messages.

For an overview of the relationship between the user process, kernel, and the control entry point, see “Overview of Device Control” on page 50.

The prototype of the entry point is

```
int pxioctl(dev_t dev, int cmd, void *arg,
            int mode, cred_t *crp, int *rvalp);
```

The argument values are

- dev*        A *dev\_t* value from which you can extract the major and minor device numbers.
- cmd*        The request value specified in the `ioctl()` call.
- arg*        The optional argument value specified in the `ioctl()` call, or NULL if none was specified.
- mode*       Flag bits specifying the `open()` mode, as associated with the file descriptor passed to the `ioctl()` system function.
- crp*        A *cred\_t* object—an opaque structure for use in authentication, describing the process that is in-context. Standard access privileges to the special device file have already been verified.
- \*rvalp*     The integer result to be returned to the user process.

It is up to the device driver to interpret the *cmd* and *arg* values in the light of the *mode* and other arguments. When the *arg* value is a pointer to data in the process address space, the driver uses the `copyin()` kernel function to copy the data into kernel space, and the `copyout()` function to return updated values. (See the `copyin(D3)` and `copyout(D3)` reference pages, and also “Transferring Data” on page 194.)

## Choosing the Command Numbers

The command numbers supported by *pxioctl()* are arbitrary; but the recommended practice is to make sure that they are different from those of any other driver. One method to achieve this is suggested in the *ioctl(D2)* reference page.

## Supporting 32-Bit and 64-Bit Callers

The *ioctl()* entry point may need to interpret a structure prepared in the user process. In a 64-bit system, the user process can be either a 32-bit or a 64-bit program. For discussion of this issue, see “Handling 32-Bit and 64-Bit Execution Models” on page 173

## User Return Value

The kernel returns 0 to the *ioctl()* system function unless the *pxioctl()* function returns an error code. In the event of an error, the kernel returns the code the driver places in *\*rvalp*, if any, or -1. To ensure that the user process sees a specific error code, set the code in *\*rvalp*, and return that value.

## Data Transfer Entry Points

The *pxread()* and *pxwrite()* entry points are supported by character device drivers and pseudo-device drivers that allow reading and writing. They are called by the kernel when the user process calls the *read()*, *readv()*, *write()*, or *writv()* system function.

The *pxstrategy()* entry point is required of block device drivers. It is called by the kernel when either a filesystem or the paging subsystem needs to transfer a block of data.

## Entry Points *read()* and *write()*

The *pxread()* and *pxwrite()* entry points are similar to each other—only the direction of data transfer differs. The prototypes of the functions are

```
int pxread (dev_t dev, uio_t *uiop, cred_t *crp);
int pxwrite(dev_t dev, uio_t *uiop, cred_t *crp);
```

The arguments are

- dev*            A *dev\_t* value from which you can extract both the major and minor device numbers.
- \*uiop*          A *uiop\_t* object—a structure that defines the user’s buffer memory areas.
- crp*            A *cred\_t* object—an opaque structure for use in authentication. Standard access privileges to the special device file have already been verified.

### Data Transfer for a PIO Device

A character device driver using PIO transfers data in the following steps:

1. If there is a possibility of a timeout, start a timeout delay (see “Waiting for Time to Pass” on page 216).
2. Initiate the device operation as required.
3. Transfer data between the device and the buffer represented by the *uiop\_t* (see “Transferring Data Through a *uiop\_t* Object” on page 197).
4. If it is necessary to wait for an interrupt, put the process to sleep (see “Waiting and Mutual Exclusion” on page 206).
5. When data transfer is complete, or when an error occurs, clear any pending timeout and return the final status of the operation. If the return code is 0, the final state of the *uiop\_t* determines the byte count returned by the **read()** or **write()** call.

### Calling Entry Point **strategy()** From Entry Point **read()** or **write()**

A device driver that supports both character and block interfaces must have a **pxstrategy()** routine in which it performs the actual I/O. For example, the Silicon Graphics disk drivers support both character and block driver interfaces, and perform all I/O operations in the **pxstrategy()** function. However, the **pxread()**, **pxwrite()** and **pxioctl()** entries supported for character-type access also need to perform I/O operations. They do this by calling the **pxstrategy()** routine indirectly, using the kernel function **physiock()** or **uiophysio()** (see the **physiock(D3)** and **uiophysio(D3)** reference pages, and see “Waiting for Block I/O to Complete” on page 219).

Both the **physiock()** and **uiophysio()** functions takes care of the housekeeping needed to interface to the **pxstrategy()** entry, including the work of allocating a buffer and a *buf\_t* structure, locking buffer pages in memory and waiting for I/O completion. Both routines require the *uio\_t* to describe only a single segment of data (*uio\_iovcnt* of 1). Although they are very similar, the two functions differ in the following ways:

- **physiock()** returns EINVAL if the initial offset is not a multiple of 512 bytes. If this is a requirement of your **pxstrategy()** routine, use **physiock()**; if not, use **uiophysio()**.
- **physiock()** is compatible with SVR4, while **uiophysio()** is unique to IRIX.

Example 8-1 shows the skeleton of a hypothetical driver in which the **pxread()** entry does its work through the **pxstrategy()** entry.

**Example 8-1** Hypothetical **pxread()** entry in a Character/Block Driver

```
hypo_read (dev_t dev, uio_t *uiop, cred_t *crp)
{
    // ...validate the operation... //
    return physiock(hypo_strategy, /* our strategy entry */
        0, /* allocate temp buffer & buf_t */
        dev, /* dev_t arg for strategy */
        B_READ, /* direction flag for buf_t */
        uiop);
}
```

The **pxwrite()** entry would be identical except for passing B\_WRITE instead of B\_READ.

This dual-entry strategy is required only in a driver that supports both character and block access.

### Entry Point strategy()

A block device driver does not directly support system calls by user processes. Instead, it provides services to a filesystem such as XFS, or to the memory paging subsystem of IRIX. These subsystems call the **pxstrategy()** entry point to read data in whole blocks.

Calls to `pxstrategy()` are not directly related in time to system functions called by a user process. For example, a filesystem may buffer many blocks of data in memory, so that the user process may execute dozens or hundreds of `write()` calls without causing an entry to the device driver. When the user function closes the file or calls `fsync()`—or when for unrelated reasons the filesystem needs to free some buffers—the filesystem calls `pxstrategy()` to write numerous blocks of data.

In a driver that supports the character interface as well, the `pxstrategy()` entry can be called indirectly from the `pxread()`, `pxwrite()` and `pxioctl()` entries, as described under “Calling Entry Point `strategy()` From Entry Point `read()` or `write()`” on page 156.

The prototype of the `pxstrategy()` entry point is

```
int pxstrategy(struct buf *bp);
```

The argument is the address of a `buf_t` structure, which gives the strategy routine the information it needs to perform the I/O:

- The `dev_t` containing major and minor device numbers
- The direction of the transfer (read or write)
- The location of the buffer in kernel memory
- The amount of data to transfer
- The starting block number on the device

For more on the contents of the `buf_t` structure, see “Structure `buf_t`” on page 185 and the `buf(D4)` reference page.

The driver uses the information in the `buf_t` to validate the data transfer and programs the device to start the transfer. Then it stores the address of the `buf_t` where the interrupt handler can find it (see “Interrupt Entry Point” on page 167) and calls `biowait()` to wait for the operation to complete. For the next step, see “Completing Block I/O” on page 169 (see also the `biowait(D3)` reference page).

## Poll Entry Point

The `pxpoll()` entry point is called by the kernel when a user process calls the `poll()` or `select()` system function asking for status on a character special device. To implement it, you need to understand the IRIX implementation of `poll()`.

## Use and Operation of `poll(2)`

The IRIX version of `poll()` allows a process to wait for events of different types to occur on any combination of devices, files, and STREAMS (see the `poll(2)` and `select(2)` reference pages). It is possible for multiple processes to be waiting for events on the same device.

It is up to you as the designer of a driver to decide which of the events that are documented in `poll(2)` are meaningful for your device. Other requested events simply never happen to the device.

Much of the complexity of `poll()` is handled by the IRIX kernel, but the kernel requires the assistance of any device driver that supports `poll()`. The driver is expected to allocate and hold a `pollhead` structure (declared in `sys/poll.h`) for each minor device that it supports. Allocation is simple; the driver merely calls the `phalloc()` kernel function. (The `pxstart()` entry point is a suitable place for this call; see “Entry Point `start()`” on page 149.)

There are two phases to the operation of `poll()`. When the system function is called, the kernel calls the `pxpoll()` entry point to find out if any requested events are pending at this time. If the kernel finds any event's pending (on this or any other polled object), the `poll()` function returns to the user process. Nothing further is required.

However, when no requested event has happened, the user process expects the `poll()` function to block until an event has occurred. The kernel cannot implement this delay by repeatedly testing for events; that would be too inefficient. The kernel must rely on device drivers to notify it when an event has occurred.

### Use of `pollwakeup()`

A device driver that supports `pxpoll()` is required to notify the kernel whenever an event that the driver supports has occurred. The driver does this by calling a kernel function, `pollwakeup()`, passing the `pollhead` structure for the affected device, and bit flags for the events that have taken place. In the event that one or more user processes are blocked in a `poll()`, waiting for an event from this device, the call to `pollwakeup()` will release the sleeping processes. For an example, see “Calling `pollwakeup()`” on page 169.

### Use of pollwakeupp() Without Interrupts

If the device in question does not support interrupts, the driver cannot support **poll()** unless it can somehow get control to discover an event and report it to **pollwakeupp()**. One possibility is that the driver could simulate interrupts by setting a succession of **itimeout()** delays. On each timeout the driver would test its device for a change of status, call **pollwakeupp()** when an event has occurred; and schedule a new delay. (See “Waiting for Time to Pass” on page 216.)

### Entry Point poll()

The prototype for **pxpoll()** is as follows:

```
int pxpoll(dev_t dev, short events, int anyyet,
           short *reventsp, struct pollhead **phpp);
```

The argument values are

<i>dev</i>	A <i>dev_t</i> value from which you can extract the major and minor device numbers.
<i>events</i>	Bit-flags for the events the user process is testing, as passed to <b>poll()</b> and declared in <i>sys/poll.h</i> .
<i>reventsp</i>	A field to receive the bit-flags of events that have occurred, or to receive 0x0000 if no requested events have occurred..
<i>anyyet</i> and <i>phpp</i>	When <i>anyyet</i> is zero and no events have occurred, the kernel requires the address of the pollhead structure for this minor device to be returned in <i>phpp</i> .

Example 8-2 shows the **pxpoll()** code of a hypothetical device driver. Only three event tests are supported: POLLIN and POLLRDNORM (treated as equivalent) and POLLOUT. The device driver maintains an array of *pollhead* structures, one for each supported minor device. These are presumably allocated during initialization.

#### Example 8-2 pxpoll() Code for Hypothetical Driver

```
struct pollhead phds[MAXMINORS];
#define OUR_EVENTS (POLLIN|POLLOUT|POLLRDNORM)
hypo_poll(dev_t dev, short events, int anyyet,
           short *reventsp, struct pollhead **phpp)
```

```

{
    minor_t dminor = getemisor(dev);
    short happened = 0;
    short wanted = events & OUR_EVENTS;
    if (wanted & (POLLIN|POLLRDNORM))
    {
        if (device_has_data_ready(dminor))
            happened |= (POLLIN|POLLRDNORM);
    }
    if (wanted & POLLOUT)
    {
        if (device_ready_for_output(dminor))
            happened |= POLLOUT;
    }
    if (device_pending_error(dminor))
        happened |= POLLERR;
    if (0 == (*reventsp = happened))
    {
        if (!anyyet) *phpp = phds[dminor]
    }
    return 0;
}

```

The code in Example 8-2 begins by discarding any unsupported event flags that might have been requested. Then it tests the remaining flags against the device status. If the device has an uncleared error, the code inserts the POLLERR event. If no events were detected, and if the kernel requested it, the address of the *pollhead* structure for this minor device is returned.

## Memory Map Entry Points

A user process requests memory mapping by calling the system function **mmap()**. When the mapped object is a character device special file, the kernel calls the *pfxmmap()* or *pfxmap()* entry to validate and complete the mapping. To understand these entry points, you must understand the **mmap()** system function.

## Concepts and Use of `mmap()`

The purpose of the `mmap()` system function (see the `mmap(2)` reference page) is to make the contents of a file directly accessible as part of the virtual address space of the user process. The results depend on the kind of file that is mapped:

- When the mapped object is a normal file, the process can load and store data from the file as if it were an array in memory.
- When the mapped object is a character device special file, the process can load and store data from device registers as if they were memory variables.
- When the mapped object is a block of memory owned and prepared by a pseudo-device driver, the process gains access to some special piece of memory data that it would not normally be able to access.

In all cases, access is gained through normal load and store instructions, without the overhead of calling system functions such as `read()`. Furthermore, the same mapping can be executed by other processes, in which case the same memory, or file, or device is shared by multiple, concurrent processes. This is how shared memory segments are achieved.

### Use of `mmap()`

The `mmap()` system function takes four key parameters:

- the file descriptor for an open file, which can be either a normal disk file or a device special file
- an offset within that file at which the mapped data is to start. For a normal file, this is a file offset; for a device file, it represents an address in the address space of the device or the bus
- the length of data to be mapped
- protection flags, showing whether the mapped data is read-only or read-write

When the mapped object is a normal file, the filesystem implements the mapping. The filesystem does not call the block device driver for assistance in mapping a file. It does call the block device driver `pfstrategy()` entry to read and write blocks of file data as necessary, but the mapping of pages of data into pages of memory is controlled in the filesystem code.

When the mapped object is a device special file, the **mmap()** parameters are passed to the device driver at either its **pxmmap()** or **pxmap()** entry point. The device driver interprets the parameters in the context of the device, and uses a kernel function to create the mapping.

### Persistent Mappings

Once a device or kernel memory has been mapped into some user address space, the mapping persists until the user process terminates or calls **unmap()** (see the `unmap(2)` reference page). In particular, the mapping does not end simply because the device special file is closed. You cannot assume, in the **pxfclose()** or **pxunload()** entry points, that all mappings to devices have ended.

### Entry Point `map()`

The **pxmap()** entry point can be defined in either a character or a block driver (it is the only mapping entry point that a block driver can supply). The function prototype is

```
int pxmap(dev_t dev, vhandle_t *vt,  
          off_t off, int len, int prot);
```

The argument values are

<i>dev</i>	A <i>dev_t</i> value from which you can extract both the major and minor device numbers.
<i>vt</i>	The address of an opaque structure that describes the assigned address in the user process address space. The structure contents are subject to change.
<i>off, len</i>	The offset and length arguments passed to <b>mmap()</b> by the user process.
<i>prot</i>	Flags showing the access intentions of the user process.

The first task of the driver is to verify that the access specified in *prot* is allowed. The next task is to validate the *off* and *len* values: do they fall in the valid address space of the device?

When the device driver approves of a mapping, it uses a kernel function, **v\_mapphys()**, to establish the mapping. This function (documented in the `v_mapphys(3)` reference page) takes the *vhandle\_t*, an address in kernel cached or uncached memory, and a length.

It makes the specified region of kernel space a part of the address space of the user process.

For example, a pseudo-device driver that intends to share kernel virtual memory with user processes would first allocate the memory:

```
caddr_t *kaddr = kmem_alloc (len , KM_CACHEALIGN);
```

It would then use the address of the allocated memory with the *vhandle\_t* value it had received to map the allocated memory into the user space:

```
v_mapphys (vt, kaddr, len)
```

**Note:** There are no special precautions to take when mapping cached memory into user space, or when mapping device registers or bus addresses. However, you should almost never map *uncached memory* into user space. The effects of uncached memory access are hardware dependent and differ between multiprocessors and uniprocessors. Among uniprocessors, the IP26 CPU module has highly restrictive rules for the use of uncached memory (see “Uncached Memory Access in the IP26 CPU” on page 29). In general, mapping uncached memory makes a driver nonportable and is likely to lead to subtle failures that are hard to resolve.

Example 8-3 contains an edited fragment of code from a Silicon Graphics device driver. This pseudo-device driver, whose prefix is *flash\_*, provides access to “flash” PROM in certain computer models. It allows a user process to map the PROM into user space.

**Example 8-3** Edited Fragment of `flash_map()`

```
int flash_map(dev_t dev, vhandle_t *vt, off_t off, long len)
{
    long offset = (long) off; /*Actual offset in flash prom*/
    /* Don't allow requests which exceed the flash prom size */
    if ((offset + len) > FLASHPROM_SIZE)
        return ENOSPC;
    /* Don't allow non page-aligned offsets */
    if ((offset % NBPC) != 0)
        return EIO;
    /* Only allow mapping of entire pages */
    if ((len % NBPC) != 0)
        return EIO;
    return v_mapphys(vt, FLASHMAP_ADDR + offset, len);
}
```

When the driver allocates some memory resource associated with the mapping, and when more than one mapping can be active at a time, the driver needs to tag each memory resource so it can be located when the `pxunmap()` entry point is called. One answer is to use the `vt_gethandle()` macro defined in `sys/region.h`. This macro takes a pointer to a `vhandle_t` and returns a unique pointer-sized integer that can be used to tag allocations. No other information in `sys/region.h` is supported for driver use.

### Entry Point `mmap()`

The `pxmmap()` (note: *two* letters “m”) entry can be used only in a character device driver. The prototype is

```
int pxmmap(dev_t dev, off_t off, int prot);
```

The argument values are

- |             |  |
|-------------|--|
| <i>dev</i>  | A <code>dev_t</code> value from which you can extract both the major and minor device numbers. |
| <i>off</i>  | The offset argument passed to <code>mmap()</code> by the user process.                         |
| <i>prot</i> | Flags showing the access intentions of the user process.                                       |

The function is expected to return the page frame number (PFN) that corresponds to the offset `off` in the device address space. A PFN is an address divided by the page size. (See “Working With Page and Sector Units” on page 200 for page unit conversion functions.)

This entry point is supported only for compatibility with SVR4. When the kernel needs to map a character device, it looks first for `pxmap()`. It calls `pxmmap()` only when `pxmap()` is not available. The differences between the two entry points are as follows:

- This entry point receives no `vhandle_t` argument, so it cannot use `v_mapphys()`. It has to calculate a page frame number, which means that it has to be aware of the current page size (obtainable from the `ptob()` kernel function, see the `ptob(D3)` reference page).
- This entry point does not receive a length argument, so it has to assume a default length for every map (typically the page size).
- When a mapping is created using this entry point, the `pxunmap()` entry is not called.

## Entry Point `unmap()`

The kernel calls the `pfxunmap()` entry point when a mapping is created using the `pfxmap()` entry point. This entry should be supplied, even if it is an empty function, when the `pfxmap()` entry point is supplied. If it is not supplied, the `munmap()` system function returns the `ENODEV` error.

The `pfxunmap()` entry point is only called when the mapped region has been completely unmapped by all processes. For example, suppose a parent process calls `mmap()` to map a device. Then the parent creates one or more child processes using `sproc()`. Each child shares the address space, including the mapped segment. A process in the share group can terminate, or can explicitly `unmap()` the segment or part of the segment; these actions do not result in a call to `pfxunmap()`. Only when the last process with access to the segment has fully unmapped the segment is `pfxunmap()` called.

On entry, the kernel has completed unmapping the object from the user process address space. This entry point does not need to do anything to affect the user address space; it only needs to release any resources that were allocated to support the mapping.

The prototype is

```
int pfxunmap(dev_t dev, vhandle_t *vt);
```

The argument values are

<i>dev</i>	A <i>dev_t</i> value from which you can extract both the major and minor device numbers.
<i>vt</i>	The address of an opaque structure that describes the assigned address in the user process address space.

If the driver allocated no resources to support a mapping, no action is needed here; the entry point can consist of a “return 0” statement.

When the driver does allocate memory to support a mapping, and supports multiple mappings, the driver needs to identify the resource associated with this particular mapping in order to release it. The `vt_gethandle()` function returns a unique number based on the *vt* argument; this can be used to identify resources.

## Interrupt Entry Point

In traditional UNIX, when a hardware device presents an interrupt, the kernel locates the device driver for the device and calls the *pfxintr()* entry point (see the *intr(D2)* reference page). In current practice, a driver must register a specific interrupt handler for each device. The kernel functions for doing this are bus-specific, and are discussed in the bus-specific chapters. For example, the means of registering a PCI interrupt handler is discussed in Chapter 15, “PCI Device Drivers.” However, the discussion of interrupts that follows is still relevant to any interrupt handler.

In principle an interrupt can happen at any time. Normally an interrupt occurs because at some previous time, the driver initiated a device operation. Some devices can interrupt without a preceding command.

### Associating Interrupt to Driver

The association between an interrupt and the driver is established in different ways depending on the hardware.

- For devices on the SCSI bus, all interrupts are handled by a single, low-level driver which notifies a callback function (see Chapter 13, “SCSI Device Drivers”).
- For devices on the PCI bus, the driver registers an interrupt handler using **pci\_intr\_connect()** at the time the device is attached (“Attaching a Device” on page 386).

### Interrupt Handler Operation

When an interrupt occurs, the system is in an unknown state. As a result, the interrupt handler can use only a restricted set of kernel services, and no services that can sleep. In general, the interrupt handler implements the following tasks.

- When the driver supports multiple logical units, use *ivec* to locate the data structure for the interrupting unit.
- Determine the reason for the interrupt by interrogating the device.
- When the interrupt is a response to a device operation, note the success or failure of the command.
- If the driver top half is waiting for the interrupt, waken it.

- If the driver supports polling, and the interrupt represents a pollable event, call **pollwakeupt()**.
- If the device is not in an error state and another operation is waiting to be started, start it.

The details of each of these tasks depends on the hardware and on the design of the data structures used by the driver top half.

### **Mutual Exclusion**

In a uniprocessor system, there is only one CPU and when it is executing the interrupt handler, nothing else is executing. An interrupt handler can only be preempted by an interrupt of higher priority—which would be an interrupt for a different driver, and so would have no conflicts with this driver over the use of data.

In a multiprocessor, an interrupt can be taken on any CPU, while other CPUs continue to execute kernel or user code.

In a multiprocessor, when an interrupt must be handled by a driver that is not marked as multiprocessor-aware (see “Flag `D_MP`” on page 145), the interrupt may be received on some other CPU, but the driver interrupt entry point is always executed on CPU 0.

In a multiprocessor, when the driver is multiprocessor-aware, one or more other CPUs can execute in the driver’s top-half entry points while another CPU executes the driver’s interrupt entry point. An interrupt handler written for a multiprocessor must not assume that it has exclusive use of the driver’s data (see “Planning for Multiprocessor Use” on page 174).

It is theoretically possible in a multiprocessor for a device to interrupt; for one CPU to enter the interrupt handler; and for the device to interrupt again, resulting in multiple concurrent entries to the same interrupt handler. However, IRIX prevents this. You can assume that your interrupt handler code is entered serially, and not used concurrently by multiple CPUs.

### Performance and Latency

Speed in exiting the interrupt handler is critical to system performance. In a uniprocessor, the system is doing nothing else while it executes the handler, and it cannot respond to interrupts of a lower priority. In a multiprocessor, interrupts can be taken by different CPUs. While a CPU executes a handler, that CPU cannot respond to lower-priority interrupts, but other CPUs can be processing user-level code or responding to other interrupts.

### Completing Block I/O

In a block device driver, an I/O operation is represented by the *buf\_t* structure. The *pxstrategy()* routine starts operations and waits for them to complete (see “Entry Point *strategy()*” on page 157).

The interrupt entry point sets the residual count in *b\_resid*. It can post an error using *bioerror()*. It posts the operation complete and wakens the *pxstrategy()* routine by calling *biodone()*. If the *pxstrategy()* entry has set the address of a completion callback function in the *b\_iodone* field of the *buf\_t*, *biodone()* invokes it. (For more discussion, see “Waiting for Block I/O to Complete” on page 219.)

### Completing Character I/O

In a character device driver, the driver top half typically awaits an interrupt by sleeping on a semaphore or synchronizing variable, and the interrupt routine posts the semaphore (see “Waiting for a General Event” on page 221). Error information must be passed in driver variables according to some local convention.

### Calling *pollwakeup()*

When the interrupt represents an event that can be reported by the driver’s *pxpoll()* entry point (see “Entry Point *poll()*” on page 160), the interrupt handler must report the event to the kernel, in case some user process is waiting in a *poll()* call. Hypothetical code to do this is shown in Example 8-4.

**Example 8-4** Hypothetical Call to pollwakeup()

```
hypo_intr(int ivec)
{
    struct hypo_dev_info *pinfo;
    if (! pinfo = find_dev_info(ivec))
        return; /* not our device */
    ...
    if (pinfo->have_data_flag)
        pollwakeup (pinfo->phead, POLLIN, POLLRDNORM);
    if (pinfo->output_ok_flag)
        pollwakeup (pinfo->phead, POLLOUT);
    ...
}
```

## Support Entry Points

Certain driver entry points are used to support the operations of the kernel or the administration of the system.

### Entry Point `unload()`

The `pfxunload()` entry point is called when the kernel is about to dynamically remove a loadable driver from the running system. The prototype is

```
int pfxunload(void);
```

A driver can be unloaded either because all its devices are closed and a timeout has elapsed, or because the operator has used the `ml` command (see the `ml(1)` reference page). The kernel does not unload a driver unless the driver provides a `pfxunload()` entry point. Without this entry point, the driver can be dynamically loaded, but then remains in memory.

It is not easy to retain state information about the device over the time when the driver is not in memory. The entire text and data of a loadable driver, including static variables, are removed and reloaded. Only global variables defined in the descriptive file (see “Describing the Driver in `/var/sysgen/master.d`” on page 235) remain in memory after the driver is unloaded. Be sure not to store any addresses of driver code or driver static variables in global variables, since these addresses will be different when the driver is reloaded.

The driver may have allocated dynamic memory. This should be released, because the addresses of allocated memory will be lost when the driver is unloaded, and more will be allocated if the driver is reloaded. For example, the driver should use **phfree()** to release a pollhead structure allocated by **phalloc()** (see “Use and Operation of poll(2)” on page 159, and the **phalloc(D3)** and **phfree(D3)** reference pages). It is also the time to release any PIO maps (see “Inactivating Maps and Releasing Objects” on page 402), and to release any process handles (see “Sending a Process Signal” on page 206).

The driver is not required to unload. If the driver should not be unloaded at this time, it returns a nonzero return code to the call, and the kernel does not unload it. There are several reasons why a driver should not be unloaded.

The kernel calls **pxunload()** only when no device special files managed by the driver are open. If any device had been opened, the **pxclose()** entry has been called. However, if any device was mapped through the **pxmap()** entry, the mapping could still exist. If the driver has any resources tied up in association with a memory mapping, it should return a nonzero value to the **pxunload** call.

A driver should never permit unloading when there is any kind of pointer to the driver held in any kernel data structure. It is a frequent design error to unload when there is a live pointer to the driver. Unpredictable kernel panics often result.

One example of a live pointer to a driver is a pending callback function. Any pending **itimeout()** or **bufcall()** timers should be cancelled before returning 0 from **pxunload()**. A driver for the PCI bus can register an interrupt handler, and should unregister an interrupt handler (see “Unloading” on page 402) before it permits unloading.

### Entry Point **halt()**

The kernel calls the **pxhalt()** entry point, if one exists, while performing an orderly system shutdown (see the **halt(1)** reference page). No other driver entry points are called after this one. The prototype is simply

```
void pxhalt(void);
```

Since the system is shutting down, there is no point in returning allocated memory. The only purpose this entry point can serve is to leave the device in a safe and stable condition. For example, this is the place at which a disk driver could command the heads of the drive to move to a safe zone for power off.

The driver cannot assume that interrupts are disabled or enabled. The driver cannot block waiting for device actions, so whatever commands it issues to the device must take effect immediately.

### Entry Point `size()`

The `pxsize()` entry point is required of block device drivers. It reports the size of the device in “sector” units, where a “sector” size is declared as `NBPSCTR` in `sys/param.h` (currently 512). The prototype is

```
int pxsize(dev_t dev);
```

The device major and minor numbers can be extracted from the `dev` argument. The entry point is not called until `pxopen()` has been called. Typically the driver will calculate the size of the medium during `pxopen()`.

Since the `int` return value is 32 bits in all systems, the largest possible block device is 1,024 gigabytes ( $(2^{31} * 512) / 1,024^3$ ).

### Entry Point `print()`

The `pxprint()` entry point is called from the kernel to display a diagnostic message when an error is detected on a block device. The prototype and the complete logic of the entry point is shown in Example 8-5.

#### Example 8-5 Entry Point `pxprint()`

```
#include <sys/cmn_err.h>
#include <sys/ddi.h>
int hypo_print(dev_t dev, char *str)
{
    cmn_err(CE_NOTE, "Error on dev %d: %s\n", getemisor(dev), str);
    return 0;
}
```

## Handling 32-Bit and 64-Bit Execution Models

The `pfxioc10` entry point can be passed a data structure from the user process address space; that is, the `arg` value can be a pointer to a structure or an array of data. In order to interpret such a structure, the driver has to know the execution model for which the user process was compiled.

The execution model is specified when code is compiled. The 32-bit model (compiler option `-32` or `-n32`) uses 32-bit address values and a `long int` contains 32 bits. The 64-bit model (compiler option `-64`) uses 64-bit address values and a `long int` contains 64 bits. (The size of an unqualified `int` is 32 bits in both models.) The execution model is sometimes casually called the “ABI” (Authorized Binary Interface), but this is an improper use of that term—an ABI comprises calling conventions, public names, and structure definitions, as well as the execution model.

An IRIX kernel compiled to the 32-bit model contains 32-bit drivers and supports only 32-bit user processes. A kernel compiled to the 64-bit model contains 64-bit drivers, but it supports user processes compiled to *either* 32-bit or 64-bit models. Therefore, in a 64-bit kernel, a driver can be asked to interpret data produced by a 32-bit program.

This is true only of the `pfxioc10` and `pfxpoll0` entry points. Other driver entry points move data to and from user space as streams or blocks of bytes—not as a structure with fields to be interpreted.

Since in other respects it is easy to make your driver portable between 64-bit and 32-bit systems, you should design your driver so that it can handle the case of operating in a 64-bit kernel, receiving `ioc10` requests alternately from 32-bit and 64-bit programs.

The simplest way to do this is to define the arguments passed to the entry points in such a way that they have the same precision in either system. However, this is not always possible. To handle the general case, the driver must know to which model the user process was compiled.

You find this out by calling the `userabi0` kernel function (for which, unfortunately, there is no reference page available).

The prototype of `userabi0` (declared in `sys/ddi.h`) is

```
int userabi(__userabi_t *);
```

If there is no user process context, **userabi()** returns ESRCH. Otherwise it fills out a `__userabi_t` structure and returns 0. The structure of type `__userabi_t` (declared in `sys/types.h`) contains the fields listed below:

<code>uabi_szint</code>	Size of a user int (4).
<code>uabi_szlong</code>	Size of a user long (4 or 8).
<code>uabi_szptr</code>	Size of a user address (4 or 8).
<code>uabi_szlonglong</code>	Size of a user long long (8).

Store the value of `uabi_szptr` when opening a device. Then you can use it to choose between 32-bit and 64-bit declarations of a structure passed to `pfxiocctl()` or an address passed to `pfxpoll()`.

## Planning for Multiprocessor Use

Multiprocessor computers are a central part of the Silicon Graphics product line and will become increasingly common in the future. A device driver that is not multiprocessor-ready can be used in a multiprocessor, but it is likely to cause a performance bottleneck. A multiprocessor-ready driver, on the other hand, works well in a uniprocessor with little if any loss of speed.

### The Multiprocessor Environment

A multiprocessor has two or more CPU modules, all of the same type. The CPUs execute independently, but all share the same main memory. Any CPU can execute the code of the IRIX kernel, and it is common for two or more CPUs to be executing kernel code, including driver code, simultaneously.

### Uniprocessor Assumptions

The original UNIX architecture assumed a uniprocessor hardware environment with a hierarchy of interrupt levels. Ordinary code could be preempted by an interrupt, but an interrupt handler could only be preempted by an interrupt at a higher level.

This assumed hardware environment was reflected in the design of device drivers and kernel support functions.

- In a uniprocessor, an upper-half driver entry point such as `pfxopen()` cannot be preempted except by an interrupt. It has exclusive access to driver variables except for those changed by the interrupt handler.
- Once in an interrupt handler, no other code can possibly execute except an interrupt of a higher hardware level. The interrupt handler has exclusive access to driver variables.
- The interrupt handler can use kernel functions such as `splhi()` to set the hardware interrupt mask, blocking interrupts of all kinds, and thus getting exclusive access to all memory including kernel data structures.

All of these assumptions fail in a multiprocessor.

- Upper-half entry points can be entered concurrently on multiple CPUs. For example, one CPU can be executing `pfxopen()` while another CPU is in `pfxstrategy()`. Exclusive use of driver variables cannot be assumed.
- An interrupt can be taken on one CPU while upper-half routines or a timeout function execute concurrently on other CPUs. The interrupt routine cannot assume exclusive use of driver variables.
- Interrupt-level functions such as `splhi()` are meaningless, since at best they set the interrupt mask on the current CPU only. Other CPUs can accept interrupts at all levels. The interrupt handler can never gain exclusive access to kernel data.

The process of making a driver multiprocessor-ready consists of changing all code whose correctness depends on uniprocessor assumptions.

### Protecting Common Data

Whenever a common resource can be updated by two processes concurrently, the resource must be protected by a *lock* that represents the exclusive right to update the resource. Before changing the resource, the software acquires the lock, claiming exclusive access. After changing the resource, the software releases the lock.

The IRIX kernel provides a set of functions for creating and using locks. It provides another set of functions for creating and using *semaphore* objects, which are like locks but sometimes more flexible. Both sets of functions are discussed under “Waiting and Mutual Exclusion” on page 206.

### Sleeping and Waking

Sometimes the lock is not available—some other process executing in another CPU has acquired the lock. When this happens, the requesting process is delayed in the lock function until the lock is free. To delay, or *sleep*, is allowed for upper-half entry points, because they execute (in effect) as subroutines of user processes.

Interrupt handlers and timeout functions are not permitted to sleep. They have no process identity and so there is no mechanism for saving and restoring their state. An interrupt handler can test a lock, and can claim the lock conditionally, but if a lock is already held, the handler must have some alternate way of storing data.

### Synchronizing Within Upper-Half Functions

When designing an upper-half entry point, keep in mind that it could be executed concurrently with any other upper-half entry point, and that the one entry point could even be executed concurrently by multiple CPUs. Only a few entry points are immune:

- The *pxinit()*, *pxedtinit()*, and *pxstart()* entry points cannot be entered concurrently with each other or any other entry point (*pxstart()* could be entered concurrently with the interrupt handler).
- The *pxunload()* and *pxhalt()* entry points cannot be entered concurrently with any other entry point except for stray interrupts.
- Certain entry points have no cause to use shared data; for example, *pxsize()* and *pxprint()* normally do not need to take any precautions.
- Other upper-half entry points, and all STREAMS entry points, can be entered concurrently by multiple CPUs, when the driver is multiprocessor-aware.

You can deal with concurrency at different levels of sophistication.

### Running on CPU 0

If you do not set the *D\_MP* flag in a character driver or STREAMS driver (see “Flag *D\_MP*” on page 145), the driver is executed only on CPU 0. As a result, upper-half entry points cannot execute concurrently, and the interrupt handler cannot run in true concurrency with an upper-half routine (although it can preempt an upper-half routine as it does in a uniprocessor).

The result is that user processes are serialized for the use of the driver for any purpose. Since CPU 0 is often busy with other housekeeping activities, access to the driver can have a latency that is long and variable.

### **Serializing on a Single Lock**

You can create a single lock for upper-half serialization. Each upper-half function begins with read-only operations such as extracting the device minor number and testing and validating arguments. You allow these to execute concurrently on any CPU (the `D_MP` flag is set.)

In each entry point, when the preliminaries are complete, you acquire the single lock, and release it just before returning. The result is that processes are serialized for I/O through the driver. If the driver supports only a single device, processes would be serialized in any case, waiting for the device to operate. Since the upper half can execute on any CPU, latency is more predictable.

### **Serializing on a Lock Per Device**

When the driver supports multiple minor devices, you will normally have a data structure per device, indexed by the device minor number. Typically an upper-half routine is concerned only with one minor device. You can define a lock in the data structure for the minor device, and acquire that lock as soon as the device number is known.

This permits concurrent execution of upper-half requests for different minor devices, while serializing access to any one device.

### **Coordinating Upper-Half and Interrupt Entry Points**

Upper-half entry points prepare work for the device to do, and the interrupt routine reports the completion of the device action. In a block device driver, this communication is relatively simple. In a character driver, you have more design options. The kernel functions mentioned in the following topics are covered under “Waiting and Mutual Exclusion” on page 206.

### Coordinating Through the `buf_t`

In a block device driver, the `pxstrategy()` routine initiates a read or a write based on a `buf_t` structure (see “Entry Point `strategy()`” on page 157), and leaves the address of the `buf_t` where the interrupt routine can find it. Then `pxstrategy()` calls the `biowait()` kernel function to wait for completion.

The `pxintr()` entry point updates the `buf_t` (using `pxbioerror()` if necessary) and then uses `biodone()` to mark the `buf_t` as complete. This ends the wait for `pxstrategy()`. These kernel functions are multiprocessor-aware.

### Coordination in a Character Driver

In a character driver that supports interrupts, you design your own coordination mechanism. The simplest (and not recommended) would be based on using the kernel function `sleep()` in the upper half, and `wakeup()` in the interrupt routine. You can also use a semaphore and use `psema()` in the upper half and `vsema()` in the interrupt handler.

If you need to allow for timeouts, you have to deal with the complication that the timeout function can be called concurrently with an interrupt. When you use a semaphore, the interrupt routine can use `vsema()` to post completion, and the timeout function can use `cvsema()` to post it only if it has not already been posted.

### Converting a Uniprocessor Driver

As a general approach, you can convert a uniprocessor driver to make it multiprocessor-safe in the following steps:

1. If it currently uses the `D_OLD` flag (or has no `pxdevflag` constant), convert it to use the current interface, with a `pxdevflag` of `0x00`.
2. Make sure it works in the original uniprocessor at the current release of IRIX.
3. Test it in a multiprocessor running in CPU 0.
4. Begin adding semaphores, locks, and other exclusion and synchronization tools. Since the driver still runs serially on CPU 0, it will never wait for a lock, but the coordination between upper half and interrupt handler should work.
5. Add the `D_MP` flag and test on a multiprocessor.

In performing the conversion, you can use calls to `spl..()` functions as signs that work is needed. These functions are used for mutual exclusion in a uniprocessor, and they are all ineffective or unnecessary in a multiprocessor-safe driver.

### Example Conversion Problem

The code in Example 8-6 shows typical logic in a uniprocessor character driver.

#### Example 8-6 Uniprocessor Upper-Half Wait Logic

```
s = splvme();
flag |= WAITING;
while (flag & WAITING) {
    sleep(&flag, PZERO);
}
splx(s);
```

The upper half calls the `splvme()` function with the intention of blocking interrupts, and thus preventing execution of this driver's interrupt handler while the *flag* variable is updated. In a multiprocessor this is ineffective because at best it sets the interrupt level on the current CPU. The interrupt handler can execute on another CPU and change the variable.

The corresponding interrupt handler is sketched in Example 8-7.

#### Example 8-7 Uniprocessor Interrupt Logic

```
if (flag & WAITING) {
    wakeup(&flag);
    flag &= ~WAITING;
}
```

The interrupt handler could execute on another CPU, and test the flag after the upper half has called `splvme()` and before it has set `WAITING` in *flag*. The interrupt is effectively lost. This would happen rarely and would be hard to repeat, but it would happen and would be hard to trace.

A more reliable, and simpler, technique is to use a semaphore. The driver defines a global semaphore:

```
static sema_t sleeper;
```

A driver with multiple devices would have a semaphore per device, perhaps as an array of *sema\_t* items indexed by device minor number.

The semaphore (or array) would be initialized to a starting value of 1 in the *pfxinit()* or *pfxstart()* entry:

```
void hypo_start()  
{  
...  
    initnsema(&sleeper, 1, "sleeper");  
}
```

After the upper half started a device operation, it would await the interrupt using **psema()**:

```
psema(sleeper, PZERO);
```

The PZERO argument makes the wait immune to signals. If the driver should wake up when a signal is sent to the calling process (such as SIGINT or SIGTERM), the second argument can be PCATCH. A return value of -1 indicates the semaphore was posted by a signal, not by a **vsema()** call.

The interrupt handler would use **vsema()** or **cvsema()** to post the semaphore. The use of **cvsema()** ensures that the semaphore is not incremented past 1, in the event that it is posted from more than one location (as from a timeout or a signal handler).

---

## Device Driver/Kernel Interface

The programming interface between a device driver and the IRIX kernel is completely documented in the reference pages in volume "D." This chapter provides a survey and a summary of the API under the following headings:

- "Important Data Types" on page 182 describes the data types that are exchanged between the kernel and a driver.
- "Important Header Files" on page 188 summarizes the C header files that are frequently included in a driver source file.
- "Memory Allocation" on page 189 describes the kernel functions for general memory allocation, for allocation of objects of specific types, and for resource suballocation.
- "Transferring Data" on page 194 describes the functions for transferring data between a driver and a buffer in the address space of either the kernel or a user process.
- "Managing Virtual and Physical Addresses" on page 198 discusses the translation from virtual to physical storage locations for DMA.
- "User Process Administration" on page 205 discusses process signalling and authentication.
- "Waiting and Mutual Exclusion" on page 206 describes and summarizes a wide array of functions you can use for those purposes.

In addition to these topics, data types and functions specific to the following areas are in the chapters shown:

Debugging and logging	Chapter 11, "Testing and Debugging a Driver."
PCI bus	Chapter 15, "PCI Device Drivers"
SCSI bus	Chapter 13, "SCSI Device Drivers"
STREAMS drivers	Chapter 16, "STREAMS Drivers"

## Important Data Types

### The Device Number Types

Two numbers are carried in the inode of a *device special file*: a *major device number* of up to 9 bits, and a *minor device number* of up to 18 bits. The numbers are assigned when the device special file is created, either by the `/dev/MAKEDEV` script or by the system administrator. The contents and meaning of device numbers is discussed under “Device Representation” on page 35.

At almost every upper-half entry point, the first argument to a driver is a `dev_t` object, an unsigned integer containing the values of the major and minor numbers for the device that is to be used. The `dev_t` type is declared in `sys/types.h` along with types `major_t` and `minor_t`, which represent major and minor numbers as variables.

### Use of the Device Numbers

You typically use the major device number to learn which device driver has been called. This is important only when a device driver supports multiple interfaces, for example when one driver represents both character and block access to the same hardware.

You use the minor device number to learn which hardware unit is being accessed. This is of interest only when a driver supports multiple units. In addition, device management options can be encoded into the minor number, as described under “Minor Device Number” on page 37.

### Device Number Functions

The kernel provides several functions for manipulating device numbers, and these are summarized in Table 9-1.

**Table 9-1** Functions to Manipulate Device Numbers

Function	Header Files	Can Sleep	Purpose
<code>etoimajor(D3)</code>	<code>ddi.h</code>	N	Convert external to internal major device number.
<code>getemajor(D3)</code>	<code>ddi.h</code>	N	Get external major device number.

**Table 9-1 (continued)** Functions to Manipulate Device Numbers

Function	Header Files	Can Sleep	Purpose
getemisor(D3)	ddi.h	N	Get external minor device number.
getemisor(D3)	ddi.h	N	Get internal major device number.
getemisor(D3)	ddi.h	N	Get internal minor device number.
itoemisor(D3)	ddi.h	N	Convert internal to external major device number.
makeemisor(D3)	ddi.h	N	Make device number from major and minor numbers.

**Note:** Under no circumstances should you decode the *dev\_t* using Boolean operations to extract major and minor numbers. Use the functions listed in Table 9-1. Drivers that treat the *dev\_t* as an integer will stop working in the next release of IRIX after IRIX 6.3 for O2.

The most important of the functions in in Table 9-1 are

- **getemisor()**, which extracts the major number from a *dev\_t* and returns it as a *major\_t*
- **getemisor()**, which extracts the minor number from a *dev\_t* and returns it as a *minor\_t*
- **makeemisor()**, which combines a *major\_t* and a *minor\_t* to form a *dev\_t*

### External and Internal Numbers

The kernel uses the major device number as a subscript to index various tables. Some variants of UNIX, in order to avoid wasting space on sparse tables, translate the major device number to an internal code. Sometimes the minor number is translated too.

This internal encoding of the device number is of no interest in IRIX. If it is done, it is done only for the purpose of subscripting tables within the kernel that are not accessible to device drivers. Internal device numbers have no utility in IRIX. However, functions related to internal device numbers are included for compatibility with SVR4.

If you are writing a new device driver specifically for IRIX, use only external device numbers. If you are porting a device driver that uses the **getemisor()**, **getemisor()**, **etoemisor()** and **etoemisor()** functions, you can leave these function calls unchanged.

(But if the driver attempts to access the kernel switch tables, it is nonportable and should be changed.)

### Structure `uio_t`

The `uio_t` structure describes data transfer for a character device:

- The `pfxread()` and `pfxwrite()` entry points receive a `uio_t` that describes the buffer of data.
- Within an `pfxiocctl()` entry point, you might construct a `uio_t` to represent data transfer for control purposes.
- In a hybrid character/block driver, the `physiock()` function translates a `uio_t` into a `buf_t` for use by the `pfxstrategy()` entry point.

The fields and values in a `uio_t` are declared in `sys/uio.h`, which is included by `sys/ddi.h`. For a detailed discussion, see the `uio(D4)` reference page. Typically the contents of the `uio_t` reflect the buffer areas that were passed to a `read()`, `readv()`, `write()`, or `writew()` call (see the `read(2)` and `write(2)` reference pages).

### Data Location and the `iovec_t`

One `uio_t` describes data transfer to or from a single address space, either the address space of a user process or the kernel address space. The address space is indicated by a flag value, either `UIO_USERSPACE` or `UIO_SYSSPACE`, in the `uio_segflg` field.

The total number of bytes remaining to be transferred is given in field `uio_resid`. Initially this is the total requested transfer size.

Although the transfer is to a single address space, it can be directed to multiple segments of data within the address space. Each segment of data is described by a structure of type `iovec_t`. An `iovec_t` contains the virtual address and length of one segment of memory.

The number of segments is given in field `uio_iovcnt`. The field `uio_iov` points to the first `iovec_t` in an array of `iovec_t` structures, each describing one segment of data. The total size in `uio_resid` is the sum of the segment sizes.

For a simple data transfer, *uio\_iovcnt* contains 1, and *uio\_iov* points to a single *iovec\_t* describing a buffer of 1 or more bytes. For a complicated transfer, the *uio\_t* might describe a number of scattered segments of data. Such transfers can arise in a network driver where multiple layers of message header data are added to a message at different levels of the software.

### Use of the *uio\_t*

In the *pxread()* and *pxwrite()* entry points, you can test *uio\_segflag* to see if the data is destined for user space or kernel space, and you can save the initial value of *uio\_resid* as the requested length of the transfer.

In a character driver, you fetch or store data using functions that both use and modify the *uio\_t*. These functions are listed under “Transferring Data Through a *uio\_t* Object” on page 197. When data is not immediately available, you should test for the FNDELAY or FNONBLOCK flags in *uio\_fmode*, and return when either is set rather than sleeping.

### Structure *buf\_t*

The *buf\_t* structure describes a block data transfer. It is designed to represent the transfer (in or out) of a sequence of adjacent, fixed-size blocks from a random-access device to a block of contiguous memory. The size of one device block is NBPSCTR, declared in *sys/param.h*. For a detailed discussion of the *buf\_t*, see the *buf(D4)* reference page.

The *buf\_t* is used internally in IRIX by the paging I/O system to manage queues of physical pages, and by filesystems to manage queues of pages of file data. The paging system and filesystems are the primary clients of the *pxstrategy()* entry point to a block device driver, so it is only natural that a *buf\_t* pointer is the input argument to *pxstrategy()*.

**Tip:** The *idbg* kernel debugging tool has several functions related to displaying the contents of *buf\_t* objects. See “Commands to Display *buf\_t* Objects” on page 270.

### Fields of *buf\_t*

The fields of the *buf\_t* are declared in *sys/buf.h*, which is included by *sys/ddi.h*. This header file also declares the names of many kernel functions that operate on *buf\_t* objects. (Many of those functions are not supported as part of the DDI/DKI. You should only use kernel functions that have reference pages.)

Because *buf\_t* is used by so many software components, it has many fields that are not relevant to device driver needs, as well as some fields that have multiple uses. The relevant fields are summarized in Table 9-2.

**Table 9-2** Accessible Fields of *buf\_t* Objects

Field Name	Access	Purpose and Contents
<i>b_edev</i>	read-only	<i>dev_t</i> giving device major and minor numbers.
<i>b_flags</i>	read-only	Operational flags; for a detailed list see <i>buf(D4)</i> .
<i>b_forw, b_back, av_forw, av_back</i>	read-write	Queuing pointers, available for driver use within the <i>pxstrategy()</i> routine.
<i>b_un.b_addr</i>	read-only	Sometimes the kernel virtual address of the buffer, depending on the <i>b_flags</i> setting <i>BP_ISMAPPED</i> .
<i>b_bcount</i>	read-only	Number of bytes to transfer.
<i>b_blkno</i>	read-only	Starting logical block number on device (for a disk, relative to the partition that the device represents).
<i>b_iodone</i>	read-write	Address of a driver internal function to be called on I/O completion.
<i>b_resid</i>	read-write	Number of bytes not transferred, set at completion to 0 unless an error occurs.
<i>b_error</i>	read-write	Error code, set at completion of I/O.

No other fields of the *buf\_t* are designed for use by a driver. In Table 9-2, “read-only” access means that the driver should never change this field in a *buf\_t* that is owned by the kernel. When the driver is working with a *buf\_t* that the driver has allocated (see “Allocating *buf\_t* Objects and Buffers” on page 192) the driver can do what it likes.

### Using the Logical Block Number

The logical block number is the number of the 512-byte block in the device. The “device” is encoded by the minor device number that you can extract from *b\_edev*. It might be a complete device surface, or it might be a partition within a larger device (for example, the IRIX disk device drivers support different minor device numbers for different disk partitions).

The `pxstrategy()` routine may have to translate the logical block number based on the driver's information about device partitioning and device geometry (sector size, sectors per track, tracks per cylinder).

### Buffer Location and `b_flags`

The data buffer represented by a `buf_t` can be in one of two places, depending on bits in `b_flags`.

When the macro `BP_ISMAPPED(buf_t-address)` returns true, the buffer is in kernel virtual memory and its virtual address is in `b_un.b_addr`.

When `BP_ISMAPPED(buf_t-address)` returns false, the buffer is described by a chain of `pfdat` structures (declared in `sys/pfdat.h`, but containing no fields of any use to a device driver). In this case, `b_un.b_addr` contains only an offset into the first page frame of the chain. See "Managing Buffer Virtual Addresses" on page 202 for a method of mapping an unmapped buffer.

### Lock and Semaphore Types

The header files `sys/sema.h` and `sys/types.h` declare the data types of locks of different types, including the following:

<code>lock_t</code>	Basic lock, or spin-lock, used with <code>LOCK()</code> and related functions
<code>mutex_t</code>	Sleeping lock, used for mutual exclusion between upper-half instances.
<code>sema_t</code>	Semaphore object, used for general locking.
<code>mrlock_t</code>	Reader-writer locks, used with <code>RW_RDLOCK()</code> and related functions.
<code>sv_t</code>	Synchronization variable, used with <code>SV_WAIT</code> and related functions

These lock types should be treated as opaque objects because their contents can change from release to release (and in fact their contents are different in IRIX 6.2 from previous releases).

The families of locking and synchronization functions contain functions for allocating, initializing, and freeing each type of lock. See "Waiting and Mutual Exclusion" on page 206.

## Important Header Files

The header files that are frequently needed in device driver source modules are summarized in Table 9-3.

**Table 9-3** Header Files Often Used in Device Drivers

Header File	Reason for Including
<i>sys/buf.h</i>	The <i>buf_t</i> structure and related constants and functions (included by <i>sys/ddi.h</i> ).
<i>sys/cmn_err.h</i>	The <b>cmn_err()</b> function.
<i>sys/conf.h</i>	The constants used in the <i>pxdevflags</i> global.
<i>sys/ddi.h</i>	Many kernel functions declared. Also includes <i>sys/types.h</i> , <i>sys/uidio.h</i> , and <i>sys/buf.h</i> .
<i>sys/debug.h</i>	Defines the ASSERT macro and others.
<i>sys/dmmap.h</i>	Data types and kernel functions related to DMA mapping.
<i>sys/edt.h</i>	Declares the <i>edt_t</i> type passed to <i>pxedtinit()</i> .
<i>sys/eisa.h</i>	EISA-bus hardware constants and EISA kernel functions.
<i>sys/errno.h</i>	Names for all system error codes.
<i>sys/file.h</i>	Names for file mode flags passed to driver entry points.
<i>sys/immu.h</i>	Types and macros used to manage virtual memory and some kernel functions.
<i>sys/kmem.h</i>	Constants like KM_SLEEP used with some kernel functions.
<i>sys/ksynch.h</i>	Functions used for sleep-locks.
<i>sys/log.h</i>	Types and functions for using the system log.
<i>sys/major.h</i>	Names for assigned major device numbers.
<i>sys/map.h</i>	Types and functions used for suballocation using <b>rmalloc()</b> .
<i>sys/mman.h</i>	Constants and flags used with <b>mmap()</b> and the <i>pxmmap()</i> entry point.
<i>sys/param.h</i>	Constants like PZERO used with some kernel functions.
<i>sys/PCI/pciio.h</i>	PCI bus interface functions and constants.
<i>sys/pio.h</i>	VME PIO functions.

**Table 9-3 (continued)** Header Files Often Used in Device Drivers

Header File	Reason for Including
<i>sys/poll.h</i>	Types and functions for pollhead allocation and poll callback.
<i>sys/scsi.h</i>	Types and functions used to call the inner SCSI driver.
<i>sys/sem.h</i>	Types and functions related to semaphores, mutex locks, and basic locks.
<i>sys/stream.h</i>	STREAMS standard functions and data types.
<i>sys/strmp.h</i>	STREAMS multiprocessor functions.
<i>sys/sysmacros.h</i>	Macros for conversion between bytes and pages, and similar values.
<i>sys/system.h</i>	Kernel functions related to system operations.
<i>sys/types.h</i>	Common data types and types of system objects (included by <i>sys/ddi.h</i> ).
<i>sys/uio.h</i>	The <i>uio_t</i> structure and related functions (included by <i>sys/ddi.h</i> ).
<i>sys/vmereg.h</i>	VME bus hardware constants and VME-related functions.

## Memory Allocation

A device or STREAMS driver can allocate memory statically, as global variables in the driver module, and this is a good way to allocate any object that is always needed and has a fixed size.

When the number or size of an object can vary, but can be determined at initialization time, the driver can allocate memory in the *pxinit()*, *pxedinit()*, or *pxstart()* entry point.

You can allocate memory dynamically in an upper-half entry point. When this is necessary, it should be done in an entry point that is called infrequently, such as *pxopen()*. The reason is that memory allocation is subject to unpredictable delays.

Memory allocation should never be attempted in an interrupt routine. The resources that might be needed at interrupt time should be obtained and set aside by an upper-half entry point before the interrupt is made possible.

## General-Purpose Allocation

There are two groups of general-purpose functions used to allocate and release memory.

- **kmem\_alloc()** and two associated functions supply a complete set of services for allocating kernel virtual memory.
- **kern\_malloc()** and two associated functions are an obsolete mechanism for allocating kernel virtual memory.

The functions you can use to dynamically allocate kernel virtual memory are summarized in Table 9-4.

**Table 9-4** Functions for Kernel Virtual Memory

Function Name	Header Files	Can Sleep?	Purpose
kmem_alloc(D3)	kmem.h & types.h	Y	Allocate space from kernel free memory.
kmem_free(D3)	kmem.h & types.h	N	Free previously allocated kernel memory.
kmem_zalloc(D3)	kmem.h & types.h	Y	Allocate and clear space from kernel free memory.
kern_calloc(D3)	system.h & types.h	Y	Allocate space from kernel memory and clear it.
kern_free(D3)	system.h & types.h	N	Free kernel memory space.
kern_malloc(D3)	system.h & types.h	Y	Allocate kernel virtual memory.

The most important of these functions is **kmem\_alloc()**. You use it to allocate blocks of virtual memory at any time. It offers these important options, controlled by a flag argument:

- Sleeping or not sleeping when space is not available. You specify not-sleeping when in a lower-half routine or when holding a basic lock, but then you must be prepared to deal with a return value of NULL.

- Physically-contiguous memory. The memory allocated is virtual, and when it spans multiple pages, the pages are not necessarily adjacent in physical memory. You need physically contiguous pages when doing DMA with a device that cannot do scatter/gather. However, contiguous memory is harder to get as the system runs, so it is best to obtain it in an initialization routine.
- Cache-aligned memory. By requesting memory that is a multiple of a cache line in size, and aligned on a cache-line boundary, you ensure that DMA operations will affect the fewest cache lines (see “Setting Up a DMA Transfer” on page 201).

The **kmem\_zalloc()** function takes the same options, but offers the additional service of zero-filling the allocated memory.

Calls to the “kern” group of functions should be replaced as follows:

<b>kern_malloc(<i>n</i>)</b>	Change to <b>kmem_alloc(<i>n</i>,KM_SLEEP)</b> .
<b>kern_calloc(<i>n</i>,<i>s</i>)</b>	Change to <b>kmem_zalloc(<i>n</i>*<i>s</i>,KM_SLEEP)</b>
<b>kern_free(<i>p</i>)</b>	Change to <b>kmem_free(<i>p</i>)</b>

### Allocating Objects of Specific Kinds

The kernel provides a number of functions with the purpose of allocating and freeing objects of specific kinds. Many of these are variants of **kmem\_alloc()** and **kmem\_free()**, but others use special techniques suited to the type of object.

### Allocating pollhead Objects

Table 9-5 summarizes the functions you use to allocate and free the *pollhead* structure that is used within the *pxpoll()* entry point (see “Entry Point *poll()*” on page 160). Typically you would call **phalloc()** while initializing each minor device, and call **phfree()** in the *pxunload()* entry point.

**Table 9-5** Functions for Allocating pollhead Structures

Function Name	Header Files	Can Sleep?	Purpose
phalloc(D3)	ddi.h & kmem.h & poll.h	Y	Allocate and initialize a pollhead structure.
phfree(D3)	ddi.h & poll.h	N	Free a pollhead structure.

### Allocating Semaphores and Locks

There are symmetrical pairs of functions to allocate and free all types of lock and synchronization objects. These functions are summarized together with the other locking functions under “Waiting and Mutual Exclusion” on page 206.

### Allocating buf\_t Objects and Buffers

The argument to the *pxstrategy()* entry point is a *buf\_t* structure that describes a buffer (see “Entry Point *strategy()*” on page 157 and “Structure *buf\_t*” on page 185).

Ordinarily, both the *buf\_t* and the buffer are allocated and initialized by the kernel or the filesystem that calls *pxstrategy(0)*. However, some drivers need to create a *buf\_t* and associated buffer for special uses. The functions summarized in Table 9-6 are used for this.

**Table 9-6** Functions for Allocating *buf\_t* Objects and Buffers

Function Name	Header Files	Can Sleep?	Purpose
<i>geteblk(D3)</i>	<i>ddi.h</i>	Y	Allocate a <i>buf_t</i> and a buffer of 1024 bytes.
<i>ngeteblk(D3)</i>	<i>ddi.h</i>	Y	Allocate a <i>buf_t</i> and a buffer of specified size.
<i>breuse(D3)</i>	<i>ddi.h</i>	N	Return a buffer header and buffer to the system.
<i>getrbuf(D3)</i>	<i>ddi.h</i>	Y	Allocate a <i>buf_t</i> with no buffer.
<i>freerbuf(D3)</i>	<i>ddi.h</i>	N	Free a <i>buf_t</i> with no buffer.

To allocate a *buf\_t* and its associated buffer in kernel virtual memory, use either ***geteblk(0)*** or ***ngeteblk(0)***. Free this pair of objects using ***breuse(0)***, or by calling ***biodone(0)***.

You can allocate a *buf\_t* to describe an existing buffer—one in user space, statically allocated in the driver, or allocated with ***kmem\_alloc(0)***—using ***getrbuf(0)***. Free such a *buf\_t* using ***freerbuf(0)***.

## Suballocation Functions

The functions summarized in Table 9-7 are used to manage suballocation of any resource.

**Table 9-7** Functions for Suballocation

Function Name	Header Files	Can Sleep?	Purpose
<i>rmalloc(D3)</i>	<i>map.h</i> & <i>types.h</i>	N	Allocate space from a private space management map.
<i>rmalloc_wait(D3)</i>	<i>map.h</i> & <i>types.h</i>	Y	Allocate resources from a space management map.
<i>rmallocmap(D3)</i>	<i>map.h</i> & <i>types.h</i>	N	Allocate and initialize a private space management map.

**Table 9-7 (continued)** Functions for Suballocation

Function Name	Header Files	Can Sleep?	Purpose
<code>rmfree(D3)</code>	<code>map.h</code> & <code>types.h</code>	N	Release resources into a space management map.
<code>rmfreemap(D3)</code>	<code>map.h</code> & <code>types.h</code>	N	Free a private space management map.

You use these functions as a convenient, efficient set of subroutines for allocating some resource—for example, disk sectors—that you obtain by other means. The expected sequence of use is as follows.

1. During driver initialization, or possibly in `pxopen()`, use `rmmallocmap()` to allocate a map. A map is a data structure large enough to keep track of as many objects as you will create. Initially the map reflects no available resources.
2. Use `rmfree()` to release existing resources into the map. For example, while opening a disk drive, you could use `rmfree()` to release all unused sectors into a sector map.
3. When a resource is needed in an upper-half routine, use `rmmalloc()` or `rmmalloc_wait()` to acquire it. The index number of the first allocated item is returned.
4. When a resource is released in any entry point, use `rmfree()` to note the available items and to wake up any upper-half process waiting in `rmmalloc_wait()`.
5. On device close or when the driver is unloaded, use `rmfreemap()` to release the map itself.

## Transferring Data

The device driver executes in the kernel virtual address space, but it must transfer data to and from the address space of a user process. The kernel supplies two kinds of functions for this purpose:

- functions that transfer data between driver variables and the address space of the current process
- functions that transfer data between driver variables and the buffer described by a `uio_t` object

**Warning:** The use of an invalid address in kernel space with any of these functions causes a kernel panic.

All functions that reference an address in user process space can sleep, because the page of process space might not be resident in memory. As a result, such functions cannot be used in an interrupt handler, or while holding a basic lock.

## General Data Transfer

The kernel supplies functions for clearing and copying memory within the kernel virtual address space, and between the kernel address space and the address space of the user process that is the current context. These general-purpose functions are summarized in Table 9-8.

**Table 9-8** Functions for General Data Transfer

Function Name	Header Files	Can Sleep?	Purpose
bcopy(D3)	ddi.h	N	Copy data between address locations in the kernel.
bzero(D3)	ddi.h	N	Clear memory for a given number of bytes.
copyin(D3)	ddi.h	Y	Copy data from a user buffer to a driver buffer.
copyout(D3)	ddi.h	Y	Copy data from a driver buffer to a user buffer.
fubyte(D3)	system.h & types.h	Y	Load a byte from user space.
fuword(D3)	system.h & types.h	Y	Load a word from user space.
hwcpin(D3)	system.h & types.h	N	Copy data from device registers to kernel memory.
hwcpout(D3)	system.h & types.h	N	Copy data from kernel memory to device registers.
subyte(D3)	system.h & types.h	Y	Store a byte to user space.
suword(D3)	system.h & types.h	Y	Store a word to user space.

### Block Copy Functions

The **bcopy()** and **bzero()** functions are used to copy and clear data areas within the kernel address space, for example driver buffers or work areas. These are optimized routines that take advantage of available hardware features.

The **bcopy()** function is not appropriate for copying data between a buffer and a device; that is, for copying between virtual memory and the physical memory addresses that represent a range of device registers (or indeed any uncached memory). The reason is that **bcopy()** uses doubleword moves and any other special hardware features available, and devices may not be able to accept data in these units. The **hwcpin()** and **hwcpout()** functions copy data in 16-bit units; use them to transfer bulk data between device space and memory. (Use simple assignment to move single words or bytes.)

The **copyin()** and **copyout()** functions take a kernel virtual address, a process virtual address, and a length. They copy the specified number of bytes between the kernel space and the user space. They select the best algorithm for copying, and take advantage of memory alignment and other hardware features.

If there is no current context, or if the address in user space is invalid, or if the address plus length is not contained in the user space, the functions return -1. This indicates an error in the request passed to the driver entry point, and the driver normally returns an EFAULT error.

### Byte and Word Functions

The functions **fubyte()**, **subyte()**, **fuword()**, and **suword()** are used to move single items to or from user space. When only a single byte or word is needed, these functions have less overhead than the corresponding **copyin()** or **copyout()** call. For example you could use **fuword()** to pick up a parameter using an address passed to the *pxioctl()* entry point. When transferring more than a few bytes, a block move is more efficient.

## Transferring Data Through a `uio_t` Object

A `uio_t` object defines a list of one or more segments in the address space of the kernel or a user process (see “Structure `uio_t`” on page 184). The kernel supplies three functions for transferring data based on a `uio_t`, and these are summarized in Table 9-9.

**Table 9-9** Functions Moving Data Using `uio_t`

Function	Header Files	Can Sleep?	Purpose
<code>uiomove(D3)</code>	<code>ddi.h</code>	Y	Copy data using <code>uio_t</code> .
<code>ureadc(D3)</code>	<code>ddi.h</code>	Y	Copy a character to space described by <code>uio_t</code> .
<code>uwritec(D3)</code>	<code>ddi.h</code>	Y	Return a character from space described by <code>uio_t</code> .

The **`uiomove()`** function moves multiple bytes between a buffer in kernel virtual space—typically, a buffer owned by the driver—and the space or spaces described by a `uio_t`. The function takes a byte count and a direction flag as arguments, and uses the most efficient mechanism for copying.

The **`ureadc()`** and **`uwritec()`** functions transfer only a single byte. You would use them when transferring data a byte at a time by PIO. When moving more than a few bytes, **`uiomove()`** is faster.

All of these functions modify the `uio_t` to reflect the transfer of data:

- `uio_resid` is decremented by the amount moved
- In the `iovec_t` for the current segment, `iov_base` is incremented and `iov_len` is decremented
- As segments are used up, `uio_iov` is incremented and `uio_iovcnt` is decremented

The result is that the state of the `uio_t` always reflects the number of bytes remaining to transfer. When the **`pfxread()`** or **`pfxwrite()`** entry point returns, the kernel uses the final value of `uio_resid` to compute the count returned to the **`read()`** or **`write()`** function call.

## Managing Virtual and Physical Addresses

The kernel supplies functions for querying the address of hardware registers and for performing memory mapping.

### Testing Device Physical Addresses

A family of functions, summarized in Table 9-10, is used to test a physical address to find out if it represents a usable device register.

**Table 9-10** Functions to Test Physical Addresses

Function Name	Header Files	Can Sleep?	Purpose
badaddr(D3)	system.h	N	Test physical address for input.
badaddr_val(D3)	system.h	N	Test physical address for input and return the input value received.
wbadaddr(D3)	system.h	N	Test physical address for output.
wbadaddr_val(D3)	system.h	N	Test physical address for output of specific value.
pio_badaddr(D3)	pio.h & types.h	N	Test physical address for input through a map.
pio_badaddr_val(D3)	pio.h & types.h	N	Test physical address for input through a map and return the input value received.
pio_wbadaddr(D3)	pio.h & types.h	N	Test physical address through a map for output.
pio_wbadaddr_val(D3)	pio.h & types.h	N	Test physical address through a map for output of specific value.

The functions return a nonzero value when the address is bad, that is, unusable. The allocation of a PIO map is bus-dependent and is covered in each chapter on a specific bus.

You normally use these functions in the `pxinit()` entry point to verify that an expected device is in fact present. The functions can also be useful in the `pxedtinit()` entry point.

However, that entry point is only called from a VECTOR statement, and the VECTOR statement can contain a PROBE argument that tests for valid hardware.

**Note:** These functions must not be called in an interrupt handler. Verify device addresses in the upper-half code, during initialization.

### Managing Mapped Memory

The `pfxmap()` and `pfxunmap()` entry points receive a `vhandl_t` object that describes the region of user process space to be mapped. The functions summarized in Table 9-11 are used to manipulate that object.

**Table 9-11** Functions to Manipulate a `vhandl_t` Object

Function Name	Header Files	Can Sleep?	Purpose
<code>v_getaddr(D3)</code>	<code>region.h</code> & <code>types.h</code>	N	Get the user virtual address associated with a <code>vhandl_t</code> .
<code>v_gethandle(D3)</code>	<code>region.h</code> & <code>types.h</code>	N	Get a unique identifier associated with a <code>vhandl_t</code> .
<code>v_getlen(D3)</code>	<code>region.h</code> & <code>types.h</code>	N	Get the length of user address space associated with a <code>vhandl_t</code> .
<code>v_mapphys(D3)</code>	<code>region.h</code> & <code>types.h</code>	N	Map kernel address space into user address space.

The `v_mapphys()` function actually performs a mapping between a kernel address and a segment described by a `vhandl_t` (see “Entry Point `map()`” on page 163).

The `v_getaddr()` function has hardly any use except for logging and debugging. The address in user space is normally undefined and unusable when the `pfxmap()` entry point is called, and mapped to kernel space when `pfxunmap()` is called. The driver has no practical use for this value.

The `v_getlen()` function is useful only in the `pfxunmap()` entry point—the `pfxmap()` entry point receives a length argument specifying the desired region size.

The `v_gethandle()` function returns a number that is unique to this mapping (actually, the address of a page table entry). You use this as a key to identify multiple mappings, so that the `pfxunmap()` entry point can properly clean up.

**Caution:** Be careful when mapping device registers to a user process. Memory protection is available only on page boundaries, so configure the addresses of I/O cards so that each device is on a separate page or pages. When multiple devices are on the same page, a user process that maps one device can access all on that page. This can cause system security problems or other problems that are hard to diagnose.

### Working With Page and Sector Units

In a 32-bit kernel, the page size for memory and I/O is 4 KB. In a 64-bit kernel, the memory page size is 16 KB, but because of hardware constraints such as the 4 KB span of DMA mapping registers in the Challenge and Onyx systems, a 4 KB page is used for I/O operations.

The header files `sys/immu.h` and `sys/sysmacros.h` contain constants and macros for working with page units. Some of the most useful are listed below:

<code>NBPP</code>	Number of bytes in a virtual memory page.
<code>NBPSCTR</code>	Number of bytes (512) in a standard disk "sector."
<code>IO_NBPP</code>	Number of bytes in an I/O page.
<code>IO_PNUMSHFT</code>	Number of bits to right-shift an address to get the I/O page number.
<code>IO_POFFMASK</code>	Mask to extract the I/O-page-offset value from an address.
<code>btod()</code>	Return number of 512-byte "sectors" in a byte count (rounded up)
<code>btop(x)</code>	Return number of I/O pages in a byte count (truncated)
<code>io_pnum(x)</code>	Return the I/O page number from an address <i>x</i> .
<code>io_poff(x)</code>	Return the I/O page offset from an address <i>x</i> .
<code>io_numpages(addr, len)</code>	Return the number of I/O pages that span a given address for a length.
<code>io_ctob(x)</code>	Return number of bytes in <i>x</i> I/O pages (rounded up).
<code>io_ctobt(x)</code>	Return number of bytes in <i>x</i> I/O pages (truncated).

The functions summarized in Table 9-12 are also provided as functions.

**Table 9-12** Functions to Convert Bytes to Sectors or Pages

Function Name	Header Files	Can Sleep?	Purpose
btop(D3)	ddi.h	N	Return number of I/O pages in a byte count (truncate).
btopr(D3)	ddi.h	N	Return number of I/O pages in a byte count (round up).
ptob(D3)	ddi.h	N	Convert size in I/O pages to size in bytes.

Using these functions and macros, you can make your driver independent of the size of pages. When examining an existing driver, be alert for any assumption that a virtual memory page has a particular size, or that an I/O page is the same size as a memory page, and convert the code to use portable functions and macros.

### Setting Up a DMA Transfer

There are two issues in preparing a DMA transfer:

- calculating physical addresses of the memory targets to be programmed into the device registers
- ensuring cache coherency in a uniprocessor

The functions you use to derive target addresses are different for different bus adapters and are discussed in the following chapters:

- The functions to set up DMA from a SCSI device are covered in Chapter 13, “SCSI Device Drivers.”
- The functions to set up DMA from a PCI device are covered in Chapter 15, “PCI Device Drivers.”

### Converting Virtual Addresses to Physical

There are almost no legitimate reasons for a device driver to convert a kernel virtual memory address to a physical address in IRIX 6.3 for O2 or any following release. All systems that support DMA, support the creation of DMA maps. A DMA map represents the mapping between a physical memory address and a bus virtual address. You initialize the map with a virtual buffer address. From the map you get a temporary physical address that you can program into a bus master for DMA. There is never a need for a driver to perform general translation from virtual to physical.

Previous releases of IRIX for simpler hardware supported functions `kvtophys()` and `sgset()` that returned physical addresses of buffer memory. If you find a use of these functions in an old driver, convert the driver to use DMA maps.

The function summarized in Table 9-13 can be used to get a physical address of kernel memory.

**Table 9-13** Functions Related to Physical Memory

Function Name	Header Files	Can Sleep?	Purpose
<code>kvtophys(D3)</code>	<code>ddi.h</code>	N	Get physical address of kernel data

### Managing Buffer Virtual Addresses

Functions to manipulate buffer page mappings are summarized in Table 9-14.

**Table 9-14** Functions to Map Buffer Pages

Function Name	Header Files	Can Sleep?	Purpose
<code>bp_mapin(D3)</code>	<code>buf.h</code>	Y	Map buffer pages into kernel virtual address space.
<code>bp_mapout(D3)</code>	<code>buf.h</code>	N	Release mapping of buffer pages.
<code>clrbuf(D3)</code>	<code>buf.h</code>	N	Clear the memory described by a mapped-in <code>buf_t</code> .
<code>userdma(D3)</code>	<code>buf.h</code>	Y	Bring pages of a buffer in user virtual address space into kernel memory and lock down.
<code>undma(D3)</code>	<code>buf.h</code>	N	Unlock pages locked by <code>userdma()</code> .

When a `pxstrategy()` routine receives a `buf_t` that is not mapped into memory (see “Buffer Location and `b_flags`” on page 187), it must make sure that the pages of the buffer space are in memory, and it must obtain valid kernel virtual addresses to describe the pages. The simplest way is to apply the `bp_mapin()` function to the `buf_t`. This function allocates a contiguous range of page table entries in the kernel address space to describe the buffer, creating a mapping of the buffer pages to a contiguous range of kernel virtual addresses. It sets the virtual address of the first data byte in `b_un.b_addr`, and sets the flags so that `BP_ISMAPPED()` returns true—thus converting an unmapped buffer to a mapped case.

### Managing Memory for Cache Coherency

Some kernel functions used for ensuring cache coherency are summarized in Table 9-15.

**Table 9-15** Functions Related to Cache Coherency

Function Name	Header Files	Can Sleep?	Purpose
<code>dki_dcache_inval(D3)</code>	<code>system.h</code> & <code>types.h</code>	N	Invalidate the data cache for a given range of virtual addresses.
<code>dki_dcache_wb(D3)</code>	<code>system.h</code> & <code>types.h</code>	N	Write back the data cache for a given range of virtual addresses.
<code>dki_dcache_wbinval(D3)</code>	<code>system.h</code> & <code>types.h</code>	N	Write back and invalidate the data cache for a given range of virtual addresses.
<code>flushbus(D3)</code>	<code>system.h</code> & <code>types.h</code>	?	Make sure contents of the write buffer are flushed to the system bus

The functions for cache invalidation are essential when doing DMA on a uniprocessor. They cost very little to use in a multiprocessor, so it does no harm to call them in every system. You call them as follows:

- Call `dki_dcache_inval()` prior to doing DMA input. This ensures that when you refer to the received data, it will be loaded from real memory.
- Call `dki_dcache_wb()` prior to doing DMA output. This ensures that the latest contents of cache memory are in system memory for the device to load.
- Call `dki_dcache_wbinval()` prior to a device operation that samples memory and then stores new data.

The **flushbus()** function is needed because in some systems the hardware collects output data and writes it to the bus in blocks. When you write a small amount of data to a device through PIO, delay, then write again, the writes could be batched and sent to the device in quick succession. Use **flushbus()** after PIO output when it is followed by PIO input from the same device. Use it also between any two PIO outputs when the device is supposed to see a delay between outputs.

### DMA Buffer Alignment

In some systems, the buffers used for DMA must be aligned on a boundary the size of a *cache line* in the current CPU. Although not all system architectures require cache alignment, it does no harm to use cache-aligned buffers in all cases. The size of a cache line varies among CPU models, but if you obtain a DMA buffer using the **KMEM\_CACHEALIGN** flag of **kmem\_alloc()**, the buffer is properly aligned. The buffer returned by **getebk()** (see “Allocating buf\_t Objects and Buffers” on page 192) is cache-aligned.

Why is cache alignment necessary? Suppose you have a variable, *X*, adjacent to a buffer you are going to use for DMA write. If you invalidate the buffer prior to the DMA write, but then reference the variable *X*, the resulting cache miss brings part of the buffer back into the cache. When the DMA write completes, the cache is stale with respect to memory. If, however, you invalidate the cache after the DMA write completes, you destroy the value of the variable *X*.

### Maximum DMA Transfer Size

The maximum size for a single DMA transfer is set by the system tuning variable *maxdmasz*, settable with the *systune* command (see the *systune(1)* reference page). A single I/O operation larger than this produces the error **ENOMEM**.

The unit of measure for *maxdmasz* is the page, which varies with the kernel. Under IRIX 6.2, a 32-bit kernel uses 4 KB pages while a 64-bit kernel uses 16 KB pages. In both systems, *maxdmasz* is shipped with the value 1024 decimal, equivalent to 4 MB in a 32-bit kernel and 16 MB in a 64-bit kernel.

In Challenge and Onyx systems, *maxdmasz* can be set as high as 64 MB. However, it is not usually possible to allocate a DMA map for a single transfer that large.

## User Process Administration

The kernel supplies a small group of functions, summarized in Table 9-16, that help a driver upper-half routine learn about the current user process.

**Table 9-16** Functions for User Process Management

Function Name	Header Files	Can Sleep?	Purpose
<code>drv_getparm(D3)</code>	<code>ddi.h</code>	N	Retrieve kernel state information.
<code>drv_priv(D3)</code>	<code>ddi.h</code>	N	Test for privileged user.
<code>drv_setparm(D3)</code>	<code>ddi.h</code>	N	Set kernel state information.
<code>proc_ref(D3)</code>	<code>ddi.h</code>	N	Obtain a reference to a process for signaling.
<code>proc_signal(D3)</code>	<code>ddi.h</code> & <code>signal.h</code>	N	Send a signal to a process.
<code>proc_unref(D3)</code>	<code>ddi.h</code>	N	Release a reference to a process.

**Note:** In older drivers you may find direct reference to a user structure. That is no longer available. Any reference to a user structure should be eliminated or replaced by one of the functions in Table 9-16.

Use `drv_getparm()` to retrieve certain miscellaneous bits of information including the process ID of the current process. In a character device driver, the current process is the user process that caused entry to the driver, for example by calling the `open()`, `ioctl()`, or `read()` system functions. In a block device driver, the current process has no direct relationship to any particular user; it is usually a daemon process of some kind.

The `drv_setparm()` function is primarily of use to terminal drivers.

The `drv_priv()` function tests a `cred_t` object to see if it represents a privileged user. A `cred_t` object is passed in to several driver entry points, and the address of the current one can be retrieved `drv_getparm()`.

## Sending a Process Signal

In traditional UNIX kernels, a device driver identified the current user process by the address of the *proc\_t* structure that the kernel uses to represent a process. Direct use of the *proc\_t* is no longer supported by IRIX. The reason is that the contents of the *proc\_t* change from release to release, and also differ between 64-bit and 32-bit kernels.

The most common use of the *proc\_t* by a driver was to send a signal to the process. This capability is still supported. To do it, take three steps:

1. Call **proc\_ref()** to get a process handle, a number unique to the current process. The returned value must be treated as an arbitrary number (in some releases of IRIX it was the *proc\_t* address, but this is not the defined behavior of the function.)
2. Use the process handle as an argument to **proc\_signal()**, sending the signal to the process.
3. Release the process handle by calling **proc\_unref()**.

The third step is important. In order to keep the process handle valid, IRIX retains information about the process to which it is related. However, that process could terminate (possibly as a result of the signal the driver sends) but until the driver announces that it is done with the handle, the kernel must try to retain process information.

It is especially important to release a process handles before unloading a loadable driver (see “Entry Point unload()” on page 170).

## Waiting and Mutual Exclusion

The kernel supplies a rich variety of functions for waiting and for mutual exclusion. In order to use these features well, you must understand the different purposes for which they are designed. In particular, you must clearly understand the distinction between *waiting* and *mutual exclusion* (or locking).

**Note:** These waiting and mutual exclusion functions have been expanded significantly in IRIX release 6.2.

## Mutual Exclusion Compared to Waiting

*Mutual exclusion* allows one entity to have exclusive use of a global resource, temporarily denying use of the resource to other entities. When software is well-designed, mutual exclusion normally does not require waiting—the resource is normally free when it is requested. A driver that calls a mutual exclusion function *expects to proceed* without delay—although there is a chance that the resource is in use, and the driver will have to wait.

The kernel offers an array of functions for mutual exclusion, and the choice among them can be critical to performance. The functions are reviewed in the following topics:

- “Basic Locks” on page 208 covers basic locks and mutex locks, the best locks for multiprocessor use.
- “Long-Term Locks” on page 210 covers sleep locks, which can be held for longer periods.
- “Reader/Writer Locks” on page 213 covers a class of locks that allow multiple, concurrent, read-only access to resources that are infrequently changed.
- “Priority Level Functions” on page 215 covers the traditional UNIX method of mutual exclusion, the `splhi()` and `splx()` functions, which have many disadvantages.

*Waiting* allows a driver to coordinate its actions with a specific event or action that occurs asynchronously. A driver can wait for a specified amount of time to pass, wait for an I/O action to complete, and so on. Therefore, a driver that calls a waiting function *expects to wait* for something to happen—although there is a chance that the expected event has already happened, and the driver will be able to continue at once.

The kernel offers several functions that allow you to wait for specific events; and also offers functions for general synchronization. These are covered in the following topics:

- “Waiting for Time to Pass” on page 216 covers timer-related functions.
- “Waiting for Memory to Become Available” on page 218 covers memory allocation waits.
- “Waiting for Block I/O to Complete” on page 219 covers waits used in the `pxstrategy()` entry point.
- “Waiting for a General Event” on page 221 covers the general-purpose functions that you can adapt to any synchronization problem.

The most general facility, the semaphore, can be used for synchronization and for locking. This topic is covered under “Semaphores” on page 224.

## Basic Locks

IRIX supports basic locks using functions compatible with SVR4. These functions are summarized in Table 9-17.

**Table 9-17** Functions for Basic Locks

Function Name	Header Files	Can Sleep?	Purpose
LOCK(D3)	ksynch.h & types.h	Y	Acquire a basic lock, waiting if necessary.
LOCK_ALLOC(D3)	ksynch.h, kmem.h & types.h	Y	Allocate and initialize a basic lock.
LOCK_DEALLOC(D3)	ksynch.h & types.h	N	Deallocate an instance of a basic lock.
LOCK_INIT(D3)	ksynch.h & types.h	N	Initialize a basic lock that was allocated statically, or reinitialize an allocated lock.
LOCK_DESTROY(D3)	ksynch.h & types.h	N	Uninitialize a basic lock that was allocated statically.
TRYLOCK(D3)	types.h & ksynch.h	N	Try to acquire a basic lock, returning a code if the lock is not currently free.
UNLOCK(D3)	types.h & ksynch.h	N	Release a basic lock.

Basic locks are objects of type *lock\_t*. Although functions are provided for allocating and freeing them, a basic lock is a very small object. Locks are typically allocated as global variables or as fields of structures.

Call `LOCK()` to seize a lock and gain possession of the resource for which it stands. Release the lock with `UNLOCK()`. These functions are optimized for mutual exclusion in the available hardware, and may be implemented differently in uniprocessors and multiprocessors. However, the programming and binary interface is the same in all systems.

The code in Example 9-1 illustrates the use of `LOCK` and `UNLOCK` in implementing a simple last-in-first-out (LIFO) queueing package. In these functions, the time between locking a queue head and releasing it is only a few microseconds.

**Example 9-1** LIFO Queue Using Basic Locks

```
typedef struct qitem {
    qitem *next; ...other fields...
} qitem_t;
typedef struct lifo {
    qitem *latest;
    lock_t grab;
} lifo_t;
void putlifo(lifo_t *q, qitem_t *i)
{
    int lockpl = LOCK(&q->grab,plhi);
    i->next = q->latest;
    q->latest = i;
    UNLOCK(&q->grab,lockpl);
}
qitem_t *poplifo(lifo_t *q)
{
    int lockpl = LOCK(&q->grab,plhi);
    qitem_t *ret = q->latest;
    q->latest = ret->next;
    UNLOCK(&q->grab,lockpl);
    return ret;
}
```

This is a typical use of basic locks: to ensure that for a brief period only one process in the system can update a queue. Basic locks are optimized for such uses, but in order to get optimal performance they are restricted to these uses. In particular, if you seize a basic lock and hold it over a function call that can sleep, the system can deadlock.

## Long-Term Locks

Sometimes you need a lock that can be held for a longer period, over a call to a function that can sleep. IRIX provides three types of such locks: mutex locks, sleep locks, and reader-writer locks.

### Using Mutex Locks

Mutex locks are designed for mutual exclusion (as the name suggests). The IRIX implementation of mutex locks is compatible with the *kmutex\_t* lock type of SunOS™, but optimized for use in Silicon Graphics hardware systems. The mutex functions are summarized in Table 9-18.

**Table 9-18** Functions for Mutex Locks

Function Name	Header Files	Can Sleep?	Purpose
MUTEX_ALLOC(D3)	types.h & kmem.h & ksynch.h	Y	Allocate and initialize a mutex lock.
MUTEX_INIT(D3)	types.h & ksynch.h	N	Initialize an existing mutex lock.
MUTEX_DESTROY(D3)	types.h & ksynch.h	N	Deinitialize a mutex lock.
MUTEX_DEALLOC(D3)	types.h & ksynch.h	N	Deinitialize and free a dynamically allocated mutex lock.
MUTEX_LOCK(D3)	types.h & kmem.h & ksynch.h	Y	Claim a mutex lock.
MUTEX_TRYLOCK(D3)	types.h & ksynch.h	N	Conditionally claim a mutex lock.
MUTEX_UNLOCK(D3)	types.h & ksynch.h	N	Release a mutex lock.
MUTEX_WAITQ(D3)	types.h & ksynch.h	N	Get the number of processes blocked by mutex lock.

**Table 9-18 (continued)** Functions for Mutex Locks

Function Name	Header Files	Can Sleep?	Purpose
MUTEX_ISLOCKED(D3)	types.h & ksynch.h	N	Test if a mutex lock is owned.
MUTEX_MINE(D3)	types.h & ksynch.h	N	Test if a mutex lock is owned by this process.

Although allocation and deallocation functions are supplied, a *mutex\_t* type is a small object that is normally allocated as a static variable or as a field of a structure. The `MUTEX_INIT()` operation prepares a statically-allocated *mutex\_t* for use.

Once initialized, a mutex lock is used to gain exclusive use of the resource with which you have associated it. The mutex lock has the following important advantages over a basic lock:

- The mutex lock can safely be held over a call to a function that sleeps.
- The mutex lock supports the inquiry functions `MUTEX_WAITQ`, `MUTEX_ISLOCKED`, and `MUTEX_MINE`.
- When a debugging kernel is used (see “Including Lock Metering in the Kernel Image” on page 248) a mutex lock can be instrumented to keep statistics of its use.

The mutex lock implementation provides *priority inheritance*. When a low-priority process owns a mutex lock and a high-priority process attempts to seize the lock and is blocked, the process holding the lock is temporarily given the higher priority of the blocked process. This hastens the time when the lock can be released, so that a low-priority process does not needlessly impede a higher-priority process.

In order to implement priority inheritance and retain high performance, the mutex lock is subject to the following restrictions:

- A mutex lock can be locked and unlocked only by an upper-half driver routine; that is, from code that has a process context. A mutex lock cannot be locked or unlocked in an interrupt routine.
- A mutex lock must be unlocked by the same process that locked it. It cannot be locked in one process identity and unlocked in another.

Because of these restrictions, a mutex lock can only be used to mediate between upper-half driver entry points. It is very effective for this purpose; you can use mutex locks to coordinate the use of global variables between upper-half entry points of a driver in a multiprocessor design.

When you need mutual exclusion between an upper-half entry point and the interrupt handler, use a basic lock. Resources that are shared with an interrupt handler should never be in use for more than a brief period. When your design requires a lock to be seized by one process and released in another, use a sleep lock or semaphore.

### Using Sleep Locks

IRIX supports sleep lock functions that are compatible with SVR4. These functions are summarized in Table 9-19.

**Table 9-19** Functions for Sleep Locks

Function Name	Header Files	Can Sleep?	Purpose
SLEEP_ALLOC(D3)	types.h & kmem.h & ksynch.h	Y	Allocate and initialize a sleep lock.
SLEEP_DEALLOC(D3)	types.h & ksynch.h	N	Deinitialize and deallocate a dynamically allocated sleep lock.
SLEEP_INIT(D3)	types.h & ksynch.h	N	Initialize an existing sleep lock.
SLEEP_DESTROY(D3)	types.h & ksynch.h	N	Deinitialize a sleep lock.
SLEEP_LOCK(D3)	types.h & ksynch.h & param.h	Y	Acquire a sleep lock, waiting if necessary until the lock is free.
SLEEP_LOCKAVAIL(D3)	types.h & ksynch.h	N	Query whether a sleep lock is available.
SLEEP_LOCK_SIG(D3)	types.h & ksynch.h & param.h	Y	Acquire a sleep lock, waiting if necessary until the lock is free or a signal is received.

**Table 9-19 (continued)** Functions for Sleep Locks

Function Name	Header Files	Can Sleep?	Purpose
SLEEP_TRYLOCK(D3)	types.h & ksynch.h	N	Try to acquire a sleep lock, returning a code if it is not free.
SLEEP_UNLOCK(D3)	types.h & ksynch.h	N	Release a sleep lock.

Although allocation and deallocation functions are supplied, a *sleep\_t* type is a small object that is normally allocated as a static variable or as a field of a structure. The SLEEP\_INIT() operation prepares a statically-allocated *sleep\_t* for use. (In IRIX 6.2, a *sleep\_t* is identical to a *sema\_t*, but this situation could change in a future release.)

A sleep lock is similar to a mutex lock in that it is used for mutual exclusion between processes, and can be held across a function call that sleeps. A sleep lock does not have either the advantages or the restrictions of a mutex lock:

- A sleep lock can be seized by one process and released by another.
- A sleep lock can be set in an upper-half entry point and released in an interrupt routine.
- A sleep lock does not provide priority inheritance. When a low-priority process holds a sleep lock, a higher-priority process can be blocked, causing a *priority inversion*.
- A sleep lock does not support the instrumentation or the query functions supported for mutex locks.

## Reader/Writer Locks

Reader/writer locks are similar to sleep locks in that they are designed for mutually exclusive control of resources for relatively long periods of time. However, Reader/Writer locks are optimized for the case in which the resource is often used by processes that only interrogate it (readers), and only rarely used by processes that modify it (writers).

Reader/writer locks compatible with SVR4 are introduced in IRIX 6.2. The functions are summarized in Table 9-20.

**Table 9-20** Functions for Reader/Writer Locks

Function Name	Header Files	Can Sleep?	Purpose
RW_ALLOC(D3)	types.h & kmem.h & ksynch.h	Y	Allocate and initialize a reader/writer lock.
RW_DEALLOC(D3)	types.h & ksynch.h	N	Deallocate a reader/writer lock.
RW_INIT(D3)	types.h & ksynch.h	N	Initialize an existing reader/writer lock.
RW_DESTROY(D3)	types.h & ksynch.h	N	Deinitialize an existing reader/writer lock.
RW_RDLOCK(D3)	types.h & ksynch.h & param.h	Y	Acquire a reader/writer lock as reader, waiting if necessary.
RW_TRYRDLOCK(D3)	types.h & ksynch.h	N	Try to acquire a reader/writer lock as reader, returning a code if it is not free.
RW_TRYWRLOCK(D3)	types.h & ksynch.h	N	Try to acquire a reader/writer lock as writer, returning a code if it is not free.
RW_UNLOCK(D3)	types.h & ksynch.h	N	Release a reader/writer lock as reader or writer.
RW_WRLOCK(D3)	types.h & ksynch.h & param.h	Y	Acquire a reader/writer lock as writer, waiting if necessary.

Although allocation and deallocation functions are supplied, a *mrlock\_t* type is a small object that is normally allocated as a static variable or as a field of a structure. The `RW_INIT()` operation prepares a statically-allocated *mrlock\_t* for use.

A process that intends to modify a resource uses `RW_WRLOCK` to claim it. This process waits until the resource is not in use by any process, then it gains exclusive access. Only one process is allowed to hold a reader/writer lock as a writer. All other processes, readers or writers, wait until the writer releases the lock.

A process that intends only to interrogate a resource uses `RW_RDLOCK` to gain access. If a writer holds the lock, the process waits. When the lock is free, or is held only by other readers, the process continues. More than one reader can hold a reader/writer lock at one time. It is also valid for a reader to “double-trip” a reader/writer lock; that is, claim it two or more times. The reader must release the lock as many times as it claimed the lock.

A reader/writer lock serves the same basic purpose as a sleep lock, but it is more efficient in a multiprocessor when there are frequent, read-only uses of a resource.

### Priority Level Functions

In traditional UNIX systems, one set of functions served all purposes of synchronization and locking: the set-priority-level, or `spl`, functions. These functions are still available in IRIX, and are summarized in Table 9-21.

**Table 9-21** Functions to Set Interrupt Levels

Function Name	Header Files	Can Sleep?	Purpose
<code>splbase(D3)</code>	<code>ddi.h</code>	N	Block no interrupts.
<code>spltimeout(D3)</code>	<code>ddi.h</code>	N	Block only timeout interrupts.
<code>spldisk(D3)</code>	<code>ddi.h</code>	N	Block disk interrupts.
<code>splstr(D3)</code>	<code>ddi.h</code>	N	Block STREAMS interrupts.
<code>spltty(D3)</code>	<code>ddi.h</code>	N	Block disk, VME, serial interrupts.
<code>splhi(D3)</code>	<code>ddi.h</code>	N	Block all I/O interrupts.
<code>spl0(D3)</code>	<code>ddi.h</code>	N	Same as <code>splbase()</code> .
<code>splx(D3)</code>	<code>ddi.h</code>	N	Restore previous interrupt level.

These functions are commonly found in device drivers being ported from uniprocessors. Such drivers rely on the use of `splhi()` to gain exclusive use of a global resource.

The `spl` functions are supported by IRIX and they are effective in a uniprocessor driver. However, in a multiprocessor, the functions affect only the interrupt handling of the current CPU. Other CPUs in the system continue to handle interrupts, including interrupts initiated by the driver that called `splhi()`.

A driver that is not multiprocessor-aware (one that does not have `D_MP` in its `pf_xdevflag` constant; see “Driver Flag Constant” on page 145) runs only in CPU 0 of a multiprocessor, so in this case the `spl` functions are still effective. Since they set the interrupt level on CPU 0 where the driver runs, and since the driver’s interrupts can only be handled on CPU 0, the use of `splhi(0)` gives the driver exclusive use of its resources.

A driver that is multiprocessor-aware uses basic locks, synchronization variables, and other tools to control access to resources, and never uses an `spl` function. This improves performance in a multiprocessor, does not harm performance in a uniprocessor, and reduces the latency of all interrupts.

## Waiting for Time to Pass

The kernel offers functions for timed delays, as summarized in Table 9-22.

**Table 9-22** Functions for Timed Delays

Function Name	Header Files	Can Sleep?	Purpose
<code>delay(D3)</code>	<code>ddi.h</code>	Y	Delay for a specified number of clock ticks.
<code>drv_hztousec(D3)</code>	<code>ddi.h</code>	N	Convert clock ticks to microseconds
<code>drv_usectohz(D3)</code>	<code>ddi.h</code>	N	Convert microseconds to clock ticks.
<code>drv_usecwait(D3)</code>	<code>ddi.h</code>	N	Busy-wait for a specified interval.
<code>dtimeout(D3)</code>	<code>ddi.h</code> & <code>ksynch.h</code>	N	Schedule a function execute on a specified processor after a specified length of time.
<code>itimeout(D3)</code>	<code>ddi.h</code> & <code>ksynch.h</code>	N	Schedule a function to be executed after a specified number of clock ticks.
<code>timeout(D3)</code>	<code>ddi.h</code> & <code>ksynch.h</code>	N	Schedule a function to be executed after a specified number of clock ticks.
<code>untimeout(D3)</code>	<code>ddi.h</code>	N	Cancel a previous <code>itimeout</code> or <code>fast_itimeout</code> request.
<code>untimeout_func(D3)</code>	<code>ddi.h</code>	N	Cancel a previous <code>itimeout</code> or <code>fast_itimeout</code> request by function name.

### Time Units

The basic time unit is the “tick.” Its value can differ between hardware platforms and between versions of IRIX. The **drvhztousec()** and **drvusectohz()** functions convert between ticks and microseconds in the current system. Use them in order to schedule a delay in a portable manner. (However, the timer function precision is the tick, not the microsecond.)

### Timer Support

Timer support is based on the idea of a “callback” function. You specify the following to **dtimeout()**, **itimeout()**, **timeout()** or **fast\_itimeout()**:

- an interval in clock ticks or fast ticks
- a function to be called at the expiration of the interval
- one or more arguments to be passed to the function
- a priority (interrupt) level at which the function should run

After a delay of at least the length requested, the function is called. The function is entered asynchronously. On a uniprocessor, it can interrupt execution of an upper-half routine. On a multiprocessor, it can execute concurrently with an upper-half routine or with an interrupt handler. You should not rely on the priority level of the function for mutual exclusion (see “Priority Level Functions” on page 215 for an explanation).

The difference between **itimeout()** and **timeout()** is that the latter takes no argument values to be passed to the function when it is called. In order to get a repeated series of timer events, start a new timeout from the callback function.

The **untimeout()** and **untimeout\_func()** functions cancel a pending timeout. In a loadable driver that has an **pfxunload()** entry point, cancel any pending timeouts before unloading.

The **STREAMS\_TIMEOUT** macro supplies similar timeout capability for a **STREAMS** driver (see “Special Considerations for Multiprocessing” on page 505).

### Short-Term Delay Support

In rare circumstances, a driver needs to pause briefly between two hardware operations. For example, the Silicon Graphics support for external interrupts in the Challenge and Onyx computers sometimes needs to set a high output level, wait for a brief, precise interval, then set a low output level.

The `drv_usecwait()` function supports this type of very short, precisely-timed delay. It “spins” for a specified number of microseconds, then returns to the caller. The CPU does nothing else during this period, so clearly a delay of more than a few microseconds can interfere with other work. Furthermore, if interrupts are disabled during the wait, the response to another interrupt is delayed also—the delay contributes directly to the “latency” of interrupt handling.

### Waiting for Memory to Become Available

Whenever you request memory of any kind, you must allow for the possibility that the memory will not be available. When you allocate memory in bulk (see “General-Purpose Allocation” on page 190) using `kmem_alloc()` you have the option of receiving a null response, or of waiting for the memory to be available.

When you request memory for specific object types (see “Allocating Objects of Specific Kinds” on page 191) there is usually no choice; the functions sleep until they can acquire an object of the requested type.

Within a STREAMS driver you have the ability to schedule a callback function to be entered when memory for a message buffer becomes available (see the `bufcall(D3)` reference page).

## Waiting for Block I/O to Complete

The `pxstrategy()` routine initiates the I/O operation to fill a buffer based on a `buf_t` structure. Then it has to wait for the I/O to complete. The functions for managing this synchronization are summarized in Table 9-23.

**Table 9-23** Functions for Synchronizing Block I/O

Function Name	Header Files	Can Sleep?	Purpose
<code>biodone(D3)</code>	<code>ddi.h</code>	N	Release buffer after I/O and wake up waiting process.
<code>bioerror(D3)</code>	<code>ddi.h</code>	N	Manipulate error fields in a <code>buf_t</code> .
<code>biowait(D3)</code>	<code>ddi.h</code>	Y	Suspend process pending completion of I/O.
<code>geterror(D3)</code>	<code>ddi.h</code>	N	Retrieve error number from a <code>buf_t</code> .
<code>physiock(D3)</code>	<code>ddi.h</code>	Y	Validate a raw I/O request and pass to a strategy function.
<code>uiophysio(D3)</code>	<code>ddi.h</code>	Y	Validate a raw I/O request and pass to a strategy function.
<code>undma(D3)</code>	<code>ddi.h</code>	?	Unlock physical memory after I/O complete
<code>userdma(D3)</code>	<code>ddi.h</code>	?	Lock physical memory in user space. small number of

### How the `strategy()` Entry Point Is Called

The `pxstrategy()` entry point is called directly from the filesystem or virtual memory management, or it can be called indirectly from a `pxread()` or `pxwrite()` entry point (see “Calling Entry Point `strategy()` From Entry Point `read()` or `write()`” on page 156).

### Strategies of the `strategy()` Entry Point

Typically the `pxstrategy()` routine must interact with its interrupt handler. The `pxstrategy()` routine can be designed in either of two ways, synchronous or asynchronous.

The synchronous *pxstrategy()* routine initiates every I/O operation. Its interrupt handler is responsible only for detecting and signalling the completion of one I/O. The *pxstrategy()* routine proceeds as follows:

1. Lock the data buffer in memory using **userdma()**.
2. Place the address of the *buf\_t* where the *pxintr()* entry point can find it.
3. Program the device (see “Setting Up a DMA Transfer” on page 201) and initiate the I/O activity.
4. Call **biowait()**.

When the interrupt handler is entered, the handler uses **bioerror()** if necessary, and **biodone()** to signal the completion of the I/O. Then it exits. The strategy code, which is waiting in the call to **biowait()**, regains control following the call to **biodone()**, and can use **geterror()** to check the results.

The asynchronous *pxstrategy()* routine only initiates the first I/O operation of a series, and never waits. It proceeds as follows:

1. Lock the data buffer in memory using **userdma()**.
2. Append the address of the *buf\_t* to a queue shared with the interrupt handler.
3. If the queue was empty, no I/O is in progress. Call a subroutine that programs the device and initiates the I/O.
4. Return to the caller. The caller (a filesystem or paging system or **uiophysio()**) waits using **biowait()**.

When the interrupt occurs, the handler proceeds as follows:

1. The first queued *buf\_t* has completed. Remove it from the queue.
2. Apply **bioerror()** if necessary, and **biodone()** to the *buf\_t*. This releases the caller of the strategy routine from **biowait()**.
3. If any operations remain in the queue, call a subroutine to program and initiate the next one.

## Waiting for a General Event

There are causes for synchronization other than time, block I/O, and memory allocation. For example, there is no defined interface comparable to **biowait()**/**biodone()** to mediate between an interrupt handler and the *pxread()* or *pxwrite()* entry points. You must design a mechanism of your own, using either a synchronization variable or the **sleep()**/**wakeup()** function pair.

### Using **sleep()** and **wakeup()**

The **sleep()** and **wakeup()** function pair are the simplest, oldest, and least efficient of the general synchronization mechanisms. They are summarized in Table 9-24.

**Table 9-24** Functions for Synchronization: *sleep/wakeup*

Function Name	Header Files	Can Sleep?	Purpose
<i>sleep(D3)</i>	ddi.h & param.h	Y	Suspend execution pending an event.
<i>wakeup(D3)</i>	ddi.h	N	Waken a process waiting for an event.

Used carefully, these functions are suitable for simple character device drivers. However, when you are writing new code or converting a driver to multiprocessing you should avoid them and use synchronization variables instead (see “Using Synchronization Variables” on page 222).

The basic concept is that the upper-layer routine calls **sleep(*n*)** in order to wait for an event that is keyed to an arbitrary address *n*. Typically *n* is a pointer to a data structure related to an I/O operation. The interrupt handler executes **wakeup(*n*)** to cause the sleeping process to resume execution.

The main reason to avoid **sleep()** is that, in a multiprocessor system, it is hard to ensure that sleeping always begins before **wakeup()** is called. The usual intended sequence of events is as follows:

1. Upper-half routine initiates a device operation that will lead to an interrupt.
2. Upper-half routine executes **sleep(*n*)**.
3. Interrupt occurs, and handler executes **wakeup(*n*)**.

In a multiprocessor-aware driver (one with `D_MP` in its `pf_xdevflag` constant; see “Driver Flag Constant” on page 145), there is a small chance that the interrupt can occur, calling `wakeup(n)`, before the `sleep(n)` call has been completed. Because `sleep(0)` has not been called, the `wakeup(0)` is lost. When the `sleep(0)` call completes, the process sleeps forever. Synchronization variables are designed to handle this case.

### Using Synchronization Variables

Synchronization variables, a feature of UNIX SVR4, are supported by IRIX beginning with release 6.2. These functions are summarized in Table 9-25.

**Table 9-25** Functions for Synchronization: Synchronization Variables

Function Name	Header Files	Can Sleep?	Purpose
SV_ALLOC(D3)	types.h & sema.h	Y	Allocate and initialize a synchronization variable.
SV_DEALLOC(D3)	types.h & sema.h	N	Deinitialize and deallocate a synchronization variable.
SV_INIT(D3)	types.h & sema.h	N	Initialize an existing synchronization variable.
SV_DESTROY(D3)	types.h & sema.h		Deinitialize a synchronization variable.
SV_BROADCAST(D3)	types.h & sema.h	N	Wake all processes sleeping on a synchronization variable.
SV_SIGNAL(D3)	types.h & sema.h	N	Wake one process sleeping on a synchronization variable.
SV_WAIT(D3)	types.h & sema.h	Y	Sleep until a synchronization variable is signalled.
SV_WAIT_SIG(D3)	types.h & sema.h	Y	Sleep until a synchronization variable is signalled or a signal is received.

A synchronization variable is a memory object of type `sv_t`, representing the occurrence of an event. You can allocate objects of this type dynamically, or declare them as static variables or as fields of structures.

One or more processes may wait for an event using `SV_WAIT()`. An interrupt handler or timer callback function can signal the occurrence of an event using `SV_SIGNAL` (to wake up only one waiting process) or `SV_BROADCAST` (to wake up all of them).

`SV_WAIT` is specifically designed to handle the difficult case that arises when the driver needs to initiate an I/O operation and then sleep, and do these things in such a way that it always begins to sleep before the `SV_SIGNAL` can possibly be issued. The procedure is done as follows:

1. The driver seizes a basic lock (see “Basic Locks” on page 208) or a mutex lock (see “Using Mutex Locks” on page 210) that is also used by the interrupt handler.  
A `LOCK()` call returns an integer that is needed later.
2. The driver initiates an I/O operation that can lead to an interrupt.
3. The driver calls `SV_WAIT`, passing the lock it holds and an integer, either the value returned by `LOCK()` or a zero if the lock is a mutex lock.
4. In one indivisible operation, `SV_WAIT` releases the lock and begins waiting on the synchronization variable.
5. The interrupt handler or other process is entered, and seizes the lock.  
This step ensures that, if the interrupt handler or other process is entered preceding the `SV_WAIT` call, it will not proceed until `SV_WAIT` has completed.
6. The interrupt handler or other process does its work and calls `SV_SIGNAL` to release the waiting driver.

This process is sketched in Example 9-2.

**Example 9-2** Skeleton Code for Use of `SV_WAIT`

```
lock_t seize_it;
sv_t wait_on_it;
initiator(...)
{
    int lock_cookie;
    for( as often as necessary )
    {
        lock_cookie = LOCK(&seize_it, PL_ZERO);
        [do something that causes a later interrupt]
        SV_WAIT(&wait_on_it, 0, &seize_it, lock_cookie);
        [interrupt has been handled]
    }
}
```

```
void handler(...)  
{  
    int lock_cookie = LOCK(&seize_it, PL_ZERO);  
    [handle the interrupt]  
    SV_SIGNAL(&seize_it);  
    UNLOCK(&seize_it);  
}
```

If it is necessary to use a semaphore as the lock, the header file `sys/semaphore.h` declares versions of `SV_WAIT` that accept a semaphore and a synchronization variable. The combination of a mutual exclusion object and a synchronization variable ensures that even in a multiprocessor, the interrupt handler cannot exit before the driver has entered a predictable wait state.

**Tip:** When a debugging kernel is used, you can display statistics about the use of a given synchronization variable. See “Including Lock Metering in the Kernel Image” on page 248.

## Semaphores

The *semaphore* is a generalized tool that can be used for both mutual exclusion and for waiting. The IRIX kernel support for semaphores is summarized in Table 9-26.

**Table 9-26** Functions for Semaphores

Function Name	Header Files	Can Sleep?	Purpose
<code>cpsema(D3)</code>	<code>sema.h</code> & <code>types.h</code>	N	Conditionally perform a “P” or wait semaphore operation.
<code>cvsema(D3)</code>	<code>sema.h</code> & <code>types.h</code>	N	Conditionally perform a “V” or release semaphore operation.
<code>freeseма(D3)</code>	<code>sema.h</code> & <code>types.h</code>	N	Free the resources associated with a semaphore.
<code>initnsema(D3)</code>	<code>sema.h</code> & <code>types.h</code>	N	Initialize a semaphore to a given value.
<code>initnsema_mutex(D3)</code>	<code>sema.h</code> & <code>types.h</code>	N	Initialize a semaphore to a value of 1.

**Table 9-26 (continued)** Functions for Semaphores

Function Name	Header Files	Can Sleep?	Purpose
psema(D3)	sema.h & types.h & param.h	Y	Perform a “P” or wait semaphore operation.
valusema(D3)	sema.h & types.h	N	Return the value associated with a semaphore.
vsema(D3)	sema.h & types.h	N	Perform a “V” or signal semaphore operation.

Conceptually, a semaphore contains an integer. The “P” operation claims the semaphore, decrementing its count by 1 (mnemonic: dePlete). If the count is 0 or less, the process waits until the count is greater than 0 before it decrements the semaphore and returns.

The “V” operation increments the semaphore count (mnemonic: reViVe) and wakens any process that is waiting.

**Tip:** When a debugging kernel is used, you can display statistics about the use of a given semaphore. See “Including Lock Metering in the Kernel Image” on page 248.

**Note:** In releases before IRIX 6.2, `initnsema_mutex()` was used to initialize a semaphore in a special way that got the performance of a basic lock in a multiprocessor. Since IRIX 6.2, this function is simply a macro that initializes the semaphore to a count of 1.

### Using a Semaphore for Mutual Exclusion

To use a semaphore for locking, initialize it to 1. (This reflects the idea that a process calling a locking function expects to continue.) When you require exclusive use of the associated resource, call `psema()`. Typically this finds a semaphore count of 1, reduces it to 0, and returns.

When you are finished with the resource, call `vsema()` to increment the semaphore count, and release any process that is blocked in a `psema()` call for the same semaphore.

For locking, a semaphore is comparable to a sleep lock. In some systems, the performance of semaphore operations may not be as good as the performance of a mutex lock. In other systems, mutex locks may be implemented using semaphores.

### Using a Semaphore for Waiting

To use a semaphore for waiting, initialize it to 0. Then call **psema()**. Because the semaphore count is 0, the process waits. When the desired event occurs, typically in the interrupt handler, call **vsema()** to release the waiting process.

This synchronization method is as reliable as a synchronization variable, but it has slightly different behavior. When a synchronization variable is used correctly (see “Using Synchronization Variables” on page 222), if the interrupt handler is entered before the SV\_WAIT call completes, the interrupt handler waits on a LOCK call.

When a semaphore is used, if the interrupt handler is entered before the **psema()** call completes, the **vsema()** operation is done immediately and the interrupt handler continues without waiting. The fact that **vsema()** was called is stored as a count within the semaphore, where **psema()** will find it. Because the semaphore can contain this state information, the interrupt handler does not have to be synchronized in time using a lock.

**Note:** In releases before IRIX 6.2, the **vpsema()** function was used in a way similar to synchronization variables are used: to release one semaphore and wait on another in an atomic operation. This function is no longer supported; replace it with synchronization variable.

## Building and Installing a Driver

After a kernel-level driver has been designed and coded, it must be compiled, linked, and installed. The topics in this chapter describe the major steps of this process, as follows:

- “Defining Device Numbers” on page 228 covers the choice of major and minor device numbers.
- “Defining Device Special Files” on page 229 describes options for creating the file or files controlled by the driver.
- “Compiling and Linking” on page 230 covers the compiler and linker options used for driver modules.
- “Configuring a Nonloadable Driver” on page 234 describes the configuration files used to set up a driver loaded at boot time.
- “Configuring a Loadable Driver” on page 239 describes the additional configuration needed for a loadable driver.

## Defining Device Numbers

The topics “Major Device Number” on page 36 and “Minor Device Number” on page 37 cover the purpose and use of the device numbers, which can be summarized as follows:

- Both numbers are encoded in the inode of a device special file in */dev*.
- The major number selects the device driver.
- The minor number specifies the logical unit, and can encode device features.
- Both numbers are passed as a parameter to driver entry points.

An important part of creating and installing a device driver is the selection of device numbers and the definition of device special files.

### Selecting a Major Number

You must select a major number to stand for your driver. The numbers that already exist are listed in *sys/major.h*. However, the major number should not be coded into the driver. Typically the driver code does not need to know its major number, and if it does, the driver should discover its major number dynamically. A method of doing this is discussed under “Variables Section” on page 237.

A driver is associated with its major number in the *master.d* configuration file. When the driver discovers this number dynamically, the system administrator is free to change major numbers in */var/sysgen/master.d* files to correct conflicts between one product and another.

It is possible to let the *lboot* command choose a major number dynamically. This is discussed under “Configuring for a Dynamic Major Number” on page 243.

## Selecting Minor Numbers

Each device minor number comprises 18 bits of information that are passed to driver entry points. The format of minor numbers can be coded into the driver. You design the content of these numbers to give your driver the information it needs about each device special file.

Examine the */dev/MAKEDEV* script to see some techniques for assembling minor numbers dynamically based on the hardware inventory and other commands.

## Defining Device Special Files

As described under “Device Special Files” on page 34, the association between a device and a driver is established by encoding the major device number in the *inode* of a device special file in the */dev* directory. Without at least one device special file, a device can never be opened.

### Static Definition of Device Special Files

The system administrator can create device special files using *mknod* or *install* (see “Making Device Files” on page 38). This can be done manually, or through an installation script, or as an exit operation of the software manager program. The device special files can be created at any time—whether or not the device driver is installed, and whether or not the hardware exists. The special device files continue to exist until the administrator removes them.

### Dynamic Definition of Device Special Files

In a more sophisticated installation you might want to have the device special files created, or recreated, dynamically each time the system boots. This is the purpose of */dev/MAKEDEV* (see “The Script MAKEDEV” on page 38); it removes and redefines device special files based on information in the hardware inventory.

Much of the logic in */dev/MAKEDEV* depends on information reported by the *hinv* command. Through IRIX 6.2, there is no OEM interface for drivers to the hardware inventory (see “Hardware Inventory” on page 31), so at this time there is no opportunity for your driver to pass information through the inventory to */dev/MAKEDEV*.

However, there are other possibilities. For example, if your driver supports a control unit with an unknown number of minor devices, you could statically create a single device special file to represent the control unit. You could write the driver so that an `ioctl()` function directed to this file returns the count of actual minor devices.

## Compiling and Linking

You compile a kernel device driver to an ELF binary (in IRIX 5.3 and before, the COFF binary format was used) not using shared libraries. The compile options differ between 32-bit and 64-bit modules.

### Using `/var/sysgen/Makefile.kernio`

The file `/var/sysgen/Makefile.kernio` is a sample Makefile for compiling kernel modules. You can include it from a Makefile of your own to set the definitions of compiler variables and compiler options appropriately for different CPUs and module types.

The `Makefile.kernio` file tests the following environment variables:

<code>CPUBOARD</code>	Set to the type of CPU used in the target system, for example IP19 or IP22.
<code>COMPILATION_MODEL</code>	Set to 64 for a 64-bit kernel module, or to 32 for a 32-bit kernel module.

The purpose of the rules in `Makefile.kernio` is to set compiler variables and compiler options into a Make variable `CFLAGS`. Other Make variables would be set by your own Makefile.

**Note:** `Makefile.kernio` is designed for nonloadable drivers. In particular it sets the compiler option `-G8`, which is valid for nonloadable drivers. A loadable driver must use `-G0`.

## Compiler Variables

The compiler variables listed in Table 10-1 are tested in system header files. They are usually defined on the compiler command line. The rules in */var/sysgen/Makefile.kernio* set definitions of the variables appropriately for different CPU types.

**Table 10-1** Compiler Variables Tested by System Header Files

Variable	Meaning
<code>_K32U32</code>	Compile for a 32-bit kernel running (only) 32-bit user programs.
<code>_K32U64</code>	Compile for a 32-bit kernel running 32-bit or 64-bit user programs (not supported in IRIX 6.2).
<code>_K64U64</code>	Compile for a 64-bit kernel running 32-bit or 64-bit user programs.
<code>_KERNEL</code>	Compile for a kernel module, not a user program.
<code>STATIC=static</code>	Use of pseudo-declarator <code>STATIC</code> is converted to real <code>static</code> .
<code>JUMP_WAR</code>	Compile workaround code for R4000 branch on end of page bug.
<code>PROBE_WAR</code>	Compile workaround code for R4000 bug which requires TLBprobe instruction to be performed in uncached mode.
<code>BADVA_WAR</code>	Compile workaround code for R4000 badvaddr bug.
<code>TFP</code>	Target machine is the R8000.
<code>R4000</code>	Target machine is the R4000.
<code>R3000</code>	Target machine is the R3000 (not supported after IRIX 5.3).
<code>IP<math>nm</math></code>	Target machine uses the <code>IP<math>nm</math></code> CPU module: <code>IP17</code> , <code>IP19</code> , <code>IP20</code> , <code>IP22</code> , <code>IP26</code> are currently supported.
<code>_PAGESZ=16384</code>	Compile for a kernel using 16K memory pages (with IRIX 6.2 this implies <code>_K64U64</code> also defined).
<code>_PAGESZ=4096</code>	Compile for a kernel using 4K memory pages (with IRIX 6.2 this implies <code>_K32U32</code> also defined).
<code>_MIPS3_ADDRSPACE</code>	Kernel for a MIPS III (R8000) machine.
<code>EVEREST</code>	Compile for a Challenge or Onyx system.
<code>MP</code>	Compile for a multiprocessor.
<code>_MP_NETLOCKS</code>	Compile network driver (only) for multiprocessor.

## Compile Options, 32-Bit Kernel

The *cc* and *ld* options shown in Table 10-2 are needed to compile and link a kernel-level driver for a 32-bit kernel in IRIX 6.2.

**Table 10-2** Compiler Options for 32-Bit Kernel Modules

Option	Purpose
<i>-non_shared</i>	Do not compile for shared libraries (dynamic linking).
<i>-elf</i>	Compile and link an ELF binary.
<i>-32 -mips2</i>	Create a 32-bit executable for the MIPS II architecture.
<i>-G 8</i>	In a nonloadable driver, use the global table for objects up to 8 bytes.
<i>-G 0</i>	In a loadable driver, do not use the global table. Refer to the <code>gp_overflow(5)</code> reference page for a discussion of the global table.
<i>-O2</i>	Maximum recommended optimization level.
<i>-r</i>	Linker to retain symbols (needed by loadable drivers only).
<i>-d</i>	Force definition of common storage even though <i>-r</i> used.
<i>-Wc,-pic0</i>	Do not allocate stack space used by shared objects.
<i>-TARGET:force_jalr</i>	Generate <b>jalr</b> instructions for subroutine calls rather than <b>jal</b> , which allows only a 26-bit target and so cannot address all kernel virtual storage.

## Compile Options, 64-Bit Kernel

The *cc* and *ld* options listed in Table 10-3 are needed to compile and link a kernel-level driver for a 64-bit kernel in IRIX 6.2.

**Table 10-3** Compiler Options for 64-Bit Kernel Modules

Option	Purpose
<i>-non_shared</i>	Do not compile for shared libraries (dynamic linking).
<i>-elf</i>	Compile and link an ELF binary.
<i>-64 -mips3</i>	Create a 32-bit executable for the MIPS III architecture. You can use <i>-mips4</i> instead, but only in a system based on the R10000 processor.
<i>-G 8</i>	In a nonloadable driver, use the global table for objects up to 8 bytes.
<i>-G 0</i>	In a loadable driver, do not use the global table. Refer to the <i>gp_overflow(5)</i> reference page for a discussion of the global table.
<i>-O2</i>	Maximum recommended optimization level.
<i>-TENV:kernel</i> <i>-TENV:misalignment=1</i>	Execution environment options for 64-bit compiler.
<i>-OPT:space</i> <i>-OPT:quad_align_branch_targets=FALSE</i> <i>-OPT:quad_align_with_memops=FALSE</i> <i>-OPT:unroll_times=0</i>	Specific optimization constraints for 64-bit compiler.

## Configuring a Nonloadable Driver

When the driver is not loadable, it is linked as part of the IRIX kernel. The following steps are needed to make the driver usable:

1. Place the driver executable file in */var/sysgen/boot*.
2. Place a descriptive file in */var/sysgen/master.d*.
3. Place a directive file in */var/sysgen/system* (or simply add a line to */var/sysgen/system/irix.sm*).
4. Run *autoconfig* to generate a new kernel.
5. Reboot the system.

Some of these steps are elaborated in the following topics.

### How Names Are Used in Configuration

The naming of a kernel-level driver begins in a file in */var/sysgen/system*, such as */var/sysgen/system/irix.sm*. Names are used as follows:

- A *USE*, *INCLUDE*, or *VECTOR* statement in */var/sysgen/system* specifies a name, for example  
`USE hypothetical`
- This statement directs *lboot* to read a file of the same name in */var/sysgen/master.d*, for example, */var/sysgen/master.d/hypothetical*.
- The file in */var/sysgen/master.d* specifies the prefix for driver entry points, for example *hypo\_*.
- The same name with the suffix *.o*, is searched for in */var/sysgen/boot*—for example, */var/sysgen/boot/hypothetical.o*. This object file is linked with the kernel.
- The public symbols in the object file are searched for names that start with the prefix, for example **hypo\_edtinit()**. These are noted in the kernel switch table so the driver can be called as needed.

## Placing the Object File in `/var/sysgen/boot`

The `/var/sysgen/boot` directory, where the kernel object modules reside, is actually a symbolic link to one of the directories `/usr/cpu/sysgen/IPnnboot`, where *nn* is the number of one of the CPU modules supported by the current release of IRIX (see “CPU Modules” on page 5). When you place the object file of a driver in `/var/sysgen/boot`, you actually place it in the CPU directory for the system in use.

## Describing the Driver in `/var/sysgen/master.d`

You describe your driver in a file with the name of the driver in `/var/sysgen/master.d`. The format of these files is described in two places: the `master(4)` reference page, and in `/var/sysgen/master.d/README`. In addition, you can examine the many examples in the distributed system.

### Descriptive Line

The first noncomment line of the master file contains a series of fields, delimited by white space, to describe the driver. These fields are:

Flags	Described in the text.
Prefix	The string of 1-14 characters that identify the public symbols of driver entry points.
Major number	The major number, or a comma-separated list of major numbers, found in device special files that are managed by this driver.
Number of sub-devices	Size of the driver’s static arrays of device information, or given as a hyphen “-” when not known.
Dependencies	A list of other modules that must be in the kernel for this driver to work—usually omitted except for SCSI drivers.

Of the many flags that can be specified in the first field, the following are most frequently used:

- b* or *c*      Block (b) or character (c) device. One or the other is essential for any device driver.
- f* or *m*      STREAMS driver (f) or module (m).
- n* or *p*      Multiprocessor-aware (n) or uniprocessor-only (p) driver. This flag must correspond to the presence or absence of `D_MP` in the `pxdevflag` global.
- s*            Software driver, either a pseudo-device or a SCSI driver.

The *s* (software-only) flag tells `lboot` not to attempt to probe for hardware. This is the case with software-only (pseudo-device) drivers, and with SCSI drivers. If `lboot` tries to probe for a SCSI device, it fails, and assumes that the device is not present, and does not include your SCSI device driver.

Additional flags (d, R, N, D) for loadable drivers are discussed later under “Configuring a Loadable Driver” on page 239.

### Listing Dependencies

The descriptive line ends with a comma-separated list of other loadable kernel modules on which this driver depends. The `lboot` command makes sure that, if it fails to load a dependent module, it also will not load this module.

In most cases, an OEM driver does not have any dependencies. However, a SCSI driver (see Chapter 13, “SCSI Device Drivers”) should list the name `scsi`, since it depends on the inner SCSI driver. A STREAMS driver might list the name of a STREAMS support module such as `clone` (see “Support for CLONE Drivers” on page 510).

It is possible for you to design a driver in the form of multiple, loadable modules. In that case, you would name your support modules in this field.

### Stubs Section

Noncomment lines that follow the descriptive line and precede a line beginning “\$” are used by library modules—not by device drivers or STREAMS drivers. Each such line specifies an entry point that this module provides, and which is used by the kernel or some other loadable module. These lines instruct *lboot* in how to create a harmless “stub” function in the event that this driver is not included in the kernel—for example, because it is specified by an EXCLUDE line in the system file. The format is discussed in the master(4) reference page.

Since a device or STREAMS driver provides only standard entry points that are accessed via the switch tables rather than by linking, drivers do not need to define any stubs.

### Variables Section

Following the descriptive line (and the stubs section, if any), you can place a line that begins with “\$” and, following this, you can write definitions of global variables using C syntax. This text is compiled into a module linked with the kernel. You refer to these variables as *extern* in the driver source code.

The advantage of defining global variables in the master file is that the initializing expressions for these variables can include values taken from the descriptive line. The following special symbols can be used:

- ##E      The integer coded as the major number in the descriptive line. The first integer, if a list of major numbers is given.
- ##C      The number of controllers (bus adapters) of this type.
- ##D      The number of sub-devices as coded in the fourth field of the descriptive line.

You can use these symbols to compile run-time values for the major device number and the number of supported sub-devices, as specified in the descriptive line of the file, without coding them as constants in the driver. In the source code you can write

```
extern major_t myMajNum;  
extern int myDevLimit;
```

In the master file you can implement the variables using the code in Example 10-1

**Example 10-1** Defining Variables in Master Descriptive File

```
$$$  
major_t myMajNum = ##E;  
int myDevLimit = ##C;
```

(In a loadable driver this technique requires one additional step; see “Master File for Loadable Drivers” on page 240.)

### Configuring a Kernel

Once you have placed the binary in */var/sysgen/boot* and the master file in */var/sysgen/master.d*, you can configure and generate a new kernel. This is done using the *autoconfig* command, which in turn calls *lboot* to actually create a new bootable file.

The *lboot* program only loads modules that are specified in a file in */var/sysgen/system*. One command is required to specify the new driver; the command is one of the following:

- VECTOR      To specify hardware details, to request a hardware probe at boot time, to load the driver and invoke *pfxedtinit()*.
- INCLUDE     To load the driver and invoke *pfxinit()*.
- USE         To load the driver and invoke *pfxinit()* only if the master file exists in *master.d*.

The form of these commands is detailed in the *system(4)* reference page. In addition, you should examine the distributed files in */var/sysgen/system*, especial *irix.sm*, which contains many comments on the meaning and use of different entries.

You could place the VECTOR, USE, or INCLUDE line for your driver in *irix.sm*. However, since *lboot* treats all files in */var/sysgen/system* as a single file, you can create a small file unique to your driver. The name of this file is not significant, but a good name is the driver name with the suffix *.sm*.

## Generating a Kernel

The *autoconfig* script invokes *lboot* to read the system files, read the master files, and link all the kernel executables. Provided there are no errors, the output is a new file */unix.install*. At boot time this file is moved to the name */unix* and used as the kernel.

During the testing period you may want to keep the original kernel file available as */unix.original*. A simple way to retain this file is to create a link to it using the *ln* command.

## Configuring a Loadable Driver

You compile and configure a loadable driver very much as you would a nonloadable driver. The differences are as follows:

- You provide an additional global variable with the public name *pfxmversion*.
- You use a few different compile options.
- You decide when the driver should be loaded, and use appropriate flags in the descriptive line in the master file.

For more background on loadable modules, see the *mload(4)* and *ml(1)* reference pages.

### Public Global Variables

To be loadable, a driver must specify a *pfxdevflag* entry point, and it may not contain the *D\_OLD* flag (see “Driver Flag Constant” on page 145).

Any loadable module must define a public name *pfxmversion*, declared as follows:

```
#include <sys/mload.h>
char *pfxmversion = M_VERSION;
```

Note the exact spelling of the variable; it is easy to overlook the letter “m” after the prefix. If the variable does not exist or is incorrectly spelled, an attempt to load the driver will fail.

## Compile Options for Loadable Drivers

Use the `-G 0` option when compiling and linking a loadable driver, since the global option table is not available to a loadable driver (see “Compile Options, 32-Bit Kernel” on page 232 and “Compile Options, 64-Bit Kernel” on page 233).

In a loadable driver, link using the `-r` and `-d` options to retain the symbol table yet generate a bss segment.

## Master File for Loadable Drivers

The file in `/var/sysgen/master.d` for a loadable driver has different flags.

### Descriptive Line

In the flags field of the descriptive line of the master file (see “Descriptive Line” on page 235), you specify that the driver is loadable, and when it should be loaded. The possible flags are as follows:

- `d` Specifies that this is a loadable driver.
- `R` Auto-register the module (discussed in text).
- `D` Load, then unload, at boot time, in order to let the driver initialize the hardware inventory.
- `N` Prevent this module from being automatically unloaded even when it has a `pfxunload()` entry point.

When the `d` flag is given for an included module, `lboot` parses the master file for the driver. Global variables defined in the variables section of the master file (see “Variables Section” on page 237) are defined and included in the kernel. However, object code of the driver is not included in the kernel, and the names of its entry points are not entered into the kernel switch table.

Such a driver has to be manually loaded with the `ml` or `lboot` command before it can be used; and it cannot be used from the miniroot.

## Registration

A loadable module is “registered” by having a stub entry placed in the *pfxopen()* column of a kernel switch table, indicating it exists but is not loaded. Registration can be done manually, after bootstrap, or automatically during bootstrap.

Registration is done automatically by for each master descriptive file that contains the *d* and *R* flags. Autoregistration is done at bootstrap phase 2. It is initiated by the script */etc/rc2/S23autoconfig*. Registration can be done manually at any time after bootstrap by using the *ml* or *lboot* command with the *reg* option (see the *ml(1M)* and *lboot(1M)* reference pages).

Once it has been registered, a driver is loaded automatically the first time a process attempts to open a device special file with the major device number of this driver. You can also load a registered driver in advance of any use with the *ml* or *lboot* command—loading implies registration.

## Loading

A registered, loadable driver can be loaded by the *lboot* and *ml* commands. When a driver is loaded, the following steps are taken:

1. The object file header is read.
2. Memory is allocated for the module’s text, data, and bss segments.
3. The module’s text and data are read.
4. The module’s text and data are relocated. References to kernel names and to global variables named in the master file are resolved.
5. The module entry points are noted in the appropriate kernel switch table.
6. The module’s *pfxinit()* or *pfxedtinit()* entry point is called, depending on whether the module is specified by an INCLUDE or a VECTOR statement in the system file (see “Initialization Entry Points” on page 147).
7. The module’s *pfxstart()* entry point, if any, is called.

Space allocated for the module's text, data, and bss is located in either *k0seg* or *k2seg*. Static buffers in loadable modules are not necessarily physically contiguous in memory.

When the *D* flag is included in the descriptive line in the descriptive file, the module is loaded at boot time and then unloaded. This gives the module a chance to update the hardware inventory.

If errors occur when a module is loaded, see the `mload(4)` reference page for a list of possible causes.

## Unloading

A module can be unloaded only when it provides an `pfxunload()` entry point (see "Entry Point `unload()`" on page 170). The *N* flag can be specified in the master file to prevent automatic unloading in any case.

A loaded module is automatically unloaded following a delay after the last close of a device it supports. The delay is configurable using *sysctl*, as the `module_unld_delay` variable (see the `sysctl(1)` reference page). You can use *ml* to specify an unloading delay for a particular module.

The *lboot* or *ml* command can be used to unload a module before the delay expires, or when the *N* flag prevents automatic unloading.

Just before unloading, the kernel invokes the driver's `pfxunload()` entry point. Then the module is removed from memory, and returned to registered status. It is up to the `pfxunload()` entry point to decide whether unloading can be permitted. Experience has shown that most of the problems with loadable drivers arise from unloading and reloading. The precautions to take are described under "Entry Point `unload()`" on page 170.

## Configuring for a Dynamic Major Number

You can place a hyphen (“-”) in the third field of the descriptive line; that is, the field for the major device number (see “Descriptive Line” on page 235). When the line also contains the *b* or *c* flag, indicating a device driver, a hyphen in the third field means that *lboot* is to assign some unused major device number dynamically when the module is registered.

Since a device driver can only be called by opening a device special file, and since a device special file has to be defined based on a major device number, how can the device special files be created?

The assigned major number can be discovered using the *ml* command. The command

```
ml list -r
```

writes a list of all registered modules, including their major numbers. The following procedure can be used to make the assignment of the major device number completely dynamic.

1. Make the device driver loadable, specifying the *d*, *R*, and *D* flags.
2. Specify a hyphen for the major number.
3. In the driver, get the major number dynamically, if indeed the driver needs to know its major number at all.
4. In a script executed from *rc2.d* (just as the */dev/MAKEDEV* script is executed), call *ml* to list registered drivers.
5. Parse the output of *ml* using UNIX utilities to extract the major device number from the line describing your driver.
6. Execute *mknod* or *install* commands to create special device files using the discovered major number.



## Testing and Debugging a Driver

As a critical system component, a driver deserves careful testing, but because it is part of the kernel, the normal testing tools are not available. This chapter describes the available testing tools and methods, in the following major topics:

- “Preparing the System for Debugging” on page 245 describes how to set up the kernel for use of the debugging tools.
- “Producing Diagnostic Displays” on page 251 covers the kernel functions your driver can use to generate diagnostic output as it executes.
- “Using symmon” on page 254 describes the use of the standalone debugger.
- “Using idbg” on page 264 describes some uses of the kernel-display command.

### Preparing the System for Debugging

The standalone debugger *symmon* is a key tool for driver programming. It must be installed in the volume header of the boot disk. In order for it to be useful you must boot a “debugging” kernel, that is, one that retains symbols, and contains the display modules, that are used by debugging tools. Normally these modules and symbols are eliminated to save space. You modify the *irix.sm* file to enable debugging, and then generate a new kernel.

All these steps should be performed before you attempt to install your device driver.

## Placing *symmon* in the Volume Header

The *symmon* standalone debugger resides in the volume header of a disk—not in a normal IRIX filesystem. The volume header is disk partition 0. It always contains a label record (*sgilabel*), and the standalone shell *sash* that manages the bootstrap operation. On some systems it may also contain the *ide* program, a PROM-level diagnostic program. If *symmon* is to be available, it, too, must be placed in the volume header.

Normally you acquire *symmon* by installing the debugging kernel feature (*oe.sw.kdebug*) in the IRIX Developer Option software distribution. You can verify that this feature has been installed by executing the command

```
versions ooe.sw.kdebug
```

The response should confirm the presence of this component (it does not show *symmon* by name). When you install the kernel debug feature, the *symmon* program file is copied to the volume header of the current boot disk automatically.

You can verify the presence of *symmon* in the volume header through the use of *dvhtool* (described in the *dvhtool(1)* reference page). The results should be similar to the display in Example 11-1. The response to the “l” (list) command shows that the volume header of this disk contains *sgilabel*, *ide*, *sash*, and *symmon*.

### Example 11-1 Verifying Presence of *symmon*

```
# dvhtool
Volume? (/dev/rvh) /dev/rvh
Command? (read, vd, pt, dp, write, bootfile, or quit): vd
(d FILE, a UNIX_FILE FILE, c UNIX_FILE FILE, g FILE UNIX_FILE or l)?
l
Current contents:
  File name      Length      Block #
  sgilabel       512         2
  ide            281600      278
  sash           281600      828
  symmon        248320      1378
(d FILE, a UNIX_FILE FILE, c UNIX_FILE FILE, g FILE UNIX_FILE or l)?

Command? (read, vd, pt, dp, write, bootfile, or quit): quit
#
```

In the event you need to install *symmon* in the volume header of a disk without using the software manager, you can copy the standalone program to the volume header using *dvhtool*. However, you first need to get a copy of the program in the form of a UNIX file.

Starting from a volume that currently has a copy of *symmon* (verified as in Example 11-1), use *dvhtool* to extract a copy of *symmon* into a convenient spot.

```
dvhtool -v g symmon /var/tmp/symmon.IPxx
```

There is a unique version of *symmon* for each CPU module, so it is a good idea to qualify the filename with the CPU module type. Once the program is available as a normal file, you can use *dvhtool* to install it in the volume header of some other disk.

In the event there is not enough room in partition 0 (the volume header) of the target disk, it is safe to use *dvhtool* to delete the *ide* program from the volume header. The *ide* application can be booted manually from a CDROM if it is ever required.

## Enabling Debugging in irix.sm

In order to make debugging symbols available in the kernel, you must make two changes, one required and one optional, in the file */var/sysgen/system/irix.sm*. As superuser, make a hard link to the file */var/sysgen/system/irix.sm* as *irix.sm.nondebug*. This enables you to return easily to a nondebugging kernel.

### Including Symbols in the Kernel Image

Edit */var/sysgen/system/irix.sm*. Near the end, note the lines that resemble the following:

```
* Compilation and load flags
*   To load a kernel that can be co-resident with symmon
*   (for breakpoint debugging) replace LDOPTS
*   with the following. You must also INCLUDE prf and idbg.
*
*LDOPTS: -non_shared -N -e start -G 8 -elf -woff 84 -woff 47 -woff 17
-mips2 -o32 -nostdlib -T 88069000
```

The active LDOPTS statement (the one without an initial asterisk) appears a few lines later. Remove the asterisk from the front of the debugging LDOPTS to make it active. Insert an asterisk to convert the original LDOPTS into a comment.

### Including *idbg* in the Kernel Image

The symbol-display routines used by the command-line kernel display tool, *idbg*, are contained in optional kernel modules. (See “Using *idbg*” on page 264.) You can change */var/sysgen/system/irix.sm* so that support for *idbg* is always present in the kernel. Alternatively, you can load these modules manually with *ml* before you use them (see the *ml(1)* reference page).

If you are entering an extended debugging period, make the modules permanent. Look for the lines in */var/sysgen/system/irix.sm* that resemble the following:

```
*
* Kernel debugging tools (see profiler(1M) and idbg(1M))
*
EXCLUDE: idbg
EXCLUDE: dmiidbg, grioidbg, xfsidbg, xlvidbg, cachefsidbg
USE : prf
```

Change these lines, if necessary, so that both *idbg* and *prf* are marked USE, not EXCLUDE. Verify that the file */var/sysgen/boot/idbg.o* exists. It is normally installed with the debugging kernel feature.

Parts of the *idbg* support that are unique to particular filesystems are in the other modules listed in this area of *irix.sm*. Modules such as *xlvidbg* are useful to Silicon Graphics developers but are not likely to be helpful to developers of third-party drivers. However, it does no harm to change those modules from EXCLUDE to USE also.

### Including Lock Metering in the Kernel Image

In addition to the display support included by the *idbg* modules, you can include a module that supports lock metering. This module keeps statistics on the use of each semaphore, basic lock, and reader/writer lock, and displays the statistics through *idbg* commands. To enable this facility, locate the lines in */var/sysgen/system/irix.sm* that resemble the following:

```
* Required kernel modules
...
* ksync - kernel synchronization routines (mutex_lock, sv_wait,
psema...)
*   or
*   ksync_metered - metered kernel synchronization routines
...
*
```

```
KERNEL: kernel
INCLUDE: os, disp, mem, zero
INCLUDE: ksync
EXCLUDE: ksync_metered
```

Reverse the state of the two “ksync” lines so that ksync is excluded and ksync\_metered is included.

## Generating a Debugging Kernel

Run the *autoconfig* command to generate a new kernel that will reflect the changes made in *irix.sm*. The result is a new kernel file, */unix.install*, that will be renamed to */unix* and used when the system is booted. This kernel can support *idbg* but is not yet ready for standalone debugging with *symmon*.

The *setsym* command copies the symbol table from a kernel file and stores it as data within the kernel, so that *symmon* can find it. After *autoconfig* has created */unix.install*, apply the *setsym* command to it, as follows:

```
#setsym /unix.install
```

If this command returns an error message about “symbol table overflow,” it means you have neglected to activate the debugging LDOPTS statement in */var/sysgen/irix.sm*.

**Tip:** You can use *setsym* with the *-d* option to generate a list of all symbols in the kernel being modified. The list is very long; direct it to a file for later reference.

At this time, you may wish to create a link to the current, nondebugging kernel so you can retrieve it easily. You can also return to a nondebugging kernel by restoring the original *irix.sm* file and running *autoconfig* again.

## Specifying a Separate System Console

In order to use the standalone debugger, you must have an ASCII terminal as a separate system console device. Install a terminal next to the system or workstation and connect it to the first serial port. Verify its operation as a terminal by logging in to the system. You may have to modify the file */etc/inittab* so that the line for the alternate console is active (see the *inittab(4)* reference page).

When you have verified that the terminal works, use the *nvr* command to change the nonvolatile RAM variable console from a letter “g” to a letter “d,” as follows:

```
# nvr console
g
# nvr console d
# nvr console
d
```

The *nvr* command is used to report and change the contents of the nonvolatile RAM variables used by the boot PROM and standalone shell (see the *nvr*(1) reference page).

### Verifying the Debugging Tools

After performing the preceding steps, restart the system. Messages from *sash* appear on the attached terminal, rather than on the graphics screen. If *symmon* is present, it announces itself on the console terminal also.

To verify operation of *idbg*, issue the *idbg* command and display the process list:

```
# idbg
idbg> plist
active process list:
34:672:"xdm" pri(60) SLEEP flags: load uload siglck recalc sv
0:0:"sched" ndpri(39) SLEEP flags: sys nwake load uload sv
31:193:"inetd" pri(60) SLEEP flags: load uload siglck recalc sv
...
```

To verify operation of *symmon*, press control-A at the console terminal. The prompt string *DBG:* should appear. At this time the system is frozen and no longer responds to mouse or keyboard input. Type the letter c (for continue) and press return. The system returns to life.

## Producing Diagnostic Displays

Normally a device or STREAMS driver produces display output in only two cases:

- To advise the operator or administrator of a serious problem.
- To display debugging information during software development.

Both of these purposes are served by the **cmn\_err()** function. It brings to a kernel-level module the abilities that a user-level process gets from **printf()** and **syslog()**.

### Using **cmn\_err**

The details of **cmn\_err()** usage are in the **cmn\_err(D3)** reference page. The function prototype and the constant values it uses are declared in *sys/cmnerr.h*.

In summary, **cmn\_err()** takes two or more arguments:

- A severity code that specifies how the message should be treated when it is written to the system log.
- A message string, which can have substitution points in the style of **printf()**.
- As many numeric values as are needed to substitute into the message string.

The first character of the message string specifies the destination of the message, either an in-memory buffer or the system log, or both.

### Displaying to the System Log

The message is sent to the the system log daemon whenever the first message character (after substitution) is not an exclamation mark ("!"). The message is written only to the system log when the first message character is a circumflex ("^").

This is basically the same service that a user-level process receives from the **syslog()** function. (Compare the **syslog(3)** and **cmn\_err(D3)** reference pages, and examine the *sys/cmnerr.h* header file; the relationship is clear.) The first argument to **cmn\_err()** is a severity code which corresponds to one of the severity codes supported by **syslog()**: **CE\_WARN** equals **LOG\_WARN**, and so on.

Use **cmn\_err()** to write log messages to record serious errors (with **CE\_ALERT** severity) or to advise the administrator of conditions that should be changed (using **CE\_NOTE**).

### Displaying to the Circular Message Buffer

The message is stored in the next available position in a circular buffer in kernel memory whenever the first message character (after substitution) is not a circumflex (“^”). The message is stored only in the memory buffer when the first message character is an exclamation mark (“!”).

The name of the circular buffer (as a symbol to *idbg* or *symmon*) is *putbuf*. The contents of *putbuf* can be displayed with the *pb* command of either *idbg* or *symmon* (see “Using *symmon*” on page 254 and “Using *idbg*” on page 264), or in a post-mortem dump using *icrash* (see “Using *icrash*” on page 271). Use **cmn\_err()** to store debugging trace data in the circular buffer, and extract it after a stop or breakpoint with *symmon*, or use *idbg* to look at it while the system is running.

The size of the buffer can be configured by changing the kernel variable *putbufsz*, as shown in the dialog in Example 11-2.

#### Example 11-2 Setting Kernel *putbuf* Size

```
# systune -i
Updates will be made to running system and /unix.install
systune-> putbufsz
      putbufsz = 1024 (0x400)
systune-> putbufsz 4096
      putbufsz = 1024 (0x400)
      Do you really want to change putbufsz to 4096 (0x1000)? (y/n)
Y
In order for the change in parameter putbufsz to become effective,
reboot the system
systune-> quit
```

### Using **cmn\_err()** Through Macros

The inventive C programmer can think of many ways to invoke **cmn\_err()** using macros. One method is illustrated in the example driver displayed in Chapter 12, “Driver Example.” It contains the code shown in Example 11-3.

**Example 11-3** Debugging Macros Using `cmn_err()`

```
#ifdef DEBUG
#define DBGMSG0(s) cmn_err(CE_DEBUG,s)
#define DBGMSG1(s,x) cmn_err(CE_DEBUG,s,x)
#define DBGMSG2(s,x,y) cmn_err(CE_DEBUG,s,x,y)
#define DBGMSG3(s,x,y,z) cmn_err(CE_DEBUG,s,x,y,z)
#else
#define DBGMSG0(s)
#define DBGMSG1(s,x)
#define DBGMSG2(s,x,y)
#define DBGMSG3(s,x,y,z)
#endif
```

Calls to the `DBGMSG` macros evaluate to nothing unless the `DEBUG` variable is defined to the compiler. Macro use could be more elaborate. For example, suppose that you want every debugging macro to start with the same string, and that you do not want to have the tedium of coding the newline at the end of every debug message text. You could write macros like the one shown in Example 11-4.

**Example 11-4** More Elaborate Debugging Macro

```
#ifdef DEBUG
char *_dbg_prefix = "mydriver: %s\n";
#define DBGPREFIX(s) cmn_err(CE_DEBUG,_dbg_prefix,s)
...
#else
#define DBGPREFIX(s)
...
#endif
```

**Using `printf()`**

You can call the `printf()` function from a kernel module. The kernel version of `printf()` is basically a call to `cmn_err()` with severity `CE_CONT`. In general it is better to use `cmn_err()` explicitly.

## Using ASSERT

The `assert()` macro is familiar to many C programmers; it terminates a program with a message if its argument evaluates to false (see the `assert(3X)` reference page). This normal `assert()` macro does not work in a kernel module because the normal C library is not available. However, a similar function is available as the `ASSERT()` macro in the header file `sys/debug.h`.

The `ASSERT()` macro compiles to null code unless the compiler variable `DEBUG` is not only defined, but defined as `YES`. When it compiles to executable code, `ASSERT()` tests its argument. If the argument evaluates to false, a kernel panic is forced.

Clearly `ASSERT()` must be used with care, testing conditions that are truly essential to the integrity of the system. When reporting conditions that are merely operational errors, use a call to `cmn_err()` with the `CE_WARN` option.

## Using symmon

The `symmon` program is a standalone debug monitor that can display and modify memory, and stop, start, and trace execution, without using any kernel facilities. Using `symmon` you can set breakpoints in your driver, single-step its execution, and display the contents of driver and kernel variables.

The facilities of `symmon` are unsophisticated compared to the high-level debuggers you might use to debug a user-level application. For example, `symmon` does not understand C syntax, so it cannot display data structures as structures. Execution tracing is done at the level of machine instructions, not at the level of C statements.

However, you can use `symmon` to examine the operations of a kernel module in a running system, and resume execution of the system. This is an invaluable facility when debugging a new driver.

## How symmon Is Entered

When the system boots a debugging kernel with *symmon* installed, control can pass into the debug monitor under several different circumstances:

- Early in the bootstrap process, if certain environment variables are set in the stand-alone shell (see “Entering symmon at Boot Time” on page 256).
- Whenever a control-A character is typed at the system console terminal.
- Whenever a breakpoint is reached or a watchpoint is tripped (see “Commands to Control Execution Flow” on page 259).
- Whenever a kernel module calls the kernel function `debug(uchar_t *msg)`.
- When a non-maskable interrupt (NMI) is detected.
- When a kernel panic is detected or forced with `cmn_err()`.

When *symmon* gains control, it displays its “DBG:” prompt at the console terminal and waits for a command.

To resume execution at the point of interruption, enter the *c* (continue) command.

## Using symmon in a Uniprocessor Workstation

In a single-processor workstation such as the Indy or Indigo<sup>2</sup>, no IRIX execution takes place while *symmon* is running. The mouse and keyboard are unresponsive. (One keystroke may be stored in the keyboard hardware to be processed when the system resumes execution.) As a result, time-dependent processes can fail; for example, the system clock is not updated. Network interrupts are not taken, so if the workstation is acting as an NFS server, it will appear to be dead to other systems.

## Using symmon in a Multiprocessor Workstation

In a multiprocessor, the CPU that was interrupted runs *symmon* and nothing else. For example, the CPU that executes the breakpoint, or the CPU that handles the DUART interrupt that returns the control-A character, or the CPU in which `debug()` was called, comes under the control of *symmon*. Other CPUs continue to execute normally. However, if the *symmon* CPU holds a lock, other CPUs may come to a halt waiting for the lock to be released.

The *symmon* breakpoint table is shared by all CPUs. A breakpoint set from one CPU can be taken by another CPU, or by multiple other CPUs. It is possible to run multiple instances of *symmon* concurrently. The output from all instances of *symmon* is multiplexed onto the system console terminal. However, only one CPU at a time issues the DBG: prompt. Use the *cpu* command with no argument to find out which CPU is prompting. Use the *cpu* command with a *cpu* number to switch to a different CPU. (See “Commands to Control Execution Flow” on page 259.)

**Tip:** To make multiprocessor execution more predictable, reduce the number of available CPUs. Use the PROM monitor to disable one or more CPUs before bootstrap.

### Entering *symmon* at Boot Time

You can cause the kernel to stop during initialization and enter *symmon* during the bootstrap process. In order to do this, you must use the miniroot to set environment variables.

1. Restart the system, for example by giving the commands *sync* and *halt*. Eventually, the 5-item PROM menu is displayed at the console terminal.
2. Select item 5, “Enter the Command Monitor.”
3. Set one or both of the environment variables *dbgstop* and *symstop* to 1, using commands such as the following:  

```
>> setenv symstop 1
```
4. Return to the PROM menu by entering the command *exit*.
5. Select menu item 1, “Start System.”

In either case, *symmon* seizes the system and displays its DBG: prompt at the system console during bootstrap. When the *dbgstop* variable is set, *symmon* takes control of the system very early in the bootstrap process. Symbolic names are not initialized at this point. However, breakpoints can be set and memory can be displayed using explicit addresses.

When the *symstop* variable is set, *symmon* takes control after symbols are defined, but before driver initialization is begun. At this stop, you can display memory and set breakpoints based on entry point names of your driver.

## Commands of symmon

The exact set of commands supported by *symmon* changes from release to release and from CPU model to CPU model. Many *symmon* commands are useful only to Silicon Graphics engineers who are debugging hardware and kernel problems. For a complete list of commands, see the *symmon(1M)* reference page, or enter *symmon* and give the *help* command. You can use control-S and control-Q on the console terminal to pause the scrolling display.

The commands described in this section are generally useful and are available on all CPU models under IRIX 6.2. These commands can be grouped into the following categories:

- Conversion between symbols and memory addresses.
- Execution control, including commands for stopping, starting, and setting breakpoints.
- Display and modification of memory, including the display of machine registers and of system data structures such as the *buf\_t* and *proc\_t* objects.
- Management of the virtual memory system and the TLB.

## Syntax of Command Elements

The *symmon* commands all have the same form: a keyword, usually followed by one or more arguments separated by spaces.

Many commands take an address value. An address argument value can have one of the following forms:

Decimal number	A number starting with 1-9 is a decimal number, for example <i>4095</i> .
Octal number	A number starting with 0 and a digit is an octal number, for example <i>033</i> .
Hex number	A number starting with 0x is a hexadecimal number, for example <i>0xffff8000</i> .
Binary number	A number starting with 0b is a binary number, for example <i>0b0100</i> .

Symbol	A word starting with a non-digit is looked up as a symbol in the kernel symbol table, and its address is the value; for example <i>dk_open</i> .
Register	A word starting with "\$" is taken as a register name, and the contents of the register as of the last interrupt is used as the argument value; for example <i>\$a2</i> .
Value and offset	A value plus or minus a number is a value, for example <i>\$a2-0x100</i> or <i>dk_open+128</i> .

Some commands accept a range of addresses. A range can be written in one of two ways:

- As *value1:value2*, meaning an inclusive range of addresses from *value1* through *value2*, for example *prtbuf:prtbuf+4096*.
- As *value1#count2*, meaning a range of count2 bytes beginning at value1, for example *prtbuf#4096*.

The register names that *symmon* accepts and shows in various displays are the conventional names used in MIPS assembly language programming. Refer to the *MIPSpro Assembly Language Programmer's Guide* and the processor manuals listed under "Additional Reading" on page xxix.

### Commands for Symbol Conversion and Lookup

The commands summarized in Table 11-1 are used to convert between symbolic names and their corresponding addresses.

**Table 11-1** Commands for Symbol Conversion and Lookup

Command	Example	Operation
<b>hx</b> <i>name</i>	<b>hx dk_read</b> dk_read 0xffffffff882b0510	The name is looked up on the symbol table and if it is found, its address is displayed.
<b>lkaddr</b> <i>addr</i>	<b>lkaddr 0x882b0510</b> 0x882af910 lockdisptab 0x882b0510 dk_read 0x882b051c dk_write	Symbols near to the specified <i>addr</i> are listed. Use this command to find out the symbolic location of an unexpected stop.

**Table 11-1 (continued)** Commands for Symbol Conversion and Lookup

Command	Example	Operation
<b>lkup</b> <i>letters</i>	<b>hx dk_rea</b> 0x880d5f10 dk_readcap 0x882b0510 dk_read 0x332b0528 dk_readcapacity	Every symbol that contains the specified <i>letters</i> at any point is listed. There is no way to anchor the search to the beginning or end of the name.
<b>msyms</b> <i>ident</i>	<b>msyms 13</b> Symbols for module 13 (prefix tcl) tclinit 0xc0403d9c tclmversion 0xc0405fe0	The symbols for the loadable module <i>ident</i> are listed. Use the <i>ml</i> command with no arguments to list all modules and their ident numbers.
<b>nm</b> <i>addr</i>	<b>nm 0xc0403da0</b> 0xc0403da0 tclinit+0x4	The symbol nearest to the specified <i>addr</i> is listed.

**Note:** When symmon displays an address it normally shows a full 64 bits. In a 32-bit kernel, the most-significant 32 bits of a kernel virtual address are all-binary-1, from extension of the sign bit of the 32-bit address—as shown in the example of *hx* in Table 11-1. When you enter an address to a command in a 32-bit system, you need to type only the significant 32-bit value.

## Commands to Control Execution Flow

The commands summarized in Table 11-2 are used to stop, start, and single-step execution in the kernel.

**Table 11-2** Commands to Control Execution

Command	Example	Operation
<b>brk</b>	<b>brk</b>	List all breakpoints currently set.
<b>brk</b> <i>addr</i>	<b>brk dk_read</b>	Set a breakpoint at the specified <i>addr</i> .
<b>c</b>	<b>c</b>	Restart execution at the point of interruption in the current CPU.
<b>c</b> <i>cpuid</i> [ <i>cpuid</i> ]... <b>c</b> <b>all</b>	<b>c 0</b>	Restart execution in the specified CPU, or in all stopped CPUs. Available in multiprocessors only.

**Table 11-2 (continued)** Commands to Control Execution

Command	Example	Operation
<code>call addr [args]</code>	<code>call getemisor 0</code>	Call a kernel function and report the contents of the result register on return.
<code>cpu</code>	<code>cpu</code>	Displays the cpu ID of the currently-executing CPU. Available in multiprocessors only.
<code>cpu cpuid</code>	<code>cpu 0</code>	Force <i>symmon</i> execution to the specified CPU. That CPU must be executing <i>symmon</i> . Other CPUs executing <i>symmon</i> wait. Available in multiprocessors only.
<code>goto addr</code>	<code>goto getemisor</code>	Set a temporary breakpoint at <i>addr</i> and then continue execution as for the <i>c</i> command (in effect “go until <i>addr</i> is reached”).
<code>quit</code>	<code>quit</code>	Return to the boot PROM, forcing an instant reboot.
<code>s [count]</code>	<code>s 8</code>	Single-step through 1 or <i>count</i> instructions, displaying each instruction and the register contents it uses. A branch and the instruction in “delay slot” following it count as 1. Steps into subroutines.
<code>S [count]</code>	<code>S 8</code>	Single-step through 1 or <i>count</i> instructions as for the <i>s</i> command, but do not step into subroutines.
<code>unbrk n</code>	<code>unbrk 2</code>	Remove break point number <i>n</i> . Use <i>brk</i> with no argument to list break points by number.
<code>wpt {r w rw} physaddr</code>	<code>wpt r 0x0841f608</code>	Set a hardware watchpoint on a physical address.

**Tip:** One way to force a memory dump from *symmon* is the command `call dumpsys`.

Following a break or a watchpoint, use the *bt* command to display the stack history and use *printreg* to display the registers (see “Commands to Display Memory” on page 262).

The hardware watchpoint used by the *wpt* command uses hardware registers in the MIPS R4000 and R10000 processors (the R8000 does not support the watchpoint registers). When a read or write access is addressed to any byte in the doubleword specified by the physical address, *symmon* gains control and displays the instruction that is attempting the access on the console terminal.

The argument of *wpt* must be a physical memory address and a multiple of 8. Use *tlbvtov* to get the physical equivalent of an address in a user address space (see “Commands to Manage Virtual Memory” on page 261). In a 32-bit kernel, the physical equivalent of an address in kernel space is obtained by changing the most significant hex digit to 0.

## Commands to Manage Virtual Memory

The commands summarized in Table 11-3 are used to display and manage the virtual memory translation system.

**Table 11-3** Commands to Manage Virtual Memory

Command	Example	Operation
cacheflush <i>range</i>	<b>cacheflush \$6:\$6+4096</b>	Flush both the instruction and data caches when they contain data that falls in <i>range</i> .
tlbdump [ <i>lo:hi</i> ]	<b>tlbdump 1:3</b>	Display the contents of the TLB registers. When a range of numbers is given, the registers from <i>lo</i> through <i>hi</i> -1 are displayed.
tlbflush [ <i>lo:hi</i> ]	<b>tlbflush</b>	Flush (nullify) the TLB registers specified. The registers are reloaded as required during subsequent execution.
tlbpid	<b>tlbpid</b> Current dbgmon pid = 79	Display the process slot number of the process whose context is in the TLB.
tlbvtov <i>addr</i>	<b>tlbptov 0xffffc000</b>	Display the TLB register that maps <i>addr</i> .

## Commands to Display Memory

The commands summarized in Table 11-4 are used to display memory or variables.

**Table 11-4** Commands to Display Memory

Command	Example	Operation
bt [ <i>frames</i> ]	<b>bt 4</b>	Display the calling function, the arguments, and the name of the called function for up to <i>frames</i> stack frames. Most useful after a break or interrupt.
dis <i>range</i>	<b>dis getemisor</b>	Disassemble and display the instructions over the specified range.
dump [-b -h -w] [-o -d -x -c] <i>range</i>	dump 0xc0000000	Display memory over a specified range. The options -b, -h, and -w specify how memory is grouped, as units of 1, 2, or 4 bytes. The options -o, -d, -x, and -c specify translation into octal, decimal, hex and character.
kp [ <i>routine</i> ]	<b>kp plist</b>	Invoke a kernel print routine loaded with the idbg kernel module. If no routine is given, all available names are displayed.
printregs	<b>printregs</b>	Display all the registers as they were when the debugger was entered.
string <i>range</i> [ <i>max</i> ]	<b>string \$v1 0x80</b>	Display memory as an ASCII string in quotes. Display stops at the first null byte, or, when <i>max</i> is specified, after at most <i>max</i> bytes.

The display routines available to the *kp* command are discussed under “Using idbg” on page 264. The names that *idbg* accepts as commands are all available under *symmon* through the *kp* command.

Use the *dump* command under *symmon*. Under *idbg*, use the *hd* command for the same purpose.

## Utility Commands

The commands summarized in Table 11-5 are general-purpose utilities.

**Table 11-5** Utility Commands

Command	Example	Operation
calc	<b>calc</b>	Starts a simple stack-oriented calculator (see text).
clear	<b>clear</b>	Clear the screen of the system console terminal.
help	<b>help</b>	List one-line summaries of all available commands. Use control-S and control-Q to control the scrolling of the display.
g [-b -h -w -d] [ <i>addr</i>   <i>\$regname</i> ]	<b>g \$a1</b> 0x882fadf8:4294967295 0xffffffff	Display one byte, halfword, word or doubleword (default word) of memory, or the contents of one register at the time symmon was entered, in decimal and hex.
p [-b -h -w -d] [ <i>addr</i>   <i>\$regname</i> ] <i>value</i>	<b>p -w 0xc0000000 4095</b>	Write a byte, halfword, word, or doubleword (default word) into a saved register or into memory at the specified address.

## Using idbg

The *idbg* command is a utility that provides much of the display capability of *symmon* but from a command line, without stopping the system. Many details of *idbg* use are covered in the *idbg(1M)* reference page. Keep in mind that all *idbg* commands are available under the standalone debugger through the *kp* command (see “Commands to Display Memory” on page 262).

### Loading and Invoking idbg

Superuser privilege is required to invoke *idbg*, because it maps kernel memory. The command is ineffective unless its support modules have been made part of the kernel. This can be done permanently by changing the *irix.sm* file (see “Including *idbg* in the Kernel Image” on page 248). Alternatively, you can load the needed modules dynamically using the *ml* command, as follows:

```
# ml ld -i /var/sysgen/boot/idbg.o
```

Dynamic loading is discussed at more length in the *idbg(1M)* and *ml(1M)* reference pages.

When the support modules are loaded, *idbg* can be invoked in three styles.

### Invoking idbg for Interactive Use

Invoking the command with no arguments causes it to enter interactive mode, prompting for one command after another from standard input, as shown in Example 11-5.

#### Example 11-5 Invoking *idbg* Interactively

```
# idbg
idbg> plist 187
pid 187 is in proc slot 31
idbg> quit
#
```

The command terminates when the command *quit* is entered, or when control-D (standard input end of file) is typed.

### Invoking idbg with a Log File

Invoking the command with the *-r* option and a filename causes it to write all its output to the specified file, as shown in Example 11-6.

#### Example 11-6 Invoking idbg with a Log File

```
# idbg -r /var/tmp/idbg.save
idbg> plist 187
pid 187 is in proc slot 31
idbg> proc 31
proc: slot 31 addr 0x8832db30 pid 187 ppid 1 uid 0 abi IRIX5
  SLEEP flags: load unload siglck recalc sv
...
idbg> ^D
# cat /var/tmp/idbg.save
pid 187 is in proc slot 31
proc: slot 31 addr 0x8832db30 pid 187 ppid 1 uid 0 abi IRIX5
  SLEEP flags: load unload siglck recalc sv
...
#
```

You can use this method to collect a series of displays in a single file as you test a driver.

### Invoking idbg for a Single Command

You can invoke *idbg* with a command on the command line. The output of the single command is written to standard output, where it can be captured or piped to another program. Example 11-7 shows one simple use of this feature.

#### Example 11-7 Invoking idbg for a Single Command

```
# idbg plist | fgrep -c tcsh
3
#
```

Since the displays of *idbg* are very rich, there are endless opportunities to use this mode to generate data within shell scripts, and to process it using tools such as *awk* and *perl*. Using *perl* you could write an intelligent display routine that showed the status of your driver's private data structures using your own terminology and display format.

## Commands of *idbg*

Almost all *idbg* commands are concerned with displaying kernel memory data in different ways. There are commands to display almost every type of kernel data.

The vocabulary of commands changes from release to release, and can change within releases by software patches. Also, the commands available depend on which support modules are loaded; for example lock and semaphore meters cannot be displayed unless the *ksynch\_meter* module is loaded (see “Including Lock Metering in the Kernel Image” on page 248). Only a few commands are listed in the *idbg(1M)* reference page.

The commands summarized in this book are generally useful and available on all platforms in the current release of IRIX. For a complete (but cursory) list, use the command itself.

```
# idbg help | lp
```

In general, commands take zero or one argument. Typically the argument is a number, which can be any of the following:

- A kernel symbol, optionally +offset
- A number in hexadecimal (starting with 0x)
- A number in octal (starting with 0)
- A number in decimal.

The number is interpreted in the context of the command: sometimes it represents a process ID (*pid*), sometimes a process “slot” number or a buffer number. Often commands treat positive numbers as slot numbers or table indexes, while negative numbers are treated as addresses in kernel space.

## Commands to Display Memory and Symbols

The commands summarized in Table 11-6 are used to display memory based on specific addresses or symbols, and to display the addresses for kernel symbols.

**Table 11-6** Commands to Display Memory and Symbols

Command	Operation
<code>dsym addr [length]</code>	Dump memory by words, starting at <i>addr</i> . When a word of memory data is reasonably close to the value of a kernel symbol, the symbol plus offset is displayed instead of the hex value.
<code>hd addr [length]</code>	Dump memory in bytes, with ASCII translation, starting at <i>addr</i> . When <i>length</i> is given, it is a count of words (not bytes) to be displayed.
<code>pb</code>	Display the strings in the circular putbuf (see “Displaying to the Circular Message Buffer” on page 252).
<code>string addr [max]</code>	Display memory as an ASCII string. Display stops at the first null byte, or, when <i>max</i> is specified, after at most <i>max</i> bytes.

When you display the circular buffer, there is no special indication to show which line is the newest. You have to deduce the boundary between the newest and oldest lines from the content.

## Commands to Display Process Information

The commands summarized in Table 11-7 are concerned with displaying the status of processes. Processes are recorded in an array of “slots.” The *plist* command gives the process slot number for a given process ID. The other commands take slot numbers.

**Table 11-7** Commands to Display Process Information

Command	Operation
<code>eframe [ addr   slot ]</code>	Display the contents of an exception frame. With no argument, displays the last exception taken for the current process. Else displays the exception associated with the process specified either by address (negative number) or process table slot number (positive number)
<code>pchain slot</code>	Display the slot numbers of sibling processes to the process in <i>slot</i> .

**Table 11-7 (continued)**      Commands to Display Process Information

Command	Operation
<code>plist [ 0   pid ]</code>	With no argument, displays a one-line summary of every active process slot, including slot number and process ID. When the argument is 0, displays all inactive process slots. With a nonzero PID, displays the slot containing that process.
<code>ptree [ addr   pid ]</code>	With a <i>pid</i> (number greater than zero), finds the process structure for that process. Else uses the process structure at <i>addr</i> . Displays the command name and command arguments for that process and for all processes that descend from it.
<code>proc [ addr   slot ]</code>	Display all the fields of a process structure specified either by address (negative number) or process table slot number (positive number).
<code>signal [ addr   slot ]</code>	Display information about pending signals for the process specified either by address (negative number) or process table slot number (positive number)
<code>slpproc [ -2   -4   -8 ]</code>	Displays a summary of all processes with <code>p_stat</code> of <code>SSLEEP</code> or <code>SXBRK</code> . When an argument is given, its absolute value is used as a mask: 2 ignores processes in <code>wait()</code> ; 4 ignores processes without upages; 8 ignores processes on a sleep semaphore.
<code>ubt slot</code>	Display a backtrace of the call stack of the sleeping process in the specified slot.
<code>user [ addr   slot ]</code>	Display the user area associated with the process specified either by address (negative number) or process table slot number (positive number)

## Commands to Display Locks and Semaphores

The commands summarized in Table 11-8 display the state of semaphores and locks of different kinds, including metering information when the metered-lock module is included in the kernel.

**Table 11-8** Commands to Display Locks and Semaphores

Command	Operation
<code>;lock <i>addr</i></code>	Display the state of the spinlock at <i>addr</i> . This command is available only in multiprocessor systems.
<code>mrlock <i>addr</i></code>	Display the state of the reader/writer lock at <i>addr</i> .
<code>mutex <i>addr</i></code>	Display the state of the mutual exclusion lock at <i>addr</i> .
<code>sema <i>addr</i></code>	Display the state of the semaphore at <i>addr</i> .
<code>smeter <i>addr</i></code>	Display metering information about the semaphore at <i>addr</i> . When <i>addr</i> is positive, it is taken as an index to the semaphore metering array.
<code>sv <i>addr</i></code>	Display the state of the synchronizing variable at <i>addr</i> , including waiting processes and metering information.

## Commands to Display I/O Status

The commands summarized in Table 11-9 can be used to display the status of an I/O device or driver.

**Table 11-9** Commands to Display I/O Status

Command	Operation
<code>file [<i>addr</i>]</code>	When <i>addr</i> is omitted, displays a summary of all entries of the kernel table of open files. When <i>addr</i> is the address of a file structure, displays only that entry.
<code>scsi <i>addr</i></code>	Display the contents of the <code>scsi_request</code> structure at <i>addr</i> .
<code>uio <i>addr</i></code>	Display the contents of the <code>uio_t</code> object at <i>addr</i> .

## Commands to Display buf\_t Objects

The commands summarized in Table 11-10 are used to display the state of *buf\_t* objects and the queue of *buf\_t* objects maintained by the kernel.

**Table 11-10** Commands to Display buf\_t Objects

Command	Operation
<i>buf [addr]</i>	If <i>addr</i> is omitted, print the entire buffer chain. When <i>addr</i> is supplied as the address of a <i>buf_t</i> , dump that structure.
<i>findbuf blkno</i>	Display any <i>buf_t</i> in the buffer chain with <i>b_blkno</i> containing <i>blkno</i> .
<i>qbuf emminor</i>	Find and display all <i>buf_t</i> objects that are queued to the device with external minor number <i>emminor</i> .

## Commands to Display STREAMS Structures

The commands summarized in Table 11-11 are concerned with displaying STREAMS data structures such as message buffers.

**Table 11-11** Commands to Display STREAMS Structures

Command	Operation
<i>datab addr</i>	Display the contents of the STREAMS data block at <i>addr</i> .
<i>mbuf addr</i>	Display the contents of the STREAMS <i>mbuf</i> structure at <i>addr</i> .
<i>modinfo addr</i>	Display the contents of the module info structure at <i>addr</i> .
<i>msgb addr</i>	Display the contents of the STREAMS message block at <i>addr</i> .
<i>qband addr</i>	Display the contents of the <i>qband_t</i> object at <i>addr</i> .
<i>qinfo addr</i>	Display the contents of the <i>qinit</i> structure at <i>addr</i> .
<i>strh addr</i>	Display the contents of the <i>stdata</i> structure at <i>addr</i> .
<i>strfq addr</i>	Display the contents of the <i>queue_t</i> object at <i>addr</i> .

## Commands to Display Network-Related Structures

The commands summarized in Table 11-12 display data structures that are related in one way or another to networking and network device drivers.

**Table 11-12** Commands to Display Network-Related Structures

Command	Operation
<i>ifnet addr</i>	Display the contents of the <i>ifnet</i> object at <i>addr</i> .
<i>rawcb addr</i>	Display the contents of the <i>rawcb</i> structure at <i>addr</i> .
<i>rawif addr</i>	Display the contents of the <i>rawif</i> structure at <i>addr</i> .
<i>sock addr</i>	Display the <i>sockbuf</i> structure at <i>addr</i> . When <i>addr</i> is positive, it is taken as a physical address; otherwise it is a kernel address.

## Using icrash

The *icrash* program is a post-mortem analysis tool for system crashes. You can use *icrash* to generate a wide variety of reports and displays based on a kernel panic dump from a crashed system. Study the *icrash(1M)* reference page for the current release. For example, you can display the *putbuf* message buffer using the *stat* command of *icrash*.



## Driver Example

This chapter displays the code of a complete device driver. The driver implements a “RAM drive,” a block of memory that simulates a disk drive. Since it has no hardware dependencies, this example driver can be used for experimentation in any IRIX system.

**Note:** It is not sensible to use a RAM drive in a system like IRIX, where there is an effective implementation of virtual memory. The RAM drive only occupies memory that is better used as buffers for the paging system. This driver is useful only as a test-bed for experiments with the driver-kernel interface and with *symmon* and other debugging tools. Do not use this driver in a production system.

### Installing the Example Driver

Use the following steps to install and test the example driver. Each step is expanded in the following topics.

1. Obtain the source code files.
2. Compile the source to obtain an object file.
3. Set up the appropriate configuration files.
4. Reboot the system and verify driver operation.
5. Install special device files and make and mount a filesystem.

## Obtaining the Source Files

The example driver consists of the following four source files:

<i>ramdrive.c</i>	Source code of the executable module.
<i>ramdrive.h</i>	Header file containing a few declarations.
<i>ramdrive</i>	Descriptive file to place in <i>/var/sysgen/master.d</i>
<i>ramdrive.sm</i>	Example VECTOR line for <i>/var/sysgen/system</i>

These files (and other example code from this book) are available from the Silicon Graphics FTP server, <ftp://ftp.sgi.com/sgi/>

Alternatively, a patient person could recreate the files by copying and pasting from the online manual. However, owing to bugs in the InSight viewer support for X-paste, this is an error-prone process and is not recommended in the current release.

## Compiling the Example Driver

Compile *ramdrive.c* using the techniques described under “Compiling and Linking” on page 230. An example compilation is shown in Example 12-1,

### Example 12-1 Compiling the Example Driver for a 32-bit Kernel

```
% set CFLAGS = "-DDEBUG -D_K32U32 -D_KERNEL -DSTATIC=static
-D_PAGESZ=4096 -D_MIPS2 -DIP22 -DR4000 -G 0 -non_shared -elf -xansi
-fullwarn -32 -mips2 -Wc,-pic0"
% cc $CFLAGS -c ramdrive.c
```

When the driver is compiled with the `-DDEBUG` option, all its informational displays are enabled. Without that option, it only displays messages related to serious errors during initialization.

## Configuring the Example Driver

Before you configure the example driver into the kernel, you should set the system with a debugging kernel, as described under “Preparing the System for Debugging” on page 245.

Configure the example driver to IRIX by copying files as follows:

- Copy the *ramdrive.o* file to */var/sysgen/boot*.
- Edit the *ramdrive* descriptive file and make any necessary changes. For example:
  - The file specifies major device number 77. If there is another driver in the system using this major number, change 77 to any unused number in the range 60-79 (see “Major Device Number” on page 36).
  - The fourth option, #DEV, is set to 4. This controls how many unique minor devices the driver supports.
- Copy the edited *ramdrive* descriptive file to */var/sysgen/master.d*.
- Edit the *ramdrive.sm* file and make any desired changes. For example,
  - Change a *base=* value to change the size of a RAM drive
  - Add or remove VECTOR lines to change the number of minor devices that are initialized.
- Copy the edited *ramdrive.sm* file to */var/sysgen/system*.
- Reboot the system.

If you compiled the example driver with `-DDEBUG`, it displays several informational lines to the system console from its `rd_init()`, `rd_start()`, and `rd_edtinit()` entry points. The display from `rd_edtinit()` gives the address of the *rd\_info\_t* structure. You can display that area with *idbg*. Example 12-2 shows the result of displaying a simulated volume header structure, using an address displayed from `rd_edtinit()` (the actual address will differ).

**Example 12-2** Displaying Simulated Volume Header Using idbg

```
# idbg hd 0x8841f604 0x80
0x8841f604: 0b e5 a9 41 00 00 00 01 00 00 00 00 00 00 00 00 ...A.....
0x8841f614: 00 00 00 00 00 00 00 00 00 00 00 00 02 00 00 00 .....
0x8841f624: 00 01 00 00 00 00 00 08 02 00 00 01 00 00 00 00 .....
0x8841f634: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 .....
0x8841f644: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 .....
0x8841f654: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 .....
0x8841f664: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 .....
0x8841f674: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 .....
0x8841f684: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 .....
0x8841f694: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 .....
0x8841f6a4: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 .....
0x8841f6b4: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 .....
0x8841f6c4: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 .....
0x8841f6d4: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 .....
0x8841f6e4: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 .....
0x8841f6f4: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 .....
0x8841f704: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 .....
0x8841f714: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 .....
0x8841f724: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 .....
0x8841f734: 00 00 00 00 00 00 00 00 00 00 0f ff 00 00 00 01 .....
0x8841f744: 00 00 00 03 00 00 00 00 00 00 10 00 00 00 00 03 .....
0x8841f754: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 .....
0x8841f764: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 .....
0x8841f774: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 .....
0x8841f784: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 0f ff .....
0x8841f794: 00 00 00 01 00 00 00 03 00 00 00 01 00 00 00 00 .....
0x8841f7a4: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 .....
0x8841f7b4: 00 00 10 00 00 00 00 00 00 00 00 06 00 00 00 00 .....
0x8841f7c4: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 .....
0x8841f7d4: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 .....
0x8841f7e4: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 .....
0x8841f7f4: 00 00 00 00 00 00 00 00 f0 19 16 a5 00 00 00 00 .....
```

## Creating Device Special Files

For each VECTOR line in */var/sysgen/system/ramdrive.sm* you can create a pair of device special files, one block device and one character device. Each file's inode contains the major device number established in descriptive file */var/sysgen/master.d/ramdrive*, and a minor number between 0 and 1 less than the limit established in the descriptive file.

The preferred command is *install* (see the *install(1)* reference page). The commands in Example 12-3 show the creation of the character and block devices for minor device number 0. The files created are */dev/ramblk0* and */dev/ramchr0*.

### Example 12-3 Install Command to Create Device Special File

```
# install -u root -g sys -f /dev -blk 77,0 ramblk0
# install -u root -g sys -f /dev -chr 77,0 ramchr0
```

In addition to device special files, you need to create ordinary directories that can be used as mount points for the mounted filesystem, for example

```
# mkdir /RAM0
```

The example driver supports as many minor device numbers as you specify in the descriptive file */var/sysgen/master.d/ramdrive*. You can create as many device special files that contain the example driver's major number as you like—or as few. Device special files are independent of the driver configuration until they are opened. However,

- A device with a minor number in excess of the limit set in */var/sysgen/master.d/ramdrive* cannot be initialized with a VECTOR line in */var/sysgen/system/ramdrive.sm*, nor can it be opened.
- A device whose minor number was not initialized by a VECTOR line cannot be opened.

## Verifying Driver Operation

You can verify operation of the driver by creating an EFS file system on one of its devices. As a first step, check the simulated volume header contents using *prtvtoc*, as shown in Example 12-4.

**Example 12-4** Applying prtvtoc to a RAM Drive of 2 MB

```
# prtvtoc /dev/ramchr0
ramdrive open: flag 1 copen 1 bopen 0 xopen 0
DIOCGETVH on 20185088
ioctl copy kmem 8841f604 -> usr 10006868 for 512
ramdrive close: flag 1 copen 0 bopen 0 xopen 0
* /dev/ramchr0 (bootfile "")
*      512 bytes/sector
*      8 sectors/track
*      1 tracks/cylinder
*      512 cylinders
*      0 cylinders occupied by header
*      512 accessible cylinders
* No space unallocated to partitions
Partition Type Fs      Start: sec   (cyl)      Size: sec   (cyl)  Mount
Directory
0          raw          1 ( 0.1)     4095 ( 511.9)
7          raw          1 ( 0.1)     4095 ( 511.9)
8          volhdr       0 ( 0)       1 ( 0.1)
10         volume       0 ( 0)       4096 ( 512)
```

You can make an EFS filesystem on a RAM drive device by applying *mkfs* to the character special device (*mkfs* is applied to the character device because it uses **ioctl()**, which is not supported for block devices). This is shown in Example 12-5. The voluminous debugging displays have been truncated in the display.

**Example 12-5** Making a Filesystem on a RAM Drive

```
# mkfs -t efs /dev/ramchr0
ramdrive open: flag 1 copen 1 bopen 0 xopen 0
DIOCGETVH on 20185088
ioctl copy kmem 8841f604 -> usr 10010c50 for 512
ramdrive close: flag 1 copen 0 bopen 0 xopen 0
ramdrive open: flag 1 copen 1 bopen 0 xopen 0
DIOCGETVH on 20185088
ioctl copy kmem 8841f604 -> usr 10010c50 for 512
ramdrive close: flag 1 copen 0 bopen 0 xopen 0
ramdrive open: flag 1 copen 1 bopen 0 xopen 0
DIOCGETVH on 20185088
ioctl copy kmem 8841f604 -> usr 10010c50 for 512
ramdrive close: flag 1 copen 0 bopen 0 xopen 0
ramdrive open: flag 1 copen 1 bopen 0 xopen 0
ramdrive close: flag 1 copen 0 bopen 0 xopen 0
ramdrive open: flag 3 copen 1 bopen 0 xopen 0
```

```

mkfs_efs: /dev/ramchr0: blocks=4095 inodes=616
mkfs_efs: /dev/ramchr0: sectors=8 cgfsz=4091
mkfs_efs: /dev/ramchr0: cgalgn=1 ialgn=1 ncg=1
mkfs_efs: /dev/ramchr0: firstcg=3 cgisz=154
mkfs_efs: /dev/ramchr0: bitmap blocks=1
rd_write entered for dev 20185088
rd_strategy: edev 20185088, flags 8018, blkno 3
      : offset 600, count 13400, dmaadr c026b610
      : write c026b610 to c0000600 for 13400
rd_write entered for dev 20185088
rd_strategy: edev 20185088, flags 18, blkno 0
      : offset 0, count 200, dmaadr 8b292cc0
      : write 8b292cc0 to c0000000 for 200
rd_read entered for dev 20185088
rd_strategy: edev 20185088, flags 19, blkno 3
      : offset 600, count 200, dmaadr c2c38a40
      : read c0000600 to c2c38a40 for 200
(many debugging displays omitted)
rd_write entered for dev 20185088
rd_strategy: edev 20185088, flags 8018, blkno ffe
      : offset 1ffc00, count 200, dmaadr c022def0
      : write c022def0 to c01ffc00 for 200
ramdrive close: flag 3 copen 0 bopen 0 xopen 0

```

After making a filesystem you can mount the filesystem, as shown in Example 12-6. You specify the block device when mounting an EFS filesystem, because the filesystem code calls the *pxsize()* and *pxstrategy()* driver entry points, which are supplied only by a block driver.

**Example 12-6** Mounting a RAM Drive Filesystem

```

# mount -t efs /dev/ramblk0 /RAM0
ramdrive open: flag 3 copen 0 bopen 1 xopen 0
rd_size entered for dev 20185088
rd_strategy: edev 20185088, flags 9, blkno 1
      : offset 200, count 200, dmaadr 889f5800
      : read c0000200 to 889f5800 for 200
rd_strategy: edev 20185088, flags 9, blkno 2
      : offset 400, count 200, dmaadr 889f5a00
      : read c0000400 to 889f5a00 for 200
rd_strategy: edev 20185088, flags 8, blkno 2
      : offset 400, count 200, dmaadr 889f5a00
      : write 889f5a00 to c0000400 for 200

```

```
rd_strategy: edev 20185088, flags 9, blkno 3
             : offset 600, count 1000, dmaadr 88ef6000
             : read c0000600 to 88ef6000 for 1000
rd_strategy: edev 20185088, flags 1000008, blkno 1
             : offset 200, count 200, dmaadr 889f5800
             : write 889f5800 to c0000200 for 200
# df /RAM0
Filesystem          Type  blocks   use   avail  %use Mounted on
/dev/ramblk0       efs    3937    22   3915    1  /RAM0
```

## Example Driver Source Files

The four source files of the example driver are displayed in the following topics:

- “Descriptive File” on page 280 displays the `/var/sysgen/master.d` file that describes the driver to `lboot`.
- “System File” on page 281 displays the `/var/sysgen/system` file that contains the VECTOR statements to initialize the driver.
- “Header File” on page 281 displays the driver’s header file.
- “Source File” on page 283 displays the source file.

### Descriptive File

```
*
* IRIX 6.2 Example driver "ramdrive" -- not for production use
*
* Flags used:
* b: block type device
* c: character type device (yes, both)
* d: dynamically loadable kernel module
* n: driver is semaphored
* s: software device driver
* w: driver is prepared to perform any cache write back operation
*
* External major number (SOFT) is an arbitrary choice from
* the range of numbers reserved for customer drivers.
*
* #DEV is passed in to the driver and used to configure its info array.
*
```

```
*FLAG  PREFIX  SOFT    #DEV    DEPENDENCIES
*bcdnsw rd_     77      4
bcnsw  rd_     77      4
```

```
$$$
```

```
int rd_e_major = ##E;
int rd_numdevs = ##D;
int rd_ctrlrs = ##C;
```

## System File

```
*
* Lboot config file for IRIX 6.2 example driver "ramdrive"
*   base= size of RAM drive, e.g. 0x00200000 is 2MB
*   ctrlr= minor device number, 0 to 3
* Store as /var/sysgen/system/ramdrive.sm
*
VECTOR: module=ramdrive ctrlr=0 base=0x00100000
*VECTOR: module=ramdrive ctrlr=1 base=0x00080000
```

## Header File

```

/*****
 *
 *          Copyright (C) 1993, Silicon Graphics, Inc.
 *
 * These coded instructions, statements, and computer programs contain
 * unpublished proprietary information of Silicon Graphics, Inc., and
 * are protected by Federal copyright law. They may not be disclosed
 * to third parties or copied or duplicated in any form, in whole or
 * in part, without the prior written consent of Silicon Graphics, Inc.
 *
 *****/

/*****
| This is ramdrive.h, containing declarations that are used in both the
| driver and in the unit-test application code. Aside from the unit-test
| module, no user-level code ever needs these declarations. The ram drive
| is accessed through the file system like any other disk.
 *****/

/*****
```

```

| The driver name for lboot and configuration is "ramdrive."
| The driver prefix is "rd_"
*****/

#define DRIVER_PFX "rd_"
#define DRIVER_NAME "ramdrive"

/*****
| MAX_RD_DEVS declares the maximum number of distinct devices supported.
| Ram drive device special files are /dev/dsk/ramblk<n> for block devices
| and /dev/rdisk/ramchr<n> for character devices. In each case <n> is the
| device minor number, between 0 and MAX_RD_DEVS-1.
*****/

#define MAX_RD_DEVS 4

/*****
| An array of MAX_RD_DEVS structures of the following type is maintained
| in the driver. VECTOR lines for up to MAX_RD_DEVS devices are written
| in /var/sysgen/system/ramdrive.sm, causing that many entries to the
| rd_edtinit() entry point, each entry initializing one structure.
|
| base:    address of allocated memory for the "drive." If NULL, this
|          minor number has not been initialized or failed initialization.
|
| size:    size of the allocated memory in bytes, always rounded down to
|          a multiple of IO_NBPP.
|
| copen:   count of successful character opens, cleared to 0 in rd_close().
| bopen:   count of successful block opens (0 or 1), cleared in rd_close().
| xopen:   nonzero when an FEXCL open has succeeded.
|
| mmap:    count of rd_map() entries, decremented in rd_unmap(). When the
|          any of copen, bopen, and mmap over all devices is nonzero, the
|          driver returns EBUSY to the rd_unload() entry point.
|
| queue:   semaphore used to serialize access for reading and writing.
|          (Note however that in a multiprocessor, a user process can
|          perform an unsynchronized write to a mapped character device.)
|
| vh:     volume header structure prepared in edtinit(), and used to
|          initialize block 0 when "formatting" the drive. The block-0
|          version can be modified by /etc/mkfs.
*****/

```

```

typedef struct rd_info {
    caddr_t    *base;
    off_t      size;
    __uint32_t  copen, bopen, xopen, mmap;
    sema_t     queue; /* requires sys/sema.h */
    struct volume_header vh; /* requires sys/dvh.h */
} rd_info_t;

```

## Source File

```

/*****
 *
 *          Copyright (C) 1993, Silicon Graphics, Inc.          *
 *
 * These coded instructions, statements, and computer programs contain *
 * unpublished proprietary information of Silicon Graphics, Inc., and *
 * are protected by Federal copyright law. They may not be disclosed *
 * to third parties or copied or duplicated in any form, in whole or *
 * in part, without the prior written consent of Silicon Graphics, Inc. *
 *
 *****/
/*****
| This sample IRIX device driver implements a "ram disk" -- a block of
| kernel memory accessed as if it were a disk. The driver supports both
| block and character interfaces and is loadable and unloadable.
|
| N N OO TTITT EEEEE It does not make sense to use a ram disk
| NN N O O T E in a system like IRIX that implements
| N N N O O T EEE :: effective virtual memory. This device
| N NN O O T E driver is useful as an example because
| N N OO T EEEEE :: it has no hardware dependencies, and so
| can be tried out in any IRIX system.
| However, this driver SHOULD NOT be employed in a production system!
| It WILL NOT give better performance. It WILL consume kernel memory
| that would be better used for buffers.
| *****/

#include <sys/ddi.h> /* gets also sys/types.h and sys/buf.h */
#include <sys/conf.h> /* for driver flags D_MP etc */
#include <sys/kmem.h> /* kmem_alloc and friends */
#include <sys/sema.h> /* the rd_info_t contains a semaphore */
#include <sys/dvh.h> /* the rd_info_t contains struct volume_header */
#include "ramdrive.h" /* declare rd_info_t, etc. */
#include <sys/edt.h> /* declare edt_t for edtinit() */

```

```

#include <sys/errno.h>      /* error codes to return */
#include <sys/cmn_err.h>    /* cmn_err() and related constants */
#include <sys/cred.h>      /* cred_t for prototypes */
#include <sys/dkio.h>      /* DKIOC* constants for ioctl */
#include <sys/param.h>     /* NBPC bytes per sector */
#include <sys/inmu.h>      /* IONBPP bytes per I/O page, btod() */
#include <sys/file.h>      /* FEXCL and other open flags */
#include <sys/open.h>      /* OTYP_CHR, OTYP_BLK */
#include <sys/region.h>    /* for vhandl_t */
#include <sys/mload.h>     /* only for M_VERSION */
/*****
| Debug display macros: one each for cmn_err calls with 0, 1, 2, or 3
| variable arguments.
| *****/
#ifdef DEBUG
#define DBGMSG0(s) cmn_err(CE_DEBUG,s)
#define DBGMSG1(s,x) cmn_err(CE_DEBUG,s,x)
#define DBGMSG2(s,x,y) cmn_err(CE_DEBUG,s,x,y)
#define DBGMSG3(s,x,y,z) cmn_err(CE_DEBUG,s,x,y,z)
#define DBGMSG4(s,x,y,z,w) cmn_err(CE_DEBUG,s,x,y,z,w)
#else
#define DBGMSG0(s)
#define DBGMSG1(s,x)
#define DBGMSG2(s,x,y)
#define DBGMSG3(s,x,y,z)
#define DBGMSG4(s,x,y,z,w)
#endif
/*****
| Driver flag: this driver is MP-safe. Also version flag for mload.
| *****/
unsigned rd_devflag = D_MP;
char *rd_mversion = M_VERSION;
/*****
| Array of rd_info_t objects, one per allowed minor device. We rely on
| the loader to ensure these static globals are zero until initialized!
| Also defined: two convenience macros for frequent expressions.
| *****/
static rd_info_t *rd_array;
#define INFOPTR(dev) &rd_array[getemisor(dev)]
#define VALIDIO(prd,off,len) (((off_t)(off) + (off_t)(len)) <= prd->size)
/*****
| rd_basic() is called from rd_edtinit() to allocate the rd_array based
| on the global rd_numdevs, an integer set to ##D in the configuration
| file /var/sysgen/master.d/ramdrive. Also display the other available
| globals for debugging purposes.

```

```

|*****|
extern int rd_e_major, rd_numdevs, rd_ctrlrs;
int
rd_basic(void)
{
    if (!rd_array)
    {
        register int size;
        DBGMSG3("ramdrive basic: ##E=%d, ##D=%d, ##C=%d\n",
            rd_e_major, rd_numdevs, rd_ctrlrs);
        if (size = rd_numdevs*sizeof(rd_info_t))
            rd_array = (rd_info_t *)kmem_zalloc(size,KM_SLEEP);
        else
            cmn_err(CE_ALERT,"ramdrive: confused");
    }
    return (0 != rd_array);
}
|*****|
| rd_init() is included solely to demonstrate that this entry point
| can be called in addition to rd_edtinit() and rd_start().
|*****|
int
rd_init(void)
{
    DBGMSG0("rd_init entry point called\n");
    return 0;
}
|*****|
| rd_start() is included solely to demonstrate that it, too can be called
| in addition to rd_edtinit() and rd_init().
|*****|
int
rd_start(void)
{
    DBGMSG0("rd_start entry point called\n");
    return 0;
}
|*****|
| rd_format() is a subroutine of both rd_edtinit() and rd_ioctl() which
| "formats" the ramdrive to zeros with a reasonable volume header.
| The volume header (set in both the info struct and "sector 0")
| describes standard SGI partitions:
|   10 == the whole "drive"
|   8 == the volume header, only one sector in this case
|   7 == all sectors except the volume header

```

```

|     0 == data ("root") same as 7
|     1 == swap contains 0 sectors
| For versimilitude we arbitrarily say we have 1 track/cylinder
| and 8 sectors/track. This assumes that nsectors is a multiple of 8,
| which is a good bet when the allocated size is a multiple of IO pages
| and sectors are 512 bytes.
| *****/
void
rd_format(register rd_info_t *prd)
{
    register struct volume_header *pvh = &prd->vh;
    register int nsectors = btod(prd->size); /* immu.h */

    bzero((void *)pvh, sizeof(struct volume_header));
    pvh->vh_magic = VHMAGIC; /* in sys/dvh.h */
    pvh->vh_rootpt = 0;
    pvh->vh_swappt = 1;
    pvh->vh_dp.dp_cyls = nsectors/8; /* number of cylinders */
    pvh->vh_dp.dp_trks0 = 1; /* tracks/cyl */
    pvh->vh_dp.dp_secs = 8; /* sectors/track */
    pvh->vh_dp.dp_secbytes = NBPSCTR; /* param.h */
    pvh->vh_dp.dp_interleave = 1;

    pvh->vh_pt[10].pt_firstlbn = 0;
    pvh->vh_pt[10].pt_nblks = nsectors;
    pvh->vh_pt[10].pt_type = PTYPE_VOLUME;
    pvh->vh_pt[ 8].pt_firstlbn = 0;
    pvh->vh_pt[ 8].pt_nblks = 1;
    pvh->vh_pt[ 8].pt_type = PTYPE_VOLHDR;
    pvh->vh_pt[ 8].pt_firstlbn = 0;
    pvh->vh_pt[ 7].pt_firstlbn = 1;
    pvh->vh_pt[ 7].pt_nblks = nsectors-1;
    pvh->vh_pt[ 7].pt_type = PTYPE_RAW;
    pvh->vh_pt[ 0] = pvh->vh_pt[ 7];
    pvh->vh_pt[ 1].pt_firstlbn = nsectors;
    pvh->vh_pt[ 1].pt_nblks = 0;
    pvh->vh_pt[ 1].pt_type = PTYPE_RAW;

    pvh->vh_csum = -vh_checksum(pvh);
    bzero(prd->base, prd->size); /* clear all sectors */
    bcopy(pvh, prd->base, sizeof(prd->vh)); /* vh in sec 0 */
}
/*****
| rd_edtinit() is called whenever the driver is loaded, once for each
| VECTOR that names this driver. A typical VECTOR line would be:

```

```

|     VECTOR module=ramdrive ctrl=2 base=0x00040000
| which says, initialize minor number 2 for a size of 256K.
|*****
int
rd_edtinit(register edt_t *pedt)
{
    register rd_info_t *prd;
    register __psint_t size;
    register int nsectors;
    register int ctrl = pedt->e_ctrl;
    /*
    || If this is the first time, allocate the rd_array of info structures.
    || Exit immediately if that fails.
    */
    if (!rd_basic())
    {
        return ENODEV;
    }
    DBGMSG3("ramdrive edtinit bustype %d adap %d ctrl %d\n",
            pedt->e_bus_type, pedt->e_adap, pedt->e_ctrl);
    DBGMSG3("          e_space[0] iopaddr %x size %x vaddr %x\n",
            pedt->e_space[0].ios_iopaddr, pedt->e_space[0].ios_size,
            pedt->e_space[0].ios_vaddr);
    /*
    || Diagnose and reject an invalid minor dev# from VECTOR ctrl=
    */
    if (ctrl > rd_numdevs)
    {
        cmn_err(CE_ALERT, "ramdrive: ctrl=%d invalid minor dev#", ctrl);
        return ENODEV;
    }
    /*
    || Address the info structure and diagnose multiple initialization
    */
    prd = INFOPTR(ctrl);
    if (prd->base)
    {
        cmn_err(CE_ALERT, "ramdrive: duplicate VECTOR for ctrl=%d", ctrl);
        return EBUSY;
    }
    /*
    || The desired size of the ramdrive is encoded as the base=# value,
    || which is passed as the ios_vaddr value in the edt_t.
    || Diagnose 0 size (omitted base=). Round the size to a
    || multiple of *memory* (not necessarily I/O) pages.

```

```

*/
size = (__psint_t) pedt->e_space[0].ios_vaddr;
if ((0 == size)||(-1 == size))
{
    cmn_err(CE_ALERT,
        "ramdrive: no size (base=) specified for ctlr=%d",ctlr);
    return EINVAL;
}
size = (size + (NBPP-1)) & (-NBPP); /* in sys/inmu.h */
/*
|| Allocate the kernel memory. Report an error if not possible.
*/
prd->size = size;
prd->base = kmem_alloc(size,KM_SLEEP);
if (!prd->base)
{
    cmn_err(CE_ALERT,
        "ramdrive: unable to allocate %x bytes for dev %d",size,ctlr);
    return ENOMEM;
}
nsectors = btod(size); /* inmu.h bytes to disk sectors */
DBGMSG3("ramdrive: dev# %d allocated %x = %x sectors\n",
        ctlr,size,nsectors);

/*
|| Initialize the semaphore.
*/
initnsema(&prd->queue,1,"ramdrive");
/*
|| Initialize the "volume."
*/
rd_format(prd);
DBGMSG2("          info at 0x%x  vh at 0x%x \n",
        prd, (__psint_t)(&prd->vh) );
return 0;
}
/*****
| rd_open() is called for each open() of a character device /dev/ramchr<n>,
| and during a mount of a block device /dev/ramblk<n>. We can distinguish
| between types of open from the otyp.
| *****/
int
rd_open(dev_t *pdev, int oflag, int otyp, cred_t *pcred)
{
    register rd_info_t *prd = INFOPTR(*pdev);
    register int error = 0;

```

```
/*
|| Make sure the device being opened was initialized by a VECTOR.
*/
if (!prd->base)
{
    cmn_err(CE_NOTE,"ramdrive: open of uninitialized dev %d",*pdev);
    return ENODEV;
}
/*
|| Seize the device semaphore so that prd->rd_info can be updated
|| without error on a multiprocessor.
*/
psema(&prd->queue,PZERO+1 | PCATCH);
/*
|| Implement FEXCL (exclusive) open for a privileged process only.
|| Exclusivity applies to the entire minor device, under both its
|| block and character special devices.
*/
if (oflag & FEXCL)
{
    if (drv_priv(pcred)) /* not privileged */
    {
        DBGMSG0("ramdrive: reject FEXCL with EPERM\n");
        error = EPERM;
    }
    else if (prd->copen+prd->bopen+prd->nmmmap) /* current use? */
    {
        DBGMSG0("ramdrive: reject FEXCL with EBUSY\n");
        error = EBUSY;
    }
    else
    {
        prd->xopen = oflag; /* note device open exclusively */
    }
}
else /* nonexclusive request can be blocked by exclusive open */
{
    if (prd->xopen)
    {
        DBGMSG0("ramdrive: reject normal open for exclusivity\n");
        error = EBUSY;
    }
}
if (!error)
{
```

```

    /*
    || Count the open so we don't unload with open devices.
    */
    if (otyp & OTYP_CHR)
        ++prd->copen;
    else
        ++prd->bopen;
    DBGMSG4("ramdrive open: flag %x copen %d bopen %d xopen %d\n",
            oflag, prd->copen, prd->bopen, prd->xopen);
}
vsema(&prd->queue);
return error;
}
/*****
| rd_close() is not called for each close() but for the final close of a
| given device (character or block). Clear the respective count of opens
| and note whether exclusivity is being given up. Since a close() in
| one CPU could happen concurrently with an open() in another CPU, we
| need to grab the semaphore before updating the rd_info.
| NOTE: the flag passed to close does not contain FEXCL even if it was
| given in the flag passed to open.
| *****/
int
rd_close(dev_t dev, int flag, int otyp, cred_t *pcred)
{
    register rd_info_t *prd = INFOPTR(dev);
    psema(&prd->queue, PZERO+1 | PCATCH);
    if (flag & FEXCL)
    {
        /* this is never entered */
    }
    if (otyp & OTYP_CHR)
    {
        prd->copen = 0;
    }
    else
    {
        prd->bopen = 0;
    }
    /* if all opens are closed, an exclusive one is closed */
    prd->xopen = 0;
    vsema(&prd->queue);
    DBGMSG4("ramdrive close: flag %x copen %d bopen %d xopen %d\n",
            flag, prd->copen, prd->bopen, prd->xopen);
    return 0;
}

```

```

}
/*****
| rd_ioctl() is called for ioctl(2), which can only be used on a character
| device. Disk ioctl command numbers for are in sys/dkio.h.
| DIOCREADCAPACITY: supported just for fun.
| DIOCGETVH: supported because /etc/mkfs and other tools use it (which
| explains why you apply mkfs to the character, not the block, device).
| DIOCSETVH: allows a program to change the "volume header" info.
| DIOCFORMAT: clears the device contents to 0, rewrites the vol header.
|
| The DIOC(S|G)ETVH calls use only the info in the per-device structure
| in memory. We make no attempt to keep that info in step with the
| contents of sector 0 of the simulated media. This is consistent with
| other current IRIX disk drivers. This has the implications that:
|   - you can change the driver's idea of the disk geometry on the fly,
|     without actually formatting the disk, this is useful for scsi.
|   - if you want to make a permanent change in the volume header,
|     -- one, that's a bad idea, use dvhtool(1) instead, but
|     -- two, if you insist, you need both a write to sector 0 and
|       a call to ioctl(,DIOCSETVH) to keep the driver up to date.
|
| Neither DIOCSETVH nor DIOCFORMAT hold the semaphore. You are strongly
| advised to do an exclusive open before calling them (but mkfp doesn't).
| *****/
int
rd_ioctl(dev_t dev, int cmd, caddr_t arg, int mode, cred_t *pcred, int *rval)
{
    register rd_info_t *prd = INFOPTR(dev);
    register int error = 0;
    register caddr_t kmemadr;
    register int len = 0;
    register int dir = 0; /* copyout */
    int capacity;
    switch(cmd)
    {
    case DIOCGETVH:
        {
            kmemadr = (caddr_t>(&prd->vh);
            len = sizeof(prd->vh);
            DBGMSGL("DIOCGETVH on %d\n",dev);
            break;
        }
    case DIOCREADCAPACITY:
        {
            capacity = prd->size/NBPSCTR;

```

```
        kmemadr = (caddr_t)&capacity;
        len = sizeof(capacity);
        DBGMSG2("DIOCREADCAPACITY on %d = %d\n",
                dev,capacity);
        break;
    }
case DIOCSETVH:
    {
        kmemadr = (caddr_t)&prd->vh;
        len = sizeof(prd->vh);
        dir = 1; /* copyin */
        DBGMSG1("DIOCSETVH on %d done\n",dev);
        break;
    }
case DIOCFORMAT:
    {
        rd_format(prd);
        DBGMSG1("DIOCFORMAT done on %d!\n",dev);
        break;
    }
default:
    {
        DBGMSG2("ramdrive invalid ioctl %x on %d\n",cmd,dev);
        error = EINVAL;
    }
} /* switch(cmd) */
/*
|| Perform the copy to or from user space if needed.
*/
if ((!error) && (len))
{
    if (!dir)
    {
        DBGMSG3("ioctl copy kmem %x -> usr %x for %d\n",
                kmemadr, arg, len);
        error = copyout(kmemadr,arg,len);
    }
    else
    {
        DBGMSG3("ioctl copy usr %x -> kmem %x for %d\n",
                arg, kmemadr, len);
        error = copyin(arg,kmemadr,len);
    }
}
#ifdef DEBUG
    if (error)
```

```

        DBGMSG1("error %d on ioctl copy\n",error);
#endif
    }
    *rval = error; /* ensure user gets correct code */
    return error;
}
/*****
| I/O Operations:
|
| rd_strategy() performs all actual I/O.  Called directly by file systems
| to read and write full I/O page units aligned on I/O page boundaries.
| Called indirectly to implement character I/O in any length and alignment.
|
| rd_read() and rd_write are called by read()/write() to a character
| device.  They defer to rd_strategy via uiophysio().  This is consistent
| with the operation of other IRIX disk drivers.
|
| The strategy code simply does a bcopy.  This is highly unrealistic.
| A real device driver would have to deal with efficient sequencing of
| track numbers and with asynchronous interrupts.
| *****/
int
rd_strategy(register struct buf *pbuf)
{
    register rd_info_t *prd = INFOPTR(pbuf->b_edev);
    register __psint_t offset = pbuf->b_blkno * NBPSCTR;
    register __psint_t count = pbuf->b_bcount;
    register caddr_t target = (caddr_t)((__psint_t)prd->base)+offset;
    DBGMSG3("rd_strategy: edev %d, flags %x, blkno %x\n",
            pbuf->b_edev,pbuf->b_flags,pbuf->b_blkno);
    DBGMSG3("          : offset %x, count %x, dmaadr %x\n",
            offset,count,(caddr_t)pbuf->b_dmaaddr);
    if (!VALIDIO(prd,offset,count))
    {
        DBGMSG0("rejecting strategy with ENOSPC\n");
        pbuf->b_error = ENOSPC;
        iodone(pbuf);
        return 0;
    }
    /*
    || Ensure that pbuf->b_dmaaddr is a valid kernel address.
    || This is never needed when called via uiophysio, only when
    || called from the file system or paging subsystem.  (Goodness!
    || wouldn't it be fun to use a ramdrive for swapping?)
    || NOTE: while a simple bp_mapin() call works, this approach

```

```
    || would impose unnecessary overhead in a real driver when
    || the device does not support scatter/gather.
    */
    if (!BP_ISMAPPED(pbuf))
    {
        bp_mapin(pbuf);
        DBGMSG1("          : after bp_mapin dmaadr %x\n", pbuf->b_dmaaddr);
    }
    /*
    || Grab the device semaphore. Note: this ensures consistency
    || between reads and writes, but does not control modifications
    || made through memory-mapped access.
    */
    psema(&prd->queue,PZERO+1 | PCATCH);
    /*
    || Perform the "read" or "write."
    */
    if (pbuf->b_flags & B_READ)
    {
        DBGMSG3("          : read %x to %x for %x\n",
                target,pbuf->b_dmaaddr,pbuf->b_bcount);
        bcopy(target,pbuf->b_dmaaddr,pbuf->b_bcount);
    }
    else
    {
        DBGMSG3("          : write %x to %x for %x\n",
                pbuf->b_dmaaddr,target,pbuf->b_bcount);
        bcopy(pbuf->b_dmaaddr,target,pbuf->b_bcount);
    }
    vsema(&prd->queue);
    iodone(pbuf);
    return 0;
}
int
rd_read(dev_t dev, uio_t *puio, cred_t *pcred)
{
    DBGMSG1("rd_read entered for dev %d\n",dev);
    return uiophysio(rd_strategy,0,dev,B_READ,puio);
}
int
rd_write(dev_t dev, uio_t *puio, cred_t *pcred)
{
    DBGMSG1("rd_write entered for dev %d\n",dev);
    return uiophysio(rd_strategy,0,dev,B_WRITE,puio);
}
```

```

int
rd_size(dev_t dev)
{
    DBGMSG1("rd_size entered for dev %d\n",dev);
    return rd_array[getemisor(dev)].size/NBPSCTR;
}
/*****
| Memory mapping: rd_map() (one "m") is called to implement an mmap()
| request on a character device. We permit read and write mappings, which
| means that in a multiprocessor, one CPU could be updating the kernel
| memory that represents the medium while another CPU executes a read()
| on the same memory.
|
| Since a map can persist after the corresponding FD is closed, we
| keep track of mappings separately from opens.
*****/
int
rd_map(dev_t dev, vhandl_t *pvh, off_t off, int len, int prot)
{
    register rd_info_t *prd = INFOPTR(dev);
    int error;

    DBGMSG3("map request on %d at %x for %x\n",dev,off,len);
    if (VALIDIO(prd,off,len))
    {
        error = v_mapphys(pvh,prd->base+off,len);
#ifdef DEBUG
        if (error)
            DBGMSG1("v_mapphys returns %d\n",error);
#endif
    }
    else
    {
        DBGMSG0("rejecting map with ENOSPC\n");
        error = ENOSPC;
    }
    if (!error)
        ++prd->nmmmap;
    return error;
}
rd_unmap(dev_t dev, vhandl_t *pvh)
{
    register rd_info_t *prd = INFOPTR(dev);
    if (prd->nmmmap)
    {

```

```
        --prd->nmmmap;
        DBGMSG2("unmap on %d, map count now %d\n",dev,prd->nmmmap);
    }
    else
    {
        DBGMSG1("unmap on %d when map count 0 ?!?!?!?\n",dev);
    }
    return 0;
}
/*****
| Unload support: rd_unload() is called when ml(1) is asked to unload
| this driver. We test to make sure that none of our devices that have
| been initialized, are in use. When any are in use, we return EBUSY
| and so will not be unloaded.
*****/
int
rd_unload(void)
{
    int j;
    for (j = 0; j<rd_numdevs; ++j)
    {
        if (( rd_array[j].base )
            && ( rd_array[j].copen
                ||rd_array[j].bopen
                ||rd_array[j].nmmmap) )
        {
            DBGMSG1("rejecting unload because dev %d busy\n",j);
            return EBUSY;
        }
    }
    DBGMSG0("accepting unload, byeeeee\n");
    return 0;
}
```

**PART FIVE**

**SCSI Device Drivers**

**Chapter 13: SCSI Device Drivers**

Actual control of the SCSI bus is managed by one or more Host Adapter drivers. This chapter tells how SCSI device drivers use these facilities.



---

## SCSI Device Drivers

All Silicon Graphics systems support the Small Computer Systems Interface (SCSI) bus for the attachment of disks, tapes, and other devices. This chapter details the kernel-level support for SCSI device drivers.

If your aim is to control a SCSI device from a user-level process, this chapter contains some useful background information to supplement Chapter 5, "User-Level Access to SCSI Devices." If you are designing a kernel-level SCSI driver, this chapter contains essential information on kernel support. The major topics in this chapter are as follows:

- "SCSI Support in Silicon Graphics Systems" on page 300 gives an overview of the hardware and software support for SCSI.
- "Host Adapter Facilities" on page 302 documents the use of the host adapter driver to access a SCSI device.
- "Designing a SCSI Driver" on page 317 notes design differences from other driver types, and includes an example driver skeleton.
- "Example SCSI Device Driver" on page 318 lists a skeleton driver to illustrate the use of the interface.
- "Designing a Host Adapter Driver" on page 323 documents the facility for creating and installing customized host adapter drivers.
- "SCSI Reference Data" on page 325 tabulates SCSI codes and messages for reference.

In addition, you may want to review the following additional sources:

intro(7) reference page	Documents the naming conventions for disk and tape device special files.
dksc(7) reference page	Documents the Silicon Graphics disk volume partition layout and the ioctl support in the base-level SCSI drivers.
ANSI X3.131-1986 and X3T9.2/85-52 Rev 4B.	SCSI standards documents.
<a href="http://www.abekrd.co.uk/SCSI2/">http://www.abekrd.co.uk/SCSI2/</a>	Web page containing the complete SCSI-2 standard in HTML form.

## SCSI Support in Silicon Graphics Systems

All current Silicon Graphics systems rely on the SCSI bus as the primary attachment for disks and tapes. The IRIX kernel provides extensive support for OEM drivers for SCSI devices.

**Note:** As used here, the term “adapter” means a SCSI controller such as the Western Digital W93 chip, which attaches a unique chain of SCSI devices. In this sense, a SCSI adapter and a SCSI bus are the same. “Adapter number” is used instead of “bus number.”

### SCSI Hardware Support

The Silicon Graphics computer systems supported by IRIX 6.2 can attach multiple SCSI adapters, as follows:

- The Indy workstation has at least one SCSI adapter on its motherboard, and can have up to two additional adapters on a GIO option board.
- The Indigo<sup>2</sup> series supports two SCSI adapters on the motherboard.
- The Challenge S system has two SCSI adapters on the motherboard, and can have one or two additional on each of one or two additional GIO option boards, for a maximum of six adapters.

- The Challenge M system supports one SCSI adapter on the CPU board and can have up to two additional adapters on a GIO option board.
- The POWERchannel (IO3) boards used in the Crimson line support two SCSI adapters per board.
- The Power Channel-2™ (IO4) boards used in the Challenge and Onyx series support two SCSI adapters, plus many as six additional SCSI adapters on mezzanine cards, for a maximum of eight adapters per IO4. In addition, VME-SCSI adapters (*Jag* units) can be installed on the VME bus in these systems.

In all systems, DMA mapping hardware allows a SCSI adapter to treat discontinuous memory locations as if they were a contiguous buffer, providing scatter/gather support.

## IRIX Kernel SCSI Support

The IRIX kernel contains two levels of SCSI support. An inner SCSI driver, the *host adapter driver*, manages all communication with SCSI hardware adapters. The kernel-level SCSI device drivers for particular devices prepare SCSI commands and call on the host adapter driver to execute them. This design centralizes the management of SCSI adapters. Centralization is necessary because the use of the SCSI bus is multiplexed across many devices, while recovery and error-handling need central handling. In addition, use of the host adapter driver makes it simpler to write a SCSI device driver.

### Host Adapter Drivers

Different host adapter drivers are loaded, depending on the hardware in the system. Some examples of host adapter drivers are *wd93*, *wd95*, and *jag*.

The host adapter drivers support all levels of the SCSI standard: SCSI-1, the Common Command Set (CCS, superseded by SCSI-2), and SCSI-2. Not all optional features of the standard are supported. Different systems support different feature combinations (such as synchronous, fast, and wide SCSI), depending on the available hardware.

The host adapter drivers handle the low-level communication over the SCSI interface, such as programming the SCSI interface chip or board, negotiating synchronous or wide mode, and handling disconnect/reconnect.

A host adapter driver is not, strictly speaking, a proper device driver because it does not support all the entry points documented in Chapter 8, “Structure of a Kernel-Level Driver.” You can think of it as a specialized library module for SCSI-bus management or as a device driver, whichever you prefer. The software interface to the host adapter driver is documented under “Host Adapter Facilities” on page 302.

**Caution:** Connect/disconnect strategy is enabled on any SCSI bus by default (the option is controlled by a constant defined in the host adapter driver descriptive file in */var/sysgen/master.d*). When disconnect is enabled on a bus, and a device on that bus refuses to disconnect, it can cause timeouts on other devices.

### SCSI Device Drivers

SCSI device drivers handle high-level device management, primarily by setting up SCSI commands for the host adapter driver to execute, and by interpreting returned sense data. Examples of device drivers are *dksc*, *tpsc*, and *smfd*.

## Host Adapter Facilities

The principal difference between a SCSI driver and other kernel-level drivers is that, while other kinds of drivers are expected to control devices directly using PIO and DMA, a SCSI driver operates its devices indirectly, by making function calls to the host adapter driver. This section documents the functional interface to the host adapter driver.

### Purpose of the Host Adapter Driver

The reason that IRIX uses host adapter drivers is that the SCSI bus is shared among multiple devices of different types, each type controlled by a different driver. A disk, a tape, a CDROM, and a scanner could all be cabled from the same SCSI adapter. Each device has a different driver, but each driver needs to use the adapter, a single chip-set, to communicate with its device.

If IRIX allowed multiple drivers to operate the host adapter, there would be confusion and errors from the conflicting uses. IRIX puts the management of each host adapter under the control of a host adapter driver, whose job is to issue commands on its bus and report the results. The host adapter is tailored to the hardware of the particular host adapter and to the architecture of the host system.

The interface to the host adapter driver is the same no matter what type of hardware the adapter uses. This insulates the individual device drivers from details of the adapter hardware.

The driver for each type of device is responsible for preparing the SCSI command bytes for its device, for passing the command requests to the correct host adapter driver, and for interpreting sense and status data as it comes back.

## Host Adapter Concepts

IRIX 6.2 permits a total of 10 unique host adapter drivers—five supplied by Silicon Graphics and up to five from other vendors. Each host adapter is customized to manage one type of adapter hardware. Each adapter driver has an adapter type number that is declared in *sys/scsi.h*. The constant names are listed in Table 13-1.

**Table 13-1** Host Adapter Driver Classes

Driver Constant	Driver Description
SCSIDRIVER_NULL	No driver exists; invalid adapter number or nonexistent adapter.
SCSIDRIVER_WD93	The <i>wd93</i> driver, for adapters based on the Western Digital WD93 chip set.
SCSIDRIVER_JAG	The <i>jag</i> driver, for adapters based on the VME-SCSI bridge used in the Challenge and Onyx systems.
SCSIDRIVER_WD95	The <i>wd95</i> driver, for adapters based on the Western Digital WD95 chip set.
SCSIDRIVER SCIP	The <i>scip</i> driver, for adapters based on the augmented WD95 chip set used in Challenge and Onyx systems.
SCSIDRIVER_QL	The <i>ql</i> driver, for adapters based on the QLogic chip set.
SCSIDRIVER_3RD_PARTY_START	First number available for OEM host adapter drivers.
SCSIDRIVER_3RD_PARTY_END	Last number available for OEM host adapter drivers.

**Caution:** The constant names listed in Table 13-1 compile to different values in different hardware systems. For this reason, you should avoid using these names in your driver; if you use one, your driver object file has to be recompiled for each CPU type.

The *lboot* command loads a host adapter driver for each unique type of adapter in the system. *lboot* is directed by VECTOR statements in the */var/sysgen/system/irix.sm* file (see “Configuring a Kernel” on page 238).

You can examine VECTOR lines in */var/sysgen/system/irix.sm* to see how many adapters your system has, and which of the host adapter drivers listed in Table 13-1 is loaded for each one.

The adapter number, the target number, and the logical unit number are important parameters to all the functions of the host adapter driver.

### **Target Numbers**

The purpose of a host adapter driver is to carry communications between a device driver and a *target*. A target is a device on the SCSI chain that responds to SCSI commands. A target can be a single device, or it can be a controller that in turn manages other devices.

A target is identified by a number between 0 and 15. Normally this number is configured into the device with switches or jumpers. The SCSI controller, usually target number 0 but 7 for the jag controller, cannot be used as a target.

The target number must be conveyed to the device driver somehow. The target numbers of Silicon Graphics disk and tape devices are passed in the device minor number.

Not all adapters support the range of 0-15 targets. The Jaguar VME-SCSI unit contains two independent adapters, each supporting target numbers 0-7.

### **Logical Unit Numbers (LUNs)**

When the target is a controller, it manages one or more sub-devices, each one a *logical unit* of that target. The logical unit being addressed is identified by a logical unit number (LUN). When the target is a single device, its LUN is 0.

## Overview of Host Adapter Functions

IRIX 6.2 permits a total of 10 unique host adapter drivers, but each of the ten must provide the same functional interface, which is based on simple concepts. The interface to host adapter drivers is declared in *sys/scsi.h*. Each adapter driver must provide the functions listed in Table 13-2.

**Table 13-2** Host Adapter Function Summary

Function	Header Files	Can Sleep?	Purpose
<code>scsi_info(D3)</code>	<code>scsi.h</code>	Y	Issue the SCSI Inquiry command and return the results.
<code>scsi_alloc(D3)</code>	<code>scsi.h</code>	Y	Open a connection between a driver and a target device.
<code>scsi_free(D3)</code>	<code>scsi.h</code>	Y	Release connection to target device.
<code>scsi_command(D3)</code>	<code>scsi.h</code>	Y	Transmit a SCSI command on the bus and return results.
<code>scsi_abort()</code>	<code>scsi.h</code>	Y?	Transmit a SCSI ABORT command (see caution).
<code>scsi_dump()</code>	<code>scsi.h</code>	Y	?

The normal sequence of operations is as follows:

1. In the `pxopen()` entry point (or, rarely, in an initialization entry point), the device driver calls `scsi_info()` to test the device characteristics. The results verify that the target device exists and is of the expected type.
2. In the `pxopen()` entry point, the device driver calls `scsi_alloc()` to set up communications with the target device. This allocates resources in the host adapter driver.
3. In the `pxstrategy()` or `pxioctl()` entry points, the device driver constructs SCSI command strings and calls `scsi_command()` to have them executed.
4. In the `pxclose()` entry point, the device driver calls `scsi_free()` to release any held resources related to this device.

**Caution:** The program interface to the `scsi_abort()` and `scsi_dump()` functions is subject to change. There is no reference page for these functions. The `scsi_reset()` function that formerly existed has been removed.

## How the Host Adapter Functions Are Found

A SCSI device driver can be asked to manage devices on different adapters. But different adapters can use different hardware, and be managed by different host adapter drivers. When opening one device, the device driver might need to call `scsi_alloc()` as provided by the `wd93` driver. When opening a different device, the driver might need the `scsi_alloc()` function from the `jag` driver. How can a driver locate the correct host adapter function for a given device?

The answer is provided by a set of function vector tables that are indexed by adapter number, and that yield the address of the appropriate function for that adapter.

### Using the Function Vector Tables

The function vector tables are maintained by the `scsi` driver module and filled in by each host adapter driver as it is initialized. The vector tables are declared in `sys/scsi.h`. The declaration of table `scsi_command` is as follows:

```
extern void (*scsi_command[])(struct scsi_request *req);
```

This declaration states that `scsi_command` is an array of pointers to functions. Each function in the table has the prototype

```
void function(struct scsi_request *req);
```

Each table is an array of pointers to functions. Each array is indexed by the adapter type number. If `iAdapT` is an integer variable containing the adapter type number for a device, the following statements are valid calls to the host adapter functions (the function arguments are examined in detail in the following topics):

```
#include <sys/scsi.h>
pTargInfo = (*scsi_info[iAdapT])(iAdap, iTarg, iLun);
iAllocRet = (*scsi_alloc[iAdapT])(iAdap, iTarg, iLun, 0, NULL);
(void) (*scsi_command[iAdapT])(&request);
(void) (*scsi_free[iAdapT])(iAdap, iTarg, iLun, NULL);
```

Each statement is a function call, but in each case, the name of the function is replaced by an expression that indexes the appropriate table.

### Learning the Adapter Type Number

Clearly, a SCSI driver needs to know the adapter type number for each device that it manages. Otherwise it cannot call the host adapter functions to manage that device.

The adapter type number for each adapter in the system is stored in an array maintained by the *scsi* driver. The array is declared as follows in *sys/scsi.h*:

```
extern u_char scsi_driver_table[];
```

When indexed by the number of the adapter in use, this table returns the adapter type number of the host adapter driver for that adapter.

### Learning the Adapter Number

Now all that remains is for the device driver to learn the adapter number with each device that it manages. There are two simple ways to do this.

One method is to get the number in the *edt\_t* structure. When a device is configured using a VECTOR line, the VECTOR should contain an *adapter=n* parameter. This number is stored in the *e\_adap* field of the *edt\_t* structure that is passed to the *pfxedtinit()* entry point. Code to retrieve it in a hypothetical driver is shown in Example 13-1.

#### Example 13-1 Storing the Adapter Type Number in *pfxedtinit()*

```
#include <sys/scsi.h>
typedef struct devVital_s {
    uchar devAdapNum;
    uchar devAdapType;
...} devVital_t;
void hypo_edtinit(edt_t *edt)
{
    devVital_t *pVitals;
    ...
    pVitals->devAdapNum = edt->e_adap;
    pVitals->devAdapType = scsi_driver_table[edt->e_adap];
    ...
}
```

A second method is to get it from the device minor number. For all Silicon Graphics disk and tape devices, the adapter number is encoded into both the visible name and the minor number of the device special file. You can use the bits of the minor number of any device in a similar way (see “Minor Device Number” on page 37).

Under the second plan, **getemisor()** is used to extract the minor number from the *dev\_t* value passed to each entry point (see “Device Number Functions” on page 182). The adapter number is calculated by shifting and masking the minor number. Hypothetical example code is shown in Example 13-2. The code of Example 13-2 can be extended to macros for the logical unit and control unit in obvious ways.

**Example 13-2** Extracting an Adapter Number From a Minor Device Number

```
/* Hypothetical minor bits: 00 aaaaaaaa cccuuuuu */
#define MINOR_ADAP_SHIFT 8
#define MINOR_ADAP_MASK 0x00ff
#define MINOR_ADAP(devt) (MINOR_ADAP_MASK & \
                          (getemisor(devt) >> MINOR_ADAP_SHIFT))
```

When the adapter number is known, the expression to call a host adapter function can be converted to a macro as well, possibly making the code more readable. The macro in Example 13-3 encapsulates a call to **scsi\_alloc()**. This code takes advantage of the fact that the adapter number is an argument to the function in any case.

**Example 13-3** Macro to Encapsulate a Call to **scsi\_alloc()**

```
#define SCSI_ALLOC(adap,targ,lun,opt,func) \
    (*scsi_alloc[scsi_driver_table[adap]]) \
    (adap,targ,lun,opt,func)
```

It could be argued that the double indexing in Example 13-3 imposes needless overhead. An approach with minimum overhead is to reserve space in the device-information structure for four function addresses, and to store the addresses of the host adapter functions with the other unique device information when the device is initialized.

## Using **scsi\_info()**

Before a SCSI driver tries to access a device, it must call the host adapter **scsi\_info()** function. This function issues an Inquiry command to the specified adapter, target, and logical unit. If the Inquiry is not successful—or if the adapter, target, or LUN is invalid—the return value is NULL. Otherwise, the return value is a pointer to a *scsi\_target\_info* structure.

The SCSI driver can learn the following things from a call to **scsi\_info()**:

- If the return is `NULL`, there is a serious problem with the device or the information about it. Write a descriptive log message with **cmn\_err()** and return `ENODEV`.
- The *si\_inq* field points to the Inquiry bytes returned by the device. Examine them for device-dependent information.
- The value in *si\_maxq* is the default limit on pending SCSI commands that can be queued to this host adapter driver. (You can specify a higher limit to **scsi\_alloc()**.)
- Test the bits in *si\_ha\_status* for information about the capabilities and error status of the host adapter itself. The possible bits are declared in *sys/scsi.h*. For example, `SRH_NOADAPSYNC` indicates that the specified target, or possibly the host adapter itself, does not support synchronous transfer. Not all bits are supported by all host adapter drivers.

You can also call **scsi\_info()** at other times; some of the returned information can be useful in error recovery. However, be aware that **scsi\_info()** for some host adapters is slow, and can use serialized access to hardware.

### Using **scsi\_alloc()**

Depending on its particular design, the host adapter driver may need to allocate memory for data structures, DMA maps, or buffers for sense and inquiry data, before it is ready to execute commands to a particular target. The call to **scsi\_alloc()** gives the host adapter driver the opportunity to prepare in these ways.

Because the host adapter driver may allocate virtual memory, it may sleep. Some host adapter drivers allocate all the resources they need on the first call to **scsi\_alloc()** and do little or nothing on subsequent calls.

A SCSI device driver will typically call the **scsi\_alloc()** function from the *pxopen()* entry point. However, if the driver needs to issue commands to the device at initialization time, it would call **scsi\_alloc()**, use **scsi\_command()**, and call **scsi\_free()**, all within the *pxinit()* or *pxedtinit()* entry point.

A call to **scsi\_alloc()** specifies these parameters:

<i>adap, targ, lun</i>	Numbers that identify the device on the bus.
<i>option</i>	An integer comprising two parameters, a flag and a count.
<i>callback</i>	Address of a function to be called whenever sense data is gotten from the device.

The option parameter may include the `SCSIALLOC_EXCLUSIVE` flag to request exclusive use of the target. If another driver has allocated a path to the same device, **scsi\_alloc()** returns `EBUSY`. For example, a tape device driver might require exclusive access, while a disk device driver would not.

The option parameter may include `SCSIALLOC_NOSYNC` to specify that this device should not, or cannot, use synchronous transfer mode. That setting can be overridden for single commands by a flag to **scsi\_command()** (see Table 13-4 on page 312).

The option parameter can also include a small integer value indicating the maximum queue depth (the number of SCSI commands the driver would like to start before any have completed). The call to **scsi\_info()** returns the default queue depth that will be used if you do not include a nonzero value (typically 1).

The callback function address can be specified as `NULL`. The specified callback function is called only when sense data is gotten from the allocated device. Only one driver that allocates a path to a device can specify a callback function. If the path is not held exclusively, any other drivers must specify a null address for their callback functions.

A call to **scsi\_alloc()** might resemble the following:

```
extern void sense_callback(char *pSense);
ret = scsi_alloc[scsi_driver_table[myAdapNum]](
    myAdapNum, myTargNum, myLun,
    SCSIALLOC_NOSYNC | 16, /* flag + max queue depth */
    sense_callback);
```

### Using **scsi\_free()**

A SCSI driver typically calls **scsi\_free()** from the **pfxclose()** entry point. That is the time when the driver knows that no processes have the device open, so the host adapter should be allowed to release any resources it is holding just for this device.

In addition, `scsi_free()` releases the device for use by other drivers, if the driver had allocated it for exclusive use.

## Using `scsi_command()`

A SCSI device driver sends SCSI commands to its device by storing information in a `scsi_request` structure and passing the structure to the `scsi_command()` function for the adapter. The host adapter driver schedules the command on the SCSI bus that it manages, and returns to the caller. When the command completes, a callback function is invoked.

**Tip:** When debugging a driver using a debugging kernel (see “Preparing the System for Debugging” on page 245), you can display the contents of a `scsi_request` structure using `symmon` or `idbg` (see “Commands to Display I/O Status” on page 269).

## Input to `scsi_command()`

The device driver prepares the `scsi_request` fields shown in Table 13-3.

**Table 13-3** Input Fields of the `scsi_request` Structure

Field Name	Contents
<code>sr_ctrl</code>	The adapter number.
<code>sr_target</code>	The target number.
<code>sr_lun</code>	The logical unit number.
<code>sr_tag</code>	If this target supports the SCSI-2 tagged-queue feature, and this command is directed to a queue, this field contains the queue tag message. Constant names for queue messages are in <code>sys/scsi.h</code> : <code>SC_TAG_SIMPLE</code> and two others.
<code>sr_command</code>	Address of the bytes of the SCSI command to issue.
<code>sr_cmdlen</code>	The length of the string at <code>*sr_command</code> . Constants for the common lengths are in <code>sys/scsi.h</code> : <code>SC_CLASS0_SZ</code> (6), <code>SC_CLASS1_SZ</code> (10), and <code>SC_CLASS2_SZ</code> (12).
<code>sr_flags</code>	Flags for data direction and DMA mapping, see Table 13-4.
<code>sr_timeout</code>	Number of ticks (HZ units) to wait for a response before timing out. The host adapter driver supplies a minimum value if this field is zero or too small.

**Table 13-3 (continued)** Input Fields of the `scsi_request` Structure

Field Name	Contents
<code>sr_buffer</code>	Address of first byte of data. Must be zero when <code>sr_bp</code> is supplied and <code>SRF_MAPBP</code> is specified in <code>sr_flags</code> .
<code>sr_bufalen</code>	Length of data or buffer space.
<code>sr_sense</code>	Address of space for sense data, in case the command ends in a check condition.
<code>sr_senseslen</code>	Length of the sense area.
<code>sr_notify</code>	Address of the callback function, called when the command is complete. A callback address is required on all commands.
<code>sr_bp</code>	Address of a <code>buf_t</code> object, when the command is called from a block driver's <code>pxstrategy()</code> entry point and buffer mapping is requested in <code>sr_flags</code> .
<code>sr_dev</code>	Address of additional information that could be useful in the callback routine <code>*sr_notify</code> .

The callback function address in `sr_notify` must be specified. (Device drivers for versions of IRIX previous to 5.x may set a NULL in this field; that is no longer permitted.)

The possible flag bits that can be set in `sr_flags` are listed in Table 13-4.

**Table 13-4** Values for the `sr_flags` Field of a `scsi_request`

Flag Constant	Purpose
<code>SRF_DIR_IN</code>	Data will be received in memory. If this flag is absent, the command sends data from memory to the device.
<code>SRF_FLUSH</code>	The data cache for the buffer area should be flushed (for output) or marked invalid (for input) prior to the command. This flag should be used whenever the buffer is local to the driver, not mapped by a <code>buf_t</code> object. It causes no extra overhead in systems that do not require cache flushing.
<code>SRF_MAPUSER</code>	Set this flag when doing I/O based on a <code>buf_t</code> and <code>B_MAPUSER</code> is set in <code>b_flags</code> .
<code>SRF_MAP</code>	Set this flag when doing I/O based on a <code>buf_t</code> and the <code>BP_ISMAPPED</code> macro returns nonzero.
<code>SRF_MAPBP</code>	The <code>sr_bp</code> field points to a <code>buf_t</code> in which <code>BP_ISMAPPED</code> returns false. The host adapter driver maps in the buffer.

**Table 13-4 (continued)** Values for the `sr_flags` Field of a `scsi_request`

Flag Constant	Purpose
<code>SRF_AEN_ACK</code>	This request is an acknowledgment of an AEN (Asynchronous Event Notification) message from the target. Following an AEN, any command without this flag is rejected with status <code>SC_ATTEN</code> .
<code>SRF_NEG_SYNC</code>	Attempt to negotiate synchronous transfer mode for this command. Ignored by some host adapter drivers. Overrides <code>SCSIALLOC_NOSYNC</code> (see "Using <code>scsi_alloc()</code> " on page 309).
<code>SRF_NEG_ASYNC</code>	Attempt to negotiate asynchronous mode for this command. Ignored unless the device is currently using synchronous mode.

When none of the three flag values beginning `SR_MAP` are supplied, the `sr_buffer` address must be a physical memory address. The `SR_MAPUSER` and `SR_MAPBP` flags are normally used when the command is issued from a `pxstrategy()` entry point in order to read or write a buffer controlled from a `buf_t` object.

### Command Execution

The host adapter driver validates the contents of the `scsi_request` structure. If the contents are valid, it queues the command for transmission on the adapter and returns. If they are invalid, it sets a status flag (see Table 13-6), calls the `sr_notify` function, and returns.

In any event, the `sr_notify` function is called when the command is complete. This function can be called from the host adapter interrupt handler, so it must assume that it is called in interrupt state.

The device driver should wait for the notify function to be called. The usual way is to share a semaphore (see "Semaphores" on page 224), as follows:

- Prior to calling `scsi_command()`, initialize the semaphore to 0 (the semaphore is being used to wait for an event).
- Immediately after the call to `scsi_command()`, call `psema()` for the semaphore.
- In the notify function, call `vsema()` for the function.

If the request is valid, the device driver will sleep in the `psema()` call until the command completes. If the request is invalid, the semaphore will already have been posted when `psema()` is called.

When the device driver holds an exclusive lock prior to issuing the command, a synchronization variable provides an appropriate way to wait for command completion (see “Using Synchronization Variables” on page 222).

**Values Returned in a *scsi\_request* Structure**

The host adapter driver sets the results of the request in the *scsi\_request* structure. The *sr\_notify* function is the first to inspect the values summarized in Table 13-5.

**Table 13-5** Values Returned From a SCSI Command

Field Name	Purpose
<i>sr_status</i>	Software status flags, see Table 13-6.
<i>sr_scsi_status</i>	SCSI status byte, see Table 13-7.
<i>sr_ha_flags</i>	Host adapter status flags, see Table 13-8.
<i>sr_sensegotten</i>	When no sense command was issued, 0. When a sense command was issued following an error, the number of bytes of sense data received. When an error occurred during a sense command, -1
<i>sr_resid</i>	The difference between <i>sr_buflen</i> and the number of bytes actually transferred.

The *sr\_status* field should be tested first. It contains an integer value; the possible values are summarized in Table 13-6.

**Table 13-6** Software Status Values From a SCSI Request

Constant Name	Meaning
SC_GOOD	The request was valid and the command was executed. The command might still have failed; see <i>sr_scsi_status</i> .
SC_TIMEOUT	The device did not respond to selection within 250 milliseconds.
SC_HARDERR	A hardware error occurred; inspect <i>sr_senselen</i> to see how much sense data was received, if any.
SC_PARITY	SCSI bus parity error detected.
SC_MEMERR	System memory parity or ECC error detected.
SC_CMDTIME	The device responded to selection but the command did not complete before <i>sr_timeout</i> expired.

**Table 13-6 (continued)** Software Status Values From a SCSI Request

Constant Name	Meaning
SC_ALIGN	The buffer address was not aligned as required by the adapter hardware. Most Silicon Graphics adapters require word (4-byte) alignment.
SC_ATTEN	Either a unit attention was received, or this command follows an AEN and did not contain the SR_AEN_ACK flag (see Table 13-4).
SC_REQUEST	An error was detected in the input values; the command was not attempted. The error could be that <code>scsi_alloc()</code> has not been called; or it could be due to missing or incorrect values.

One or more bits are set in the `sc_scsi_status` field. This field represents the status following the requested command, when the requested command executes correctly. When the requested command ends with Check Condition status, a sense command is issued and the SCSI status following the sense is placed in `sc_scsi_status`. In other words, the true indication of successful execution of the requested command is a zero in `sr_sensegotten`, because this indicates that no sense command was attempted.

Possible values of `sc_scsi_status` are summarized in Table 13-7.

**Table 13-7** SCSI Status Bytes

Constant Name	Meaning
ST_GOOD	The target has successfully completed the SCSI command. If a check condition was returned, a sense command was issued. The <code>sr_sensegotten</code> field is nonzero when this was the case.
ST_CHECK	This bit is only set for the special case when a check condition occurred on a sense command following a check condition on the requested command. The <code>sr_sensegotten</code> field contains -1.
ST_COND_MET	Search condition was met.
ST_BUSY	The target is busy. The driver will normally delay and then request the command again.
ST_INT_GOOD	This status is reported for every command in a series of linked commands. Linked commands are not supported by Silicon Graphics host adapters.
ST_RES_CONF	A conflict with a reserved logical unit or reserved extent.

One or more bits can be set in *sr\_ha\_flags* to document a host adapter state or problem. These flags are summarized in Table 13-8.

**Table 13-8** Host Adapter Status After a SCSI Request

Constant Name	Meaning
SRH_CANTSYNC	Unable to negotiate synchronous mode.
SRH_SYNCXFR	Synchronous mode was used. If not set, asynchronous mode was used.
SRH_TRIEDSYNC	Synchronous mode negotiation was attempted; see the SRH_CANTSYNC bit for the result.
SRH_BADSYNC	When SRH_CANTSYNC is set, this bit indicates that the negotiation failed because the device cannot negotiate.
SRH_NOADAPSYNC	When SRH_CANTSYNC is set, this bit indicates that the host adapter does not support synchronous negotiation, or that the system has been configured to not use synchronous mode for this device.
SRH_WIDE	This adapter supports Wide mode.
SRH_DISC	This adapter supports Disconnect mode and is configured to use it.
SRH_TAGQ	This adapter supports tagged queueing and is configured to use it.
SRH_MAPUSER	This host adapter driver can map user addresses.

### Using `scsi_abort()`

The purpose of the `scsi_abort()` function is to issue a SCSI ABORT command to a specified target and logical unit. The prototype of the function is:

```
int (*scsi_abort[adapter-type])(struct scsi_request *req);
```

The only fields of the `scsi_request` that are input to this function are those that identify the device: `sr_ctlr`, `sr_target`, and `sr_lun`. The ABORT command is issued on the bus as soon as possible but there could be a delay if the bus is busy. Status is returned in `sr_status`. The function returns a nonzero value when the ABORT command is issued successfully, and a zero when the ABORT command fails (which probably indicates a serious bus problem).

### Using `scsi_reset()`

The purpose of `scsi_reset()` is to reset the adapter hardware and possibly the attached bus, for example by asserting the reset line on the bus for at least 25 microseconds. The prototype of the function is

```
int (*scsi_reset[adapter-type])(uchar adap);
```

The adapter number is reset and a nonzero value is returned. If the host adapter driver does not support this function, or if it is unable to reset the hardware, it returns 0.

## Designing a SCSI Driver

A kernel-level SCSI device driver has the driver architecture described in Chapter 8, “Structure of a Kernel-Level Driver,” and it uses the basic system services documented in Chapter 9, “Device Driver/Kernel Interface.” You prepare a SCSI driver and configure it into the kernel as described in Chapter 10, “Building and Installing a Driver.”

However, a SCSI driver uses additional services, including those of the host adapter driver, and its configuration is slightly different from other drivers.

### SCSI Driver Initialization

A SCSI driver can be included in the kernel through a VECTOR, INCLUDE, or USE line in the system file in `/var/sysgen/system` (see “Configuring a Kernel” on page 238). When included through a VECTOR line, the `pfxedtinit()` entry point is called for each VECTOR line given. The VECTOR can describe a logical unit or a control unit, according to your design choice. However, a VECTOR line for a high-level SCSI driver can not include a `probe` or `exprobe` operand, because the hardware is owned by the host adapter driver, not by the SCSI device driver.

When included through a USE line, a SCSI driver is initialized at its `pfxininit()` entry point. In this case, the driver must obtain the adapter number by some other means. (See “Initialization Entry Points” on page 147.)

## Opening a SCSI Device

When the `pxopen()` entry point is called, the SCSI driver uses the appropriate `scsi_info()` function to verify the device and get hardware dependent Inquiry data from it. If the device is operational, the driver calls `scsi_alloc()` to open a communications path to it.

## Accessing a SCSI Device

In general, it is simplest to put all access to a device within a `pxstrategy()` entry point, even in a character device driver. When the `pxread()`, `pxwrite()`, or `pxioctl()` entry point needs to read or write data, it can prepare a `uio_t` to describe the data, and call `uiophysio()` to direct the operation through the single `pxstrategy()` entry point.

The notify routine passed in the `sr_notify` field plays the same role as the `pxintr()` entry point in other device drivers. It is called asynchronously, when the SCSI command completes. It may not call a kernel function that can sleep. However, it does not have to be named `pxintr()`, and a SCSI driver does not have to provide a `pxintr()` entry point.

## Configuring a SCSI Driver

A SCSI driver can be either a block or a character driver, or it can support both interfaces. When preparing the descriptive file for `/var/sysgen/master.d`, you must use the `s` flag, specifying a software-only driver, and list `scsi` as a dependency, in the descriptive line in `/var/sysgen/master.d`. See “Describing the Driver in `/var/sysgen/master.d`” on page 235.

## Example SCSI Device Driver

The following example shows how a driver can communicate with a direct access SCSI device, such as a disk. This driver is simplified and does not do as much error checking as a real driver would do. Also, this example uses a single, global SCSI request structure that does not work in real drivers, since multiple reads or writes would overwrite a command in progress.

**Tip:** A more complete sample SCSI driver is available on the Developer’s Toolbox CD. See information about the Developer Program under “Developer Program” on page xxvii.

```

#include "sys/param.h"
#include "sys/types.h"
#include "sys/user.h"
#include "sys/buf.h"
#include "sys/errno.h"
#include "sys/cmn_err.h"
#include "sys/cred.h"
#include "sys/ddi.h"
#include "sys/system.h"
#include "sys/scsi.h"

int sdk_devflag = 0; /* not old, not _MP either */

#define ADAPT    0    /* SCSI host adapter */
#define TARGET  7    /* the disk will have target ID #7 */
#define LU      0    /* and logical unit #0 */
#define TIMEOUT (30*HZ) /* wait 30 secs for SCSI device to
                        respond */
#define DIRECTACCESS 0 /* First byte of inquiry cmdnd */

unchar scsi_read[] = {0x28, 0, 0, 0, 0, 0, 0, 0, 0, 0};
unchar scsi_write[] = {0x2a, 0, 0, 0, 0, 0, 0, 0, 0, 0};
int     sdk_inuse = 0;
int     sdk_driver;
struct scsi_target_info *sdk_info;
struct scsi_request sdk_req;
u_char  sdk_sensebuf[SCSI_SENSE_LEN]; /* SCSI_SENSE_LEN
                                       from scsi.h */

/* forward definitions*/
int sdk_strategy(struct buf *bp);
void sdk_notify(struct scsi_request *req);
/*
 * sdk_open - Open the SCSI device exclusively.
 *
 * Issue a SCSI inquiry command upon device and ensure
 * it is a direct access device.
 */
int
sdk_open(dev_t *devp, int flag, int otyp, cred_t *crp)
{
    if (sdk_inuse)
        return EBUSY;
    /* Get driver number */
    sdk_driver = scsi_driver_table[ADAPT];
    /*

```

```
    * Call through scsi_info to get inquiry data and to
    * find out if a device is at the address we want.
    */
    sdk_info = (*scsi_info[ sdk_driver ])( ADAPT, TARGET, LU );
    if ( sdk_info == NULL )
        return ENODEV;
    /*
    * Is it a direct access device? We could check the
    * entire inquiry buffer to ensure it is actually the
    * correct device.
    */
    if ( sdk_info->si_inq[0] != DIRECTACCESS )
        return ENXIO;
    /*
    * It's a direct access device (disk drive). Initialize
    * the connection to the host adapter driver.
    */
    if ( (*scsi_alloc[ sdk_driver ] )
          ( ADAPT, TARGET, LU, 1, NULL ) == 0 )
        return EBUSY;
    /*
    * We have successfully allocated a connection between
    * sdk and the host adapter driver. Initialize the
    * scsi_request structure, and mark the driver as being
    * in use.
    */
    sdk_inuse = 1;
    bzero( & sdk_req, sizeof( sdk_req ) );
    sdk_req.sr_ctlr = ADAPT;
    sdk_req.sr_target = TARGET;
    sdk_req.sr_lun = LU;
    sdk_req.sr_timeout = TIMEOUT;
    sdk_req.sr_sense = sdk_sensebuf;
    sdk_req.sr_senselen = sizeof( sdk_sensebuf );
    sdk_req.sr_notify = sdk_notify;

    return 0;
}
/* sdk_close - close the device and free the subchannel. */
int
sdk_close( dev_t dev, int flag, int otyp, cred_t *crp )
{
    (*scsi_free[ sdk_driver ])( ADAPT, TARGET, LU, NULL );
    sdk_inuse = 0;
    return 0;
}
```

```

}
/*
 * sdk_read - read from the SCSI device, ensuring an even
 * block count and a word-aligned address.
 */
sdk_read(dev_t dev, uio_t *uiop, cred_t *crp)

/*
 * sdk_write - write to the SCSI device, ensuring an even
 * block count and a word-aligned address.
 */
sdk_write(dev_t dev, uio_t *uiop, cred_t *crp)

/*
 * sdk_strategy - do the dirty work of the I/O.
 * Use either the SCSI read or write command as
 * appropriate.  Modify the block number and block counts
 * within the command buffer.  Simply return here;
 * physio( ) will wait for an iodone( ).
 */
int
sdk_strategy(struct buf *bp)
{
    int blkno, blkcount;
    /* Prime the subchannel communication block. */
    blkno = bp->b_blkno;
    blkcount = BTOBB(bp->b_bcount);
    sdk_req.sr_command = bp->b_flags & B_READ ?
        scsi_read : scsi_write;
    sdk_req.sr_command[2] = (char)(blkno>>24);
    sdk_req.sr_command[3] = (char)(blkno>>16);
    sdk_req.sr_command[4] = (char)(blkno>>8);
    sdk_req.sr_command[5] = (char) blkno;
    sdk_req.sr_command[7] = (char)(blkcount>>8);
    sdk_req.sr_command[8] = (char) blkcount;

    sdk_req.sr_cmdlen = SC_CLASS1_SZ;
    sdk_req.sr_flags = bp->b_flags & B_READ ? SRF_DIR_IN : 0;
    if (BP_ISMAPPED(bp)) {
        sdk_req.sr_buffer = bp->b_dmaaddr;
        sdk_req.sr_buflen = bp->b_bcount;
        sdk_req.sr_flags |= SRF_MAP;
    }
    else {
        sdk_req.sr_buffer = NULL;
    }
}

```

```
        sdk_req.sr_buflen = bp->b_bcount;
        sdk_req.sr_flags = SRF_MAPBP;
    }
    sdk_req.sr_bp = bp;    /* required for SRF_MAPBP, but a
                          * convenience in all cases */
    /* Perform the SCSI operation. */
    (*scsi_command[ sdk_driver ])( & sdk_req );
}

/*
 * sdk_notify - SCSI command completion notification routine
 *
 * Simply check for errors and wake up physio( ) with
 * an iodone( ) on the buffer.
 * Note that a more robust driver would be more thorough
 * about error handling by retrying errors, giving more
 * information about error types, etc.
 */
void
sdk_notify(struct scsi_request *req)
{
    register struct buf *bp = req->sr_bp;
    if ((req->sr_status != SC_GOOD) ||
        (req->sr_scsi_status != ST_GOOD) ||
        (req->sr_sensegotten < 0))
    {
        cmn_err(CE_NOTE,
            "sdk: Error: driver stat 0x%x, scsi stat 0x%x"
            " sensegotten %d\n", req->sr_status,
            req->sr_scsi_status, req->sr_sensegotten);
        bioerror(bp, EIO);
    }
    else if (req->sr_sensegotten > 0) {
        cmn_err(CE_NOTE, "sdk: Error: sensekey 0x%x\n",
            sdk_sensebuf[2] & 0x0F);
        bioerror(bp, EIO);
    }
    bp->b_resid = req->sr_resid;
    biodone(bp);
}
```

## Designing a Host Adapter Driver

IRIX 6.2 provides the ability to load and install third-party host adapter drivers. This section documents the special features of this type of driver.

### Overview of Host Adapter Driver Architecture

A host adapter driver is a low-level driver for a SCSI bus adapter. A host adapter driver is similar to other device drivers described in this book in many ways:

- Like other device drivers, it uses the kernel facilities described in Chapter 9, “Device Driver/Kernel Interface.”
- It is compiled and linked like other drivers (see Chapter 10, “Building and Installing a Driver.”).
- It is configured to the system using files in */var/sysgen/master.d*, and loaded by *lboot*. A host adapter driver should not be loadable; if it is loadable, it should not unload.
- Like other drivers, it can have entry points *pfxstart()* or *pfxedtinit()* for initialization, *pfxintr()* for interrupt handling, and *pfxhalt()* for shutdown.

Unlike other drivers, a host adapter driver does not provide any entry points for serving the needs of system functions, such as *pfxread()*, *pfxpoll()*, or *pfxstrategy()*. Instead, it supplies the entry points used by SCSI device drivers.

### Host Adapter Initialization

In its initialization, the host adapter driver does three things:

- initializes the adapter hardware it supports
- acquires an adapter type number
- stores pointers to its functions in the function pointer arrays

### Initializing the Hardware

If it is called at its `pfxedtinit()` entry point, the host adapter driver receives adapter hardware information in an `edt_t` structure. Integration of the driver in this way, using a VECTOR line, is preferred. It removes the need to hard-code device addresses; and it allows the use of an `exprobe` operand to load the driver only when its adapter hardware is present.

If it is not loaded by a VECTOR line, the driver must be loaded with a USE line in the system database (see “Configuring a Kernel” on page 238) and it must find out the address of its adapter hardware by other means.

The driver may also include a `pfxstart()` entry point for general initialization, including the two steps of acquiring a type number and setting up its entry point addresses.

### Acquiring a Type Number

Every host adapter driver must have a unique adapter type number. The type numbers for Silicon Graphics drivers are declared in `sys/scsi.h`. An OEM driver acquires a number dynamically, by calling the kernel function `get_driver_number()`. The prototype of this function is

```
uchar get_driver_number(void);
```

If successful, the function returns a number between `SCSIDRIVER_3RD_PARTY_START` and `SCSIDRIVER_3RD_PARTY_END`, inclusive (see Table 13-1 on page 303). If unsuccessful, it returns -1, and the driver cannot initialize.

### Storing Entry Point Addresses

After it has its type number, the driver can store the address of each of its functional entry points in the arrays used by its callers. For example it stores the address of its command execution function in the `scsi_command` array, indexed by its type number.

The driver must support the functions summarized in Table 13-2 on page 305. For each function there is a corresponding array of function pointers, in which the driver stores the address of its function, indexed by its driver type number.

## SCSI Reference Data

This section contains reference material in the following categories:

- “SCSI Error Messages” on page 325 describes the general form of messages written by host adapter drivers into the system log.
- “Adapter Error Codes (Table *scsi\_adaperrs\_tab*)” on page 326 lists the possible adapter error codes and their message strings.
- “SCSI Sense Codes (Table *scsi\_key\_msgtab*)” on page 327 lists the primary sense codes and the corresponding message strings.
- “Additional Sense Codes (Table *scsi\_addit\_msgtab*)” on page 328 lists the possible additional sense codes (ASCs) and their message strings.

### SCSI Error Messages

The host adapter drivers such as *wd93*, *wd95*, and *jag* send error messages to the system log using the **cmn\_err()** function (see “Producing Diagnostic Displays” on page 251).

These messages almost always contain the adapter number (sometimes called the bus number or controller number). They sometimes contain the number of the target device, and sometimes add the number of the logical unit that was addressed.

Messages from the *wd93* driver specify the adapter number as `BUS=n`. The target device is shown as `ID=n` and the logical unit as `LUN=n`.

Messages from the *wd95* and *jag* drivers contain one, two, or three or more decimal numbers. In all cases, the first number is the adapter number, the second is the target ID, and the third (when present) is the logical unit number.

When error messages list a sense code, refer to “SCSI Sense Codes (Table *scsi\_key\_msgtab*)” on page 327 and to “Additional Sense Codes (Table *scsi\_addit\_msgtab*)” on page 328.

When the error message reports an error from the adapter itself, refer to “Adapter Error Codes (Table *scsi\_adaperrs\_tab*)” on page 326.

## SCSI Error Message Tables

The *scsi* module contains a set of error message tables that you can use to generate error messages based on SCSI sense codes and other data. The contents of these tables is documented here for reference, and to assist in decoding messages from SCSI drivers.

Each table is an array of pointers to strings; for example the *scsi\_key\_msgtab* table is defined beginning as follows:

```
char *scsi_key_msgtab[SC_NUMSENSE] = {
    "No sense",          /* 0x0 */
    "Recovered Error",  /* 0x1 */
    ...};
```

Each of the tables is declared as extern in *sys/scsi.h*.

### Adapter Error Codes (Table *scsi\_adaperrs\_tab*)

The table with the external name *scsi\_adaperrs\_tab* contains message strings to document the adapter error codes that can be returned in the *scsi\_request.sr\_status* field (see Table 13-6). The *scsi\_adaperrs\_tab* table contains NUM\_ADAP\_ERRS entries (9, defined in *sys/scsi.h*). The first entry (index 0x0) contains a pointer to a null string. The other entries are documented in Table 13-9.

**Table 13-9** Adapter Error Codes

Adapter Error Code	Constant Name	Message Text
0x1	SC_TIMEOUT	Device does not respond to selection .
0x2	SC_HARDERR	Controller protocol error or SCSI bus reset.
0x3	SC_PARITY	SCSI bus parity error.
0x4	SC_MEMERR	Parity/ECC error in system memory during DMA.
0x5	SC_CMDTIME	Command timed out.
0x6	SC_ALIGN	Buffer not correctly aligned in memory.
0x7	SC_ATTEN	Unit attention received on another command causes retry.
0x8	SC_REQUEST	Driver protocol error.

**SCSI Sense Codes (Table *scsi\_key\_msgtab*)**

The table with the external name *scsi\_key\_msgtab* is indexed by the primary sense code. Its contents are listed in Table 13-10. The table contains SC\_NUMADDSSENSE entries (16, defined in *sys/scsi.h*), of which the last two should not occur.

**Table 13-10** Primary Sense Key Error Table

Sense Key	Message	Most Common Cause
0x0	No sense	No error information available.
0x1	Recovered error	The device recovered by itself.
0x2	Device not ready	No media, or drive not spun up.
0x3	Media error	An actual media problem.
0x4	Device hardware error	Usually a device hardware error.
0x5	Illegal request	Invalid command or data issued.
0x6	Unit attention	Device was reset or power-cycled.
0x7	Data protect error	Usually device is write-protected.
0x8	Unexpected blank media	Tried to read at end of a tape.
0x9	Vendor unique error	Varies.
0xA	Copy aborted	Copy command aborted by host (not used).
0xB	Aborted command	Target device aborted command.
0xC	Search data successful	Search data command OK (not used).
0xD	Volume overflow	Tried to write past EOT on tape.
0xE	Reserved (0xE)	0xE should not be seen.
0xF	Reserved (0xF)	0xF should not be seen.

**Additional Sense Codes (Table `scsi_addit_msgtab`)**

The table with the external name `scsi_addit_msgtab` is indexed by the Additional Sense Code (ASC) value, when one is present. The table contains `SC_NUMADDSSENSE` entries (0x71, defined in `sys/scsi.h`). Some values have no standard definition; for these, the table contains a NULL value. Therefore you should always test the table value for a valid pointer before using it to format a message.

Table 13-11 lists the contents of this message table. Undefined (NULL) table entries are omitted.

**Table 13-11** Additional Sense Code Table

ASC Value	Corresponding Message String
0x01	No index/sector signal
0x02	No seek complete
0x03	Write fault
0x04	Not ready to perform command
0x05	Unit does not respond to selection
0x06	No reference position
0x07	Multiple drives selected
0x08	LUN communication error
0x09	Track error
0x0a	Error log overflow
0x0c	Write error
0x10	ID CRC or ECC error
0x11	Unrecovered data block read error
0x12	No address mark found in ID field
0x13	No address mark found in Data field
0x14	No record found
0x15	Seek position error

**Table 13-11 (continued)** Additional Sense Code Table

<b>ASC Value</b>	<b>Corresponding Message String</b>
0x16	Data sync mark error
0x17	Read data recovered with retries
0x18	Read data recovered with ECC
0x19	Defect list error
0x1a	Parameter overrun
0x1b	Synchronous transfer error
0x1c	Defect list not found
0x1d	Compare error
0x1e	Recovered ID with ECC
0x20	Invalid command code
0x21	Illegal logical block address
0x22	Illegal function
0x24	Illegal field in CDB
0x25	Invalid LUN
0x26	Invalid field in parameter list
0x27	Media write protected
0x28	Media change
0x29	Device reset
0x2a	Log parameters changed
0x2b	Copy requires disconnect
0x2c	Command sequence error
0x2d	Update in place error
0x2f	Tagged commands cleared
0x30	Incompatible media

**Table 13-11 (continued)** Additional Sense Code Table

ASC Value	Corresponding Message String
0x31	Media format corrupted
0x32	No defect spare location available
0x33 <sup>a</sup>	Media length error
0x36	Toner/ink error
0x37	Parameter rounded
0x39	Saved parameters not supported
0x3a	Medium not present
0x3b	Forms error
0x3d	Invalid ID msg
0x3e	Self config in progress
0x3f	Device config has changed
0x40	RAM failure
0x41	Data path diagnostic failure
0x42	Power on diagnostic failure
0x43	Message reject error
0x44	Internal controller error
0x45	Select/reselect failed
0x46	Soft reset failure
0x47	SCSI interface parity error
0x48	Initiator detected error
0x49	Inappropriate/illegal message
0x4a	Command phase error
0x4b	Data phase error
0x4c	Failed self configuration

**Table 13-11 (continued)** Additional Sense Code Table

---

<b>ASC Value</b>	<b>Corresponding Message String</b>
0x4e	Overlapped commands attempted
0x53	Media load/unload failure
0x57	Unable to read table of contents
0x58	Generation (optical device) bad
0x59	Updated block read (optical device)
0x5a	Operator request or state change
0x5b	Logging exception
0x5c	RPL status change
0x5d	Self diagnostics predict unit will fail soon
0x60	Lamp failure
0x61	Video acquisition error/focus problem
0x62	Scan head positioning error
0x63	End of user area on track
0x64	Illegal mode for this track
0x70 <sup>b</sup>	Decompression error

---

a. Specified as tape only.

b. DAT only; may be in SCSI3.

## WD93 States and Phases

Some of the SCSI states and phases that can be detected by the *wd93* host adapter driver are listed in Table 13-12 for reference, in case they appear in a debugging log message. These states and phases are declared in the */usr/include/sys/wd93.h* header file. The comments in the table have been extracted from that file and supplemented with additional information.

“Out” is from the CPU to the SCSI device in these descriptions, and “receive” and “send” are also from the SCSI device point of view, since the target controls all the bus phases except for initial selection.

**Table 13-12** SCSI State Error Messages

State Message	Sense Key	Comments
ST_RESET	0x00	SCSI chip reset by reset command or power-up.
ST_SELECT	0x11	Selection of target complete (after C93SELATN).
ST_SATOK	0x16	Select-And-Transfer completed successfully, that is, all phases have completed in a normal manner.
ST_TR_DATAOUT	0x18	Transfer command done, target requesting data.
ST_TR_DATAIN	0x19	Transfer command done, target sending data.
ST_TR_STATIN	0x1b	Target is sending status in.
ST_TR_MSGIN	0x1f	Transfer command done, target sending message.
ST_TRANPAUSE	0x20	Transfer command has paused with ACK.
ST_SAVEDP	0x21	Save Data Pointers message during SAT normal state when device is disconnecting from the bus.
ST_A_RESELECT	0x27	Reselected after disconnect (93A).
ST_UNEXPDISC	0x41	Device disconnected without sending a disconnect message. This sometimes happens when devices with removable media have had the media removed during a transfer.
ST_PARITY	0x43	Command terminated due to parity error on the SCSI bus.

**Table 13-12 (continued)** SCSI State Error Messages

State Message	Sense Key	Comments
ST_PARITY_ATN	0x44	Command terminated due to parity error (ATN is asserted so that host can send a message to device; the transfer is just aborted).
ST_TIMEOUT	0x42	Time-out during Select or Reselect, that is, the device never responded to an attempt to select it; normally seen only during hardware inventory probing, but sometimes happens after a SCSI bus reset if device takes a long time to recover from the reset or is powered off.
ST_INCORR_DATA	0x47	Incorrect message or status byte.
ST_UNEX_RDATA	0x48	Unexpected receive data phase device tried to send more data than the SCSI chip is programmed to expect. This can be OK, as when a high-level request is made to transfer more data than the DMA hardware can map on a single request. In this case, simply reprogram the DMA hardware for the next chunk of data and restart the transfer (but don't send a new SCSI command to the device). When printed as part of an error message, it can sometimes be caused by a SCSI cabling problem, or (particularly with devscsi user drivers) by a mismatch in the byte count given to the driver and the byte count implied by the SCSI command sent to the device.
ST_UNEX_SDATA	0x49	Unexpected send-data phase (same as above, but device is asking for more data).
ST_UNEX_CMDPH	0x4a	Unexpected command phase
ST_UNEX_SSTATUS	0x4b	Unexpected send status phases occur at the end of SCSI command (that is, byte count remaining is 0); if they happen at other times, the chip interrupts. This can happen when you ask a device for more data than it can give you, and in this case, you just return a short I/O count to the caller. When printed as part of an error message, it usually implies a cabling or termination problem.
ST_UNEX_RMESGOUT	0x4e	Unexpected request-message-out phase; usually indicates a SCSI cabling problem.

**Table 13-12 (continued)** SCSI State Error Messages

State Message	Sense Key	Comments
ST_UNEX_SMESGIN	0x4f	Unexpected send-message-in phase. Usually indicates a SCSI cabling problem; also happens when device sends an unsolicited disconnect message when preparing to disconnect from the bus.
ST_RESELECT	0x80	WD33C93 has been reselected.
ST_93A_RESEL	0x81	Reselected while idle (93A).
ST_DISCONNECT	0x85	Disconnect has occurred.
ST_NEEDCMD	0x8a	Target is ready for a command.
ST_REQ_SMESGOUT	0x8e	REQ signal for send message out.
ST_REQ_SMESGIN	0x8f	REQ signal for send message in above 3 usually seen only during sync negotiations.

PART SIX

## Network Drivers

### Chapter 14: Network Device Drivers

Network device drivers are special in that they interface a device to the *ifnet* interface of the TCP/IP protocol stack.



---

## Network Device Drivers

A network device driver is a kernel-level driver that connects a communications device to the IRIX TCP/IP protocol stack using the *ifnet* interface established by BSD UNIX. This chapter contains these major topics:

- “Overview of Network Drivers” on page 338 gives an overview of the IRIX networking subsystem and the role of an *ifnet* driver in it.
- “Network Driver Interfaces” on page 340 summarizes the unique interfaces used by an *ifnet* driver.
- “Example *ifnet* Driver” on page 348 displays the code of a network driver, omitting all device-specific features.

**Note:** If your interest is in creating a network application based on sockets, TLI, or streams, this chapter offers little but background information. Refer to the *IRIX Network Programming Guide*, document Number 007-0810-050, for a complete review of all application-level services.

Even if your interest is in creating a kernel-level network driver, you should be familiar with the facilities documented in the *IRIX Network Programming Guide*. This chapter assumes that you are familiar with them.

## Overview of Network Drivers

A network driver is a kernel-level driver module that connects a communications device such as an Ethernet board to the IRIX implementation of TCP/IP. An overview of the IRIX networking subsystem is shown in Figure 14-1.

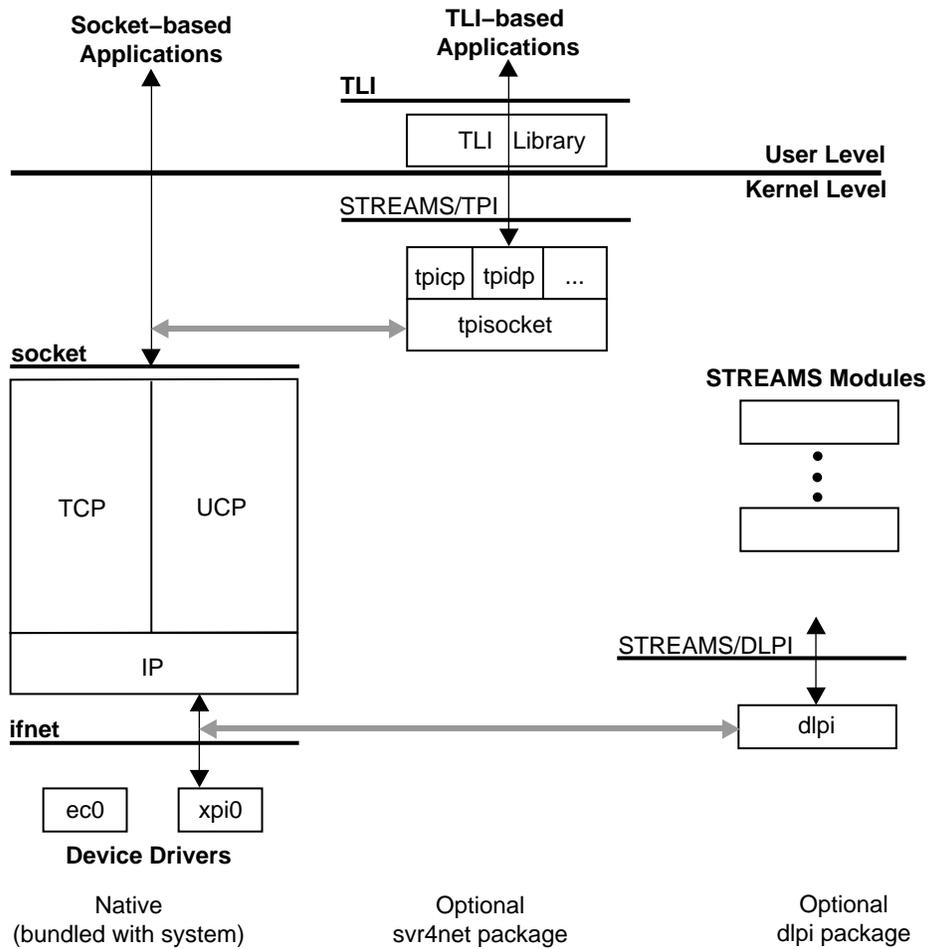


Figure 14-1 Overview of Network Architecture

## Application Interfaces

User-level processes access the network in one of three ways:

- using the BSD socket interface (top left of Figure 14-1)
- using the SVR4 TLI interface through compatibility libraries that convert TLI operations into socket operations (top center of Figure 14-1)
- using a STREAMS interface to a STREAMS-based protocol stack (top right of Figure 14-1)

These three interfaces are documented in the *IRIX Network Programming Guide*

The native socket-based TCP/IP protocol code, the socket layer, and a number of *ifnet*-based device drivers are bundled in the basic IRIX system. Socket-based applications such as *rlogin*, *rcp*, NFS client and server, and the socket-based RPC library operate directly over this native networking framework.

Compatibility support is included for applications written to the STREAMS Transport Layer Interface (TLI). *tpisocket* is a kernel library module used by protocol-specific STREAMS pseudo-drivers, such as *tpitcp*, *tpiudp*, and so on, providing a TPI interface above the native kernel sockets-based network protocol stack.

A STREAMS pseudo-driver that supports the Data Link Provider Interface (DLPI) for STREAMS-based kernel protocol stacks is delivered in the optional *dlpi* package.

## Protocol Stack Interfaces

A *protocol stack* is the software subsystem that manages data traffic according to the rules of a particular communications protocol. There are two ways in which a protocol stack can be integrated into the IRIX kernel. The TCP/IP stack creates and uses the *ifnet* interface to drivers (bottom left of Figure 14-1) and the socket interface to applications (top left of Figure 14-1).

Alternatively, a stack written to the DLPI architecture can communicate with STREAMS drivers (bottom right of Figure 14-1).

## Device Driver Interfaces

A network driver uses the methods and facilities of other kernel-level device drivers, as described in Part III, “Kernel-Level Drivers” of this book. A network driver is compiled and linked like other drivers, configured using the same configuration files, and loaded into the kernel by *lboot* like other drivers.

However, other device drivers support the UNIX filesystem, transferring data in response to calls to their *pxread()*, *pxwrite()*, or *pxstrategy()* entry points. This is not the case with a network driver; it supports protocol stacks, and it transfers data in response to calls from the *ifnet* interface.

## Network Driver Interfaces

The IRIX kernel networking design is based on the kernel networking framework in 4.3BSD. If you are familiar with the 4.3BSD kernel networking design, then you are already familiar with the IRIX kernel networking design because they are basically the same.

The IRIX networking design is based on the socket interface: *mbuf* objects are used to exchange messages within the kernel, and device drivers support the TCP/IP internet protocol suite by supporting the *ifnet* interface.

Since the BSD-based networking framework and the implementation of the TCP/IP protocol suite have changed little from previous releases of IRIX, porting your *ifnet* device driver to this release of IRIX should be straightforward.

**Note:** Although the general kernel facilities documented in Chapter 9, “Device Driver/Kernel Interface,” are standardized and stable, this is not the case with network interfaces. *The ifnet and other interfaces summarized in this topic are subject to change without notice.*

## Kernel Facilities

A network driver is structured like any kernel-level device driver, much as described in Chapter 8, “Structure of a Kernel-Level Driver,” but with the following similarities and differences:

- A network driver is loaded by *lboot* in response to either a USE or VECTOR line in a file in */var/sysgen/system* (see “Configuring a Nonloadable Driver” on page 234).
- A network driver is initialized by a call to its *pfxinit()* entry point when it is loaded, and to its *pfxattach()* entry point when its corresponding physical device is detected.
- A network driver does not need to provide any other entry points (see “Entry Point Summary” on page 142).
- A network driver does not need to provide a driver flag constant *pfxdevflag* because a network driver is always assumed to be multiprocessor-aware (see “Driver Flag Constant” on page 145).
- Although a network driver can use the kernel functions for synchronization and locking (see “Waiting and Mutual Exclusion” on page 206), it normally does not because the *ifnet* interface includes special-purpose locking facilities that are more convenient (see “Multiprocessor Considerations” on page 344).

## Principal ifnet Header Files

The software interface to network facilities is declared in the following important header files:

<i>net/if.h</i>	Basic ifnet facilities and data structures, including the <i>ifnet</i> structure, the basic driver interface object.
<i>net/if_types.h</i>	Constants for interface types, used in decoding address headers.
<i>sys/mbuf.h</i>	The <i>mbuf</i> structure with related constants and macros, and declarations of functions to allocate, manipulate, and free <i>mbuf</i> objects.
<i>net/netisr.h</i>	Declarations related to software interrupts, including <b>schednetisr()</b> to schedule an interrupt, and the IP input queue <i>ipintrq</i> .
<i>net/multi.h</i>	Routines defining a generic filter for use by drivers whose devices cannot perfectly filter multicast packets.

<i>net/soioctl.h</i>	Socket <b>ioctl()</b> function numbers, some of which reach a driver for action.
<i>net/raw.h</i>	The interface to the raw protocol family members <i>snoop</i> and <i>drain</i> .
<i>net/if_arp.h</i>	Generic ARP declarations.
<i>netinet/if_ether.h</i>	Essential declarations for Ethernet drivers, including ARP protocol for Ethernet.
<i>sys/dlsap_register.h</i>	DLPI interface declarations.

### Debugging Facilities

When your driver is operating under a debugging kernel, you can use the facilities of *symmon* and *idbg* to display a variety of network-related data structures. See “Preparing the System for Debugging” on page 245, and see “Commands to Display Network-Related Structures” on page 271.

### Information Sources

Aside from comments in header files, the complete *ifnet* interface and related interfaces have never been documented. In prior years, most people working on *ifnet* drivers have had access to the Berkeley UNIX source distribution and have been able to answer questions by referring to the code.

Referring to the code is an even more common option today, thanks to the release of 4.4BSD-Lite, a software distribution of BSD UNIX that does not require a source license, now widely available at a reasonable price. To obtain a copy, order the following:

- *4.4BSD-Lite Berkely Software Distribution CD-ROM Companion*, published by USENIX and O’Reilly & Associates; ISBN 1-56592-081-3 (US domestic) or ISBN 1-56592-092-9 (non-US).

The *ifnet* source code in this software is functionally compatible with IRIX *ifnet*, although some protocols (for example, *snoop* and *drain*) are not implemented in BSD-Lite.

Finally, the IRIX reference pages contain a wealth of detail regarding network interfaces. Some reference pages that are related to the interests of driver designers are listed in Table 14-1.

**Table 14-1** Important Reference Pages Related to Network Drivers

Reference Page	Contents
arp(7)	Operation of the ARP protocol, with details of <b>ioctl()</b> functions.
drain(7)	Operation of the drain driver, which receives unwanted packets, with details of its <b>ioctl()</b> functions.
ethernet(7)	Overview of the IRIX Ethernet drivers, including error messages and the use of VECTOR lines to configure them.
fddi(7)	Cursory overview of IRIX FDDI drivers, with naming conventions.
ifconfig(1)	Management program used to enable and disable network interfaces (drivers) and change their runtime parameters.
netintro(7)	Overview of network facilities; mentions the role of the network interface (driver); has extensive detail on routing <b>ioctl()</b> calls.
network(1)	Documents the network initialization script that runs when the system is booted up.
raw(7)	Overview of the Raw protocol family whose members are snoop and drain.
routed(1)	Documents operation of the routing daemon, including <b>ioctl()</b> use.
snoop(7)	Operation of the snoop driver, which allows inspection of packets, with details of its <b>ioctl()</b> features.
ticlts(7)	Operation and use of the ticlts, ticots, and ticotsord loopback drivers.
tokenring(7)	Overview of the IRIX token-ring drivers, including packet formats.

## Network Inventory Entries

The driver must call **add\_to\_inventory()** to add an inventory entry to the inventory database (see *sys/invent.h* and “Creating an Inventory Entry” on page 34):

<i>vhdl</i>	The vertex handle of the attached device.
<i>class</i>	INV_NETWORK
<i>type</i>	The packet type, for example INV_NET_ETHER. See <i>sys/invent.h</i> for the possible “types for class network” list.
<i>controller</i>	The kind of network controller from the “controllers for network types” list in <i>sys/invent.h</i> .
<i>unit</i>	Any distinguishing number for this device. The <i>hinv</i> command does not decode this field.
<i>state</i>	Any characteristic number for this device. The <i>hinv</i> command does not decode this field.

## Multiprocessor Considerations

Prior to IRIX 5.3, the kernel BSD framework code and TCP/IP protocol stack executed under a single kernel lock, creating a single-threaded implementation. Beginning with IRIX 5.3, the BSD framework and TCP/IP protocol suite have been multi-threaded to support symmetric multiprocessing. The code uses different kernel locks to protect different critical sections.

IRIX now supports multiple, concurrent threads of execution within the TCP/UDP/IP protocol suite, and the kernel socket layer. In addition, network device drivers run on any available CPU, concurrently with the network software, applications, and other drivers.

This means that any ifnet-based network driver must be prepared to run asynchronously and concurrently with other drivers and with the protocol stack.

## Ineffective spl() Functions

The **spl\*()** functions were the traditional UNIX method of gaining exclusive use of data. In single-threaded ifnet drivers, the **splimp()** or **splnet()** functions were used to get exclusive use of the ifnet structure.

In a multiprocessor, **spl\*()** functions like **splimp()** or **splnet()** do block interrupts on the local CPU, but they do not prevent interrupts from occurring on other processors in the system, nor do they prevent other processes on other CPUs from executing code that refers to the same data.

If you are porting a driver from a uniprocessor environment, search for any use of an **spl\*()** function and plan to replace it with effective mutual exclusion locking macros.

## Multiprocessor Locking Macros

Under BSD networking, drivers interface with the protocol stacks by queueing incoming packets on a per-protocol input queue. In a multiprocessor, each protocol input queue must be protected by the locking macros defined in the file *net/if.h*.

All the locking macros that protect the input queue are assumed to be called at the proper processor interrupt masking level, **splimp**. All input queue locking macros also take an input parameter *ifq*, which is a pointer to the protocol input queue that must be defined as a *struct ifqueue*.

## Compiler Flags for MP TCP/IP

The `_MP_NETLOCKS` and MP compiler variables must be defined in order to enable the macros necessary to run under multi-threaded TCP/IP (see “Compiler Variables” on page 231).

## Mutual Exclusion Macros

The macros for mutual exclusion defined in *net/if.h* are listed in Table 14-2.

**Table 14-2** Mutual Exclusion Macros for ifnet Drivers

Macro Prototype	Purpose
IFNET_LOCK( <i>ifp</i> , <i>s</i> )	Get exclusive use of the structure <i>*ifp</i> . <b>splimp()</b> is called to raise the interrupt level if necessary, and the returned value is saved in <i>s</i> .
IFNET_UNLOCK( <i>ifp</i> , <i>s</i> )	Release use of <i>*ifp</i> , and return to interrupt level <i>s</i> .
IFNET_LOCKNOSPL( <i>ifp</i> )	Get exclusive use of the structure <i>*ifp</i> , but do not call <b>splimp()</b> (the driver knows it is already at the appropriate level.)
IFNET_UNLOCKNOSPL( <i>ifp</i> )	Release use of <i>*ifp</i> after use of IFNET_LOCKNOSPL.
IFNET_ISLOCKED( <i>ifp</i> )	Test whether <i>*ifp</i> is locked.
IFQ_LOCK( <i>ifq</i> )	Get exclusive use of an input queue <i>*ifq</i> .
IFQ_UNLOCK( <i>ifq</i> )	Release use of <i>*ifq</i> .
IF_ENQUEUE( <i>ifq</i> , <i>mp</i> )	Lock the queue <i>*ifq</i> ; post the mbuf <i>*mp</i> ; release the queue.
IF_ENQUEUE_NOLOCK( <i>ifq</i> , <i>mp</i> )	Post the mbuf <i>*mp</i> without locking.

The variables used in Table 14-2 are as follows:

<i>ifp</i>	Address of a <i>struct ifnet</i> to be used exclusively.
<i>s</i>	Integer variable to store the current interrupt mask level.
<i>ifq</i>	Address of a <i>struct ifqueue</i> to be posted.
<i>mp</i>	Address of a <i>struct mbuf</i> to be posted.

**Macro Use**

The TCP/IP protocol stack automatically acquires the ifnet structure before calling a network driver routine through that structure. Thus the driver's **init()**, **stop()**, **start()**, **output()**, and **ioctl()** functions do not need to use IFNET\_LOCK or IFNET\_UNLOCK. Look for expressions

```
ASSERT( IFNET_ISLOCKED( ifp ) );
```

in the example driver ("Example ifnet Driver" on page 348) to see places where this is the case. Explicit use of IFNET\_LOCK is needed in the interrupt handler.

## Example ifnet Driver

The code in Example 14-1 represents the skeleton of an ifnet driver, showing its entry points, data structures, required `ioctl()` functions, address format conventions, and its use of kernel utility routines and locking primitives.

A comment beginning "MISSING:" represents a point at which a complete driver would contain code related to the device or bus it manages.

### Example 14-1 Skeleton ifnet Driver

```
/*
 * if_sk - skeleton IRIX 6.4 ifnet device driver
 *
 * This is a skeleton ifnet driver for IRIX 6.4
 * meant to demonstrate ifnet driver entry points,
 * data structures, required ioctls, address format
 * conventions, kernel utility routines, and locking
 * primitives. These kernel data structures and
 * routines are SUBJECT TO CHANGE w/o notice.
 *
 * Refer to the IRIX 6.4 Device Driver Programming Guide
 * and Device Driver Reference Pages for complete
 * information on writing PCI, GIO, VME, and EISA bus
 * device drivers for SGI systems.
 *
 * "MISSING" is used to designate places where device/bus/driver-
 * specific code sections are required.
 *
 * Locking strategy:
 * IFNET_LOCK() and IFNET_UNLOCK() acquire/release the
 * lock on a given ifnet structure. IFQ_LOCK() and
 * IFQ_UNLOCK() acquire/release the lock on a given ifqueue
 * structure. The ifnet or ifqueue lock must be held while
 * modifying any fields within the associated data
 * structure. The ifnet lock is also held to singlethread
 * portions of the device driver. The driver xxinit,
 * xxreset, xxoutput, xxwatchdog, and xxioctl entry points
 * are called with IFNET_LOCK() already acquired thus only
 * a single thread of execution is allowed in these
 * portions of the driver for each interface. It is the
 * driver's responsibility to call IFNET_LOCK() within its
 * xxintr() and other private routines to singlethread any
 * other critical sections. It is also the driver's
 * responsibility to acquire the ifq lock by calling
```

```
* IFQ_LOCK() before attempting to enqueue onto the IP
* input queue "ipintrq".
*
* Notes:
* - don't forget appropriate machine-specific cache flushing operations
*   (refer to IRIX Device Driver Programming guide)
* - declare pointers to device registers as "volatile"
*
* Caveat Emptor:
* No guarantees are made wrt correctness nor completeness
* of this source.
*
* Copyright 1996 Silicon Graphics, Inc. All rights reserved.
*/
#ident "$Revision: 2.0$"
#include <sys/types.h>
#include <sys/param.h>
#include <sys/system.h>
#include <sys/sysmacros.h>
#include <sys/cmn_err.h>
#include <sys/debug.h>
#include <sys/hwgraph.h>
#include <sys/iograph.h>
#include <sys/errno.h>
#include <sys/PCI/pciio.h>
#include <sys/idbgentry.h>
#include <sys/tcp-param.h>
#include <sys/mbuf.h>
#include <sys/immu.h>
#include <sys/sbd.h>
#include <sys/ddi.h>
#include <sys/kmem.h>
#include <sys/cpu.h>
#include <sys/invent.h>
#include <net/if.h>
#include <net/if_types.h>
#include <net/netisr.h>
#include <netinet/if_ether.h>
#include <net/raw.h>
#include <net/multi.h>
#include <netinet/in_var.h>
#include <net/soioctl.h>
#include <sys/dlsap_register.h>
/* MISSING: driver-specific header includes go here */
/*
```

```
* driver-specific and device-specific data structure
* declarations and definitions might go here.
*/
#define SK_MAX_UNITS    8
#define SK_MTU          4096
#define SK_DOG          (2*IFNET_SLOWHZ) /* watchdog duration in seconds */
#define SK_IPT          (IPT_FDDI) /* refer to <net/if_types.h> */
#define SK_INV          (INV_NET_FDDI) /* refer to <sys/invent.h> */
#define INV_FDDI_SK     (23) /* refer to <sys/invent.h> */
#define IFF_ALIVE       (IFF_UP|IFF_RUNNING)
#define iff_alive(flags) (((flags) & IFF_ALIVE) == IFF_ALIVE)
#define iff_dead(flags) (((flags) & IFF_ALIVE) != IFF_ALIVE)
#define SK_ISBROAD(addr) (!bcmp((addr), &skbroadcastaddr, SKADDRLEN))
#define SK_ISGROUP(addr) ((addr)[0] & 01)
/*
 * MISSING media-specific definitions of address size and header format.
 */
#define SKADDRLEN      (6)
#define SKHEADERLEN    (sizeof (struct skheader))
/*
 * Our fictional media has an IEEE 802-looking header..
 */
struct skaddr {
    u_int8_t sk_vec[SKADDRLEN];
};
struct skheader {
    struct skaddr sh_dhost;
    struct skaddr sh_shost;
    u_int16_t sh_type;
};
struct skaddr skbroadcastaddr = {
    0xff, 0xff, 0xff, 0xff, 0xff, 0xff
};
/*
 * Each interface is represented by a private
 * network interface data structure that maintains
 * the device hardware resource addresses, pointers
 * to device registers, allocated dma_alloc maps,
 * lists of mbufs pending transmit or reception, etc, etc.
 * We use ARP and have an 802 address.
 */
struct sk_info {
    struct arpcom si_ac; /* common ifnet and arp */
    struct skaddr si_ouraddr; /* our individual media address */
    struct mfilter si_filter; /* AF_RAW sw snoop filter */
};
```

```

    struct rawif si_rawif;      /* raw snoop interface */
    int si_flags;
    caddr_t si_regs;          /* pointer to device registers */
    vertex_hdl_t si_our_vhdl;  /* our vertex */
    vertex_hdl_t si_conn_vhdl; /* our parent vertex */
    pciio_intr_t si_intr;     /* interrupt handle */
    /* MISSING additional driver-specific data structures */
};
#define si_if si_ac.ac_if
#define sktoifp(si) (&(si)->si_ac.ac_if)
#define ifptosk(ifp)((struct sk_info *)ifp)
#define ALIGNED(addr, alignment) (((u_long)(addr) & (alignment-1)) == 0)
#define sk_info_set(v,i) hwgraph_fastinfo_set((v),(arbitrary_info_t)(i))
#define sk_info_get(v) ((struct sk_info *)hwgraph_fastinfo_get((v)))
/*
 * The start of an mbuf containing an input frame
 */
struct sk_ibuf {
    struct ifheader sib_ifh;
    struct snoopheader sib_snoop;
    struct skheader sib_skh;
};
#define SK_IBUFSZ (sizeof (struct sk_ibuf))
/*
 * Multicast filter request for SIOCADMULTI/SIOCDELMULTI .
 */
struct mfreq {
    union mkey *mfr_key; /* pointer to socket ioctl arg */
    mval_t mfr_value; /* associated value */
};
void sk_init(void);
static int sk_ifinit(struct ifnet *ifp);
int sk_attach(vertex_hdl_t conn_vhdl);
static void sk_reset(struct sk_info *si);
static void sk_intr(struct sk_info *si);
static int sk_output(struct ifnet *ifp, struct mbuf *m, struct sockaddr *dst);
static void sk_input(struct sk_info *si, struct mbuf *m, int totlen);
static int sk_ioctl(struct ifnet *ifp, int cmd, void *data);
static void sk_watchdog(struct ifnet *ifp);
static void sk_stop(struct sk_info *si);
static int sk_start(struct sk_info *si, int flags);
static int sk_add_da(struct sk_info *si, union mkey *key, int ismulti);
static int sk_del_da(struct sk_info *si, union mkey *key, int ismulti);
static int sk_dahash(char *addr);
static int sk_dlp(struct sk_info *si, int port, int encap, struct mbuf *m, int len);

```

```
static void sk_dump(int unit);
/* MISSING additional driver-specific routine prototypes */
extern struct ifqueue ipintrq; /* ip input queue */
extern struct ifnet loif; /* loopback driver if */
extern void bitswapcopy(void *, void *, int);
int sk_devflag = D_MP;
/*
 * xxinit() routine called early during boot.
 */
void
sk_init(void)
{
    /* register ourselves with the pci i/o infrastructure */
    pciio_driver_register(0x10A9, 0x0003, "sk_", 0);
    /*
     * register a handy debugging routine so we can call it
     * from idbg(1) and the kernel debugger.
     */
    idbg_addfunc("sk_dump", (void (*)())sk_dump);
}
/*
 * xxattach() routine is called by the i/o infrastructure
 * when a hardware device matches our pci vendor and device ids.
 */
int
sk_attach(vertex_hdl_t conn_vhdl)
{
    graph_error_t rc;
    vertex_hdl_t our_vhdl;
    struct sk_info *si;
    struct ifnet *ifp;
    device_desc_t sk_dev_desc;
    /* add a char device vertex to the hardware graph tree ("/hw") */
    if ((rc = hwgraph_char_device_add(conn_vhdl, "sk", "sk_",
        &our_vhdl)) != GRAPH_SUCCESS) {
        cmn_err(CE_ALERT, "skattach: hwgraph_char_device_add error %d", rc);
        return (EIO);
    }
    /* fix up device descriptor */
    sk_dev_desc = device_desc_dup(our_vhdl);
    device_desc_intr_name_set(sk_dev_desc, "sk device");
    device_desc_default_set(our_vhdl, sk_dev_desc);
    if ((si = (struct sk_info*) kmem_zalloc(sizeof (struct sk_info), KM_SLEEP)) == NULL) {
        cmn_err(CE_ALERT, "skattach: kmem_alloc failed\n");
        return (ENOMEM);
    }
}
```

```

}
/* save our vertex and our parent's vertex for later */
si->si_our_vhdl = our_vhdl;
si->si_conn_vhdl = conn_vhdl;
/* save a pointer to our sk_info structure in our vertex */
sk_info_set(our_vhdl, si);
/*
 * MISSING
 * Driver-specific actions that might go here:
 *
 * - call sk_reset to disable the device
 * - pciio_pio map in the device registers
 * - allocate a new sk_info structure
 * - allocate device host memory buffers and descriptors
 * and create any static dma mappings (pciio_dmamap_xx )
 * ...
 */
/* register our interrupt handler */
si->si_intr = pciio_intr_alloc(conn_vhdl, sk_dev_desc,
    PCIIO_INTR_LINE_A, our_vhdl);
pciio_intr_connect(si->si_intr, (intr_func_t) sk_intr, (intr_arg_t) si, (void *)0);
/*
 * MISSING your address translation protocol goes here.
 * Save a copy of our MAC address in the arpcom structure.
 */
bcopy((caddr_t)&si->si_ouraddr, (caddr_t)si->si_ac.ac_enaddr,
    SKADDRLEN);
/*
 * Initialize ifnet structure with our name, type, mtu size,
 * supported flags, pointers to our entry points,
 * and attach to the available ifnet drivers list.
 */
ifp = sktoifp(si);
ifp->if_name = "sk";
ifp->if_unit = -1;
ifp->if_type = SK_IPT;
ifp->if_mtu = SK_MTU;
ifp->if_flags = IFF_BROADCAST | IFF_MULTICAST | IFF_NOTRAILERS;
ifp->if_output = sk_output;
ifp->if_ioctl = (int (*)(struct ifnet*, int, void*))sk_ioctl;
ifp->if_watchdog = sk_watchdog;
/*
 * A note about unit numbering and when to call if_attach:
 *
 * IRIX 6.4 includes a new boot-time command ioconfig(1M)

```

```
* which walks the hardware device tree ("/hw") and allocates
* and assigns a controller number (unit number) to each
* device vertex it finds which has an inventory record.
*
* So we do everything but the if_attach() call now,
* since we don't yet have our unit number, and call
* if_attach() from our xxopen() routine when it is
* called by the ioconfig(1M) command during booting.
*/
/*
* Allocate a multicast filter table with an initial
* size of 10. See <net/multi.h> for a description
* of the support for generic sw multicast filtering.
* Use of these mf routines is purely optional -
* if you're not supporting multicast addresses or
* your device does perfect filtering or you think
* you can roll your own better, feel free.
*/
if (!mfnew(&si->si_filter, 10))
    cmn_err(CE_PANIC, "sk_edtinit: no memory for frame filter\n");
/*
* You must create an inventory record for this vertex now
* or ioconfig(1M) will not call our xxopen() routine to
* pass in an allocated unit number later.
*/
device_inventory_add(our_vhdl, INV_NETWORK, INV_NET_FDDI, 100, -1, 0);
return (0);
}
/*
* Driver xxopen() routine exists only to take unit# which
* has now been assigned to the vertex by ioconfig(1M)
* and if_attach() the device.
*/
/* ARGSUSED */
int
sk_open(dev_t *devp, int flag, int otyp, struct cred *crp)
{
    vertex_hdl_t our_vhdl;
    struct sk_info *si;
    int unit;
    our_vhdl = dev_to_vhdl(*devp);
    if ((si = sk_info_get(our_vhdl)) == NULL)
        return (EIO);
    /* if already if_attached, just return */
    if (si->si_if.if_unit != -1)
```

```
        return (0);
/* get our unit number from the vertex label */
if ((unit = device_controller_num_get(our_vhdl)) < 0) {
    cmn_err(CE_ALERT, "sk_open: vertex missing ctrlr number");
    return (EIO);
}
si->si_if.if_unit = unit;
/*
 * Install this device in the list of IRIX ifnet structures.
 */
if_attach(&si->si_if);
/*
 * Initialize the raw socket interface. See <net/raw.h>
 * and the man pages for descriptions of the SNOOP
 * and DRAIN raw protocols.
 */
rawif_attach(&si->si_rawif, &si->si_if,
             (caddr_t) &si->si_ouraddr,
             (caddr_t) &skbroadcastaddr,
             SKADDRLEN,
             SKHEADERLEN,
             structoff(skheader, sh_shost),
             structoff(skheader, sh_dhost));
return (0);
}
static int
sk_ifinit(struct ifnet *ifp)
{
    struct sk_info *si;
    si = ifptosk(ifp);
    ASSERT(IFNET_ISLOCKED(ifp));
    /*
     * Reset the device first, ask questions later..
     */
    sk_reset(si);
    /*
     * - free or reuse any pending xmit/recv mbufs
     * - initialize device configuration registers, etc.
     * - allocate and post receive buffers
     *
     * Refer to Device Driver Programming guide for
     * descriptions on use of kvtophys() (GIO) or
     * dma_map/dma_mapaddr() (VME) routines for
     * obtaining DMA addresses and system-specific
     * issues like flushing caches or write buffers.
     */
}
```

```
    */
/*
 * MISSING
 * enable if_flags device behavior (IFF_DEBUG on/off, etc.)
 */
ifp->if_timer = SK_DOG; /* turn on watchdog */
/* MISSING: turn device "on" now */
return 0;
}
/*
 * Reset the interface.
 */
static void
sk_reset(struct sk_info *si)
{
    struct ifnet *ifp = sktoifp(si);
    ifp->if_timer = 0; /* turn off watchdog */
    /*
     * MISSING
     * - reset device
     * - reset device receive descriptor ring
     * - free any enqueued transmit mbufs
     * - create device xmit descriptor ring
     */
}

static void
sk_intr(struct sk_info *si)
{
    struct ifnet *ifp;
    struct mbuf *m;
    int totlen;
#ifdef INTR_KTHREADS
    int s;
#endif
    ifp = &si->si_if;
    /*
     * Ignore early interrupts.
     */
    if (iff_dead(ifp->if_flags)) {
        sk_stop(si);
        return;
    }
    IFNET_LOCK(ifp, s); /* acquire interface lock */
    /*
```

```

    * MISSING: read and clear the device interrupt status register.
    */
/*
    * process any received packets.
    */
while (0/* MISSING: received packets available */) {
    /*
        * MISSING
        * Do device-specific receive processing here.
        * Allocate and post a replacement receive buffer.
        */
    sk_input(si, m, totlen);
}
while (0/* MISSING mbufs completed transmission */) {
    /*
        * MISSING
        * Reclaim any completed device transmit resources
        * freeing completed mbufs, checking for errors,
        * and maintaining if_opackets, if_oerrors,
        * if_collisions, etc.
        */
}
/* MISSING: process any other interrupt conditions */
IFNET_UNLOCK(ifp, s);
}
/*
    * Transmit packet.  If the destination is this system or
    * broadcast, send the packet to the loop-back device if
    * we cannot hear ourself transmit.  Return 0 or errno.
    */
static int
sk_output(
    struct ifnet    *ifp,
    struct mbuf    *m0,
    struct sockaddr *dst)
{
    struct sk_info *si = ifptosk(ifp);
    struct skheader *sh;
    struct mbuf *m, *m1;
    struct mbuf *mloop;
    struct sockaddr_sdl *sdl;
    int error;
    ASSERT(IFNET_ISLOCKED(ifp));
    mloop = NULL;
    if (iff_dead(ifp->if_flags)) {

```

```
        error = EHOSTDOWN;
        goto bad;
    }
    /*
     * If snd queue full, try reclaiming some completed
     * mbufs.  If it's still full, then just drop the
     * packet and return ENOBUFS.
     */
    if (IF_QFULL(&si->si_if.if_snd)) {
        while (0/* MISSING xmits done */) {
            /*
             * MISSING: Reclaim completed xmit descriptors.
             */
            IF_DEQUEUE_NOLOCK(&si->si_if.if_snd, m);
            m_freem(m);
        }
        if (IF_QFULL(&si->si_if.if_snd)) {
            m_freem(m0);
            si->si_if.if_odrops++;
            IF_DROP(&si->si_if.if_snd);
            return (ENOBUFS);
        }
    }
}
switch (dst->sa_family) {
case AF_INET: {
    /*
     * Get room for media header,
     * use this mbuf if possible.
     */
    if (!M_HASCL(m0)
        && m0->m_off >= MMINOFF+sizeof(*sh)
        && (sh = mtod(m0, struct skheader*))
        && ALIGNED(sh, sizeof (int))) {
        ASSERT(m0->m_off <= MSIZE);
        m1 = 0;
        --sh;
    } else {
        m1 = m_get(M_DONTWAIT, MT_DATA);
        if (m1 == NULL) {
            m_freem(m0);
            si->si_if.if_odrops++;
            IF_DROP(&si->si_if.if_snd);
            return (ENOBUFS);
        }
        sh = mtod(m1, struct skheader*);
    }
}
```

```

        m1->m_len = sizeof (*sh);
    }
    bcopy(&si->si_ouraddr, &sh->sh_shost, SKADDRLLEN);
    /*
     * translate dst IP address to media address.
     */
    if (!ip_arpresolve(&si->si_ac, m0,
        &((struct sockaddr_in *)dst)->sin_addr,
        (u_char*)&sh->sh_dhost)) {
        m_freem(m1);
        return (0); /* just wait if not yet resolved */
    }
    if (m1 == 0) {
        m0->m_off -= sizeof (*sh);
        m0->m_len += sizeof (*sh);
    } else {
        m1->m_next = m0;
        m0 = m1;
    }
    /*
     * Listen to ourselves, if we are supposed to.
     */
    if (SK_ISBROAD(&sh->sh_shost)) {
        mloop = m_copy(m0, sizeof (*sh), M_COPYALL);
        if (mloop == NULL) {
            m_freem(m0);
            si->si_if.if_odrops++;
            IF_DROP(&si->si_if.if_snd);
            return (ENOBUFS);
        }
    }
    break;
}
case AF_UNSPEC:
#define EP ((struct ether_header *)&dst->sa_data[0])
    /*
     * Translate an ARP packet using RFC-1042.
     * Require the entire ARP packet be in the first mbuf.
     */
    sh = mtod(m0, struct skheader*);
    if (M_HASCL(m0)
        || !ALIGNED(sh, sizeof (int))
        || m0->m_len < sizeof(struct ether_arp)
        || m0->m_off < MMINOFF+sizeof(*sh)
        || EP->ether_type != ETHERTYPE_ARP) {

```

```
        printf("sk_output: bad ARP output\n");
        m_freem(m0);
        si->si_if.if_oerrors++;
        IF_DROP(&si->si_if.if_snd);
        return (EAFNOSUPPORT);
    }
    ASSERT(m0->m_off <= MSIZE);
    m0->m_len += sizeof(*sh);
    m0->m_off -= sizeof(*sh);
    --sh;
    bcopy(&si->si_ouraddr, &sh->sh_shost, SKADDRLEN);
    bcopy(&EP->ether_dhost[0], &sh->sh_dhost, SKADDRLEN);
    sh->sh_type = EP->ether_type;
# undef EP
    break;
case AF_RAW:
    /* The mbuf chain contains the raw frame incl header.
    */
    sh = mtod(m0, struct skheader*);
    if (M_HASCL(m0)
        || m0->m_len < sizeof(*sh)
        || !ALIGNED(sh, sizeof(int))) {
        m0 = m_pullup(m0, SKHEADERLEN);
        if (m0 == NULL) {
            si->si_if.if_odrops++;
            IF_DROP(&si->si_if.if_snd);
            return (ENOBUFS);
        };
        sh = mtod(m0, struct skheader*);
    }
    break;
case AF_SDL:
    /*
    * Send an 802 packet for DLPI.
    * mbuf chain should already have everything
    * but MAC header.
    */
    sdl = (struct sockaddr_sdl*) dst;
    /* sanity check the MAC address */
    if (sdl->ssdl_addr_len != SKADDRLEN) {
        m_freem(m0);
        return (EAFNOSUPPORT);
    }
    sh = mtod(m0, struct skheader*);
    if (!M_HASCL(m0)
```

```

    && m1->m_off >= MMINOFF+SKHEADERLEN
    && ALIGNED(sh, sizeof(int))) {
    ASSERT(m0->m_off <= MSIZE);
    m0->m_len += SKHEADERLEN;
    m0->m_off -= SKHEADERLEN;
} else {
    m1 = m_get(M_DONTWAIT, MT_DATA);
    if (!m1) {
        m_freem(m0);
        si->si_if.if_odrops++;
        IF_DROP(&si->si_if.if_snd);
        return (ENOBUFS);
    }
    m1->m_len = SKHEADERLEN;
    m1->m_next = m0;
    m0 = m1;
    sh = mtod(m0, struct skheader*);
}
sh->sh_type = htons(ETHERTYPE_IP);
bcopy(&si->si_ouraddr, &sh->sh_shost, SKADDRLEN);
bcopy(sd1->ssdl_addr, &sh->sh_dhost, SKADDRLEN);
break;
default:
    printf("sk_output: bad af %u\n", dst->sa_family);
    m_freem(m0);
    return (EAFNOSUPPORT);
}
/*
 * Check whether snoopers want to copy this packet.
 */
if (RAWIF_SNOOPING(&si->si_rawif)
    && snoop_match(&si->si_rawif, (caddr_t)sh, m0->m_len)) {
    struct mbuf *ms, *mt;
    int len;          /* m0 bytes to copy */
    int lenoff;
    int curlen;
    len = m_length(m0);
    lenoff = 0;
    curlen = len + SK_IBUFSZ;
    if (curlen > MCLBYTES)
        curlen = MCLBYTES;
    ms = m_vget(M_DONTWAIT, MAX(curlen, SK_IBUFSZ), MT_DATA);
    if (ms) {
        IF_INITHEADER(mtod(ms, caddr_t), &si->si_if, SK_IBUFSZ);
        curlen = m_datacopy(m0, lenoff, curlen - SK_IBUFSZ,

```

```
        mtod(ms,caddr_t) + SK_IBUFSZ);
mt = ms;
for (;;) {
    lenoff += curlen;
    len -= curlen;
    if (len <= 0)
        break;
    curlen = MIN(len, MCLBYTES);
    m1 = m_vget(M_DONTWAIT, curlen, MT_DATA);
    if (0 == m1) {
        m_freem(ms);
        ms = 0;
        break;
    }
    mt->m_next = m1;
    mt = m1;
    curlen = m_datacopy(m0, lenoff, curlen,
        mtod(m1, caddr_t));
}
}
if (ms == NULL) {
    snoop_drop(&si->si_rawif, SN_PROMISC,
        mtod(m0,caddr_t), m0->m_len);
} else {
    (void)snoop_input(&si->si_rawif, SN_PROMISC,
        mtod(m0, caddr_t),
        ms,
        (lenoff > SKHEADERLEN)?
        (lenoff - SKHEADERLEN) : 0);
}
}
/*
 * Save a copy of the mbuf chain to free later.
 */
IF_ENQUEUE_NOLOCK(&si->si_if.if_snd, m0);
/*
 * MISSING
 * Alloc and initialize transmit descriptor resources
 * and kick the chip to start DMA reads for transmitting.
 */
if (error)
    goto bad;
ifp->if_opackets++;
if (mloop) {
    si->si_if.if_omcasts++;
}
```

```

        (void) looutput(&loif, mloop, dst);
    } else if (SK_ISGROUP(sh->sh_dhost.sk_vec))
        si->si_if.if_omcasts++;
    return (0);
bad:
    ifp->if_oerrors++;
    m_freem(m);
    m_freem(mloop);
    return (error);
}
/*
 * deal with a complete input frame in a string of mbufs.
 * mbuf points at a (struct sk_ibuf), totlen is #bytes
 * in user data portion of the mbuf.
 */
static void
sk_input(struct sk_info *si,
        struct mbuf *m,
        int totlen)
{
    struct sk_ibuf *sib;
    struct ifqueue *ifq;
    int snoopflags = 0;
    uint port;
    /*
     * MISSING: set local variables 'snoopflags' and 'if_ierrors' as appropriate
     */
    ifq = NULL;
    sib = mtod(m, struct sk_ibuf*);
    IF_INITHEADER(sib, &si->si_if, SK_IBUFSZ);
    si->si_if.if_abytes += totlen;
    si->si_if.if_ipackets++;
    /*
     * If it is a broadcast or multicast frame,
     * get rid of imperfectly filtered multicasts.
     */
    if (SK_ISGROUP(sib->sib_skh.sh_dhost.sk_vec)) {
        if (SK_ISBROAD(sib->sib_skh.sh_dhost.sk_vec))
            m->m_flags |= M_BCAST;
        else {
            if (((si->si_ac.ac_if.if_flags & IFF_ALLMULTI) == 0)
                && !mfethermatch(&si->si_filter,
                    sib->sib_skh.sh_dhost.sk_vec, 0)) {
                if (RAWIF_SNOOPING(&si->si_rawif)
                    && snoop_match(&si->si_rawif,

```

```
        (caddr_t) &sib->sib_skh, totlen))
        snoopflags = SN_PROMISC;
    else {
        m_freem(m);
        return;
    }
    m->m_flags |= M_MCAST;
}
}
si->si_if.if_imcasts++;
} else {
    if (RAWIF_SNOOPING(&si->si_rawif)
        && snoop_match(&si->si_rawif,
            (caddr_t) &sib->sib_skh,
            totlen))
        snoopflags = SN_PROMISC;
    else {
        m_freem(m);
        return;
    }
}
/*
 * Set 'port' . For us, just sh_type.
 */
port = ntohs(sib->sib_skh.sh_type);
/*
 * do raw snooping.
 */
if (RAWIF_SNOOPING(&si->si_rawif)) {
    if (!snoop_input(&si->si_rawif, snoopflags,
        (caddr_t)&sib->sib_skh,
        m,
        (totlen>sizeof(struct skheader)
        ? totlen-sizeof(struct skheader) : 0))) {
    }
    if (snoopflags)
        return;
} else if (snoopflags) {
    goto drop; /* if bad, count and skip it */
}
/*
 * If it is a frame we understand, then give it to the
 * correct protocol code.
 */
switch (port) {
```

```

case ETHERTYPE_IP:
    ifq = &ipintrq;
    break;
case ETHERTYPE_ARP:
    arpinput(&si->si_ac, m);
    return;
default:
    if (sk_dlp(si, port, DL_ETHER_ENCAP, m, totlen))
        return;
    break;
}
/*
 * if we cannot find a protocol queue, then flush it down the
 * drain, if it is open.
 */
if (ifq == NULL) {
    if (RAWIF_DRAINING(&si->si_rawif)) {
        drain_input(&si->si_rawif,
                    port,
                    (caddr_t)&sib->sib_skh.sh_dhost.sk_vec,
                    m);
    } else
        m_freem(m);
    return;
}
/*
 * Put it on the IP protocol queue.
 */
if (IF_QFULL(ifq)) {
    si->si_if.if_iqdrops++;
    si->si_if.if_ierrors++;
    IF_DROP(ifq);
    goto drop;
}
IF_ENQUEUE(ifq, m);
schednetisr(NETISR_IP);
return;
drop:
m_freem(m);
if (RAWIF_SNOOPING(&si->si_rawif))
    snoop_drop(&si->si_rawif, snoopflags,
               (caddr_t)&sib->sib_skh, totlen);
if (RAWIF_DRAINING(&si->si_rawif))
    drain_drop(&si->si_rawif, port);
}

```

```
/*
 * See if a DLPI function wants a frame.
 */
static int
sk_dlp(struct sk_info *si,
       int port,
       int encap,
       struct mbuf *m,
       int len)
{
    dlsap_family_t *dlp;
    struct mbuf *m2;
    struct sk_ibuf *sib;
    if ((dlp = dlsap_find(port, encap)) == NULL)
        return (0);
    /*
     * The DLPI code wants the entire MAC and LLC headers.
     * It needs the total length of the mbuf chain to reflect
     * the actual data length, not to be extended to contain
     * a fake, zeroed LLC header which keeps the snoop code from
     * crashing.
     */
    if ((m2 = m_copy(m, 0, len+sizeof(struct skheader))) == NULL)
        return (0);
    if (M_HASCL(m2)) {
        m2 = m_pullup(m2, SK_IBUFSZ);
        if (m2 == NULL)
            return (0);
    }
    sib = mtod(m2, struct sk_ibuf*);
    /*
     * MISSING: The DLPI code wants the MAC address in canonical bit order.
     * Convert here if necessary.
     */
    /*
     * MISSING:
     * The DLPI code wants the LLC header, if present,
     * not to be hidden with the MAC header. Decrement
     * LLC header size from ifh_hdrlen if necessary.
     */
    if ((*dlp->dl_infunc)(dlp, &si->si_if, m2, &sib->sib_skh)) {
        m_freem(m);
        return (1);
    }
    m_freem(m2);
}
```

```
    return (0);
}
/*
 * Process an ioctl request.
 * Return 0 or errno.
 */
static int
sk_ioctl(
    struct ifnet *ifp,
    int cmd,
    void *data)
{
    struct sk_info *si;
    int error = 0;
    int flags;
    ASSERT(IFNET_ISLOCKED(ifp));
    si = ifptosk(ifp);
    switch (cmd) {
    case SIOCSIFADDR:
    {
        struct ifaddr *ifa = (struct ifaddr *)data;
        switch (ifa->ifa_addr->sa_family) {
        case AF_INET:
            sk_stop(si);
            si->si_ac.ac_ipaddr = IA_SIN(ifa)->sin_addr;
            sk_start(si, ifp->if_flags);
            break;
        case AF_RAW:
            /*
             * Not safe to change addr while the
             * board is alive.
             */
            if (!iff_dead(ifp->if_flags))
                error = EINVAL;
            else {
                bcopy(ifa->ifa_addr->sa_data,
                    si->si_ac.ac_enaddr, SKADDRLEN);
                error = sk_start(si, ifp->if_flags);
            }
            break;
        default:
            error = EINVAL;
            break;
        }
    }
    break;
}
```

```
    }
    case SIOCSIFFLAGS:
    {
        flags = ((struct ifreq *)data)->ifr_flags;
        if (((struct ifreq*)data)->ifr_flags & IFF_UP)
            error = sk_start(si, flags);
        else
            sk_stop(si);
        break;
    }
    case SIOCADMULTI:
    case SIOCDELMULTI:
    {
#define MKEY ((union mkey*)data)
        int allmulti;
        /*
         * Convert an internet multicast socket address
         * into an 802-type address.
         */
        error = ether_cvtmulti((struct sockaddr *)data, &allmulti);
        if (0 == error) {
            if (allmulti) {
                if (SIOCADMULTI == cmd)
                    si->si_if.if_flags |= IFF_ALLMULTI;
                else
                    si->si_if.if_flags &= ~IFF_ALLMULTI;
                /* MISSING enable hw all multicast adrs */
            } else {
                bitswapcopy(MKEY->mk_dhost, MKEY->mk_dhost,
                    sizeof (MKEY->mk_dhost));
                if (SIOCADMULTI == cmd)
                    error = sk_add_da(si, MKEY, 1);
                else
                    error = sk_del_da(si, MKEY, 1);
            }
        }
        break;
#undef MKEY
    }
    case SIOCADDSNOOP:
    case SIOCDELSNOOP:
    {
#define SF(nm) ((struct skheader*)&(((struct snoopfilter *)data)->nm))
        /*
         * raw protocol snoop filter. See <net/raw.h>

```

```

    * and <net/multi.h> and the snoop(7P) man page.
    */
    u_char *a;
    union mkey key;
    a = &SF(sf_mask[0])->sh_dhost.sk_vec[0];
    if (!SK_ISBROAD(a)) {
        /*
         * cannot filter on device unless mask is trivial.
         */
        error = EINVAL;
    } else {
        /*
         * Filter individual destination addresses.
         * Use a different address family to avoid
         * damaging an ordinary multi-cast filter.
         * MISSING You'll have to invent your own
         * mulicast filter routines if this doesn't
         * fit your address size or needs.
         */
        a = &SF(sf_match[0])->sh_dhost.sk_vec[0];
        key.mk_family = AF_RAW;
        bcopy(a, key.mk_dhost, sizeof (key.mk_dhost));
        if (cmd == SIOCADDSNOOP)
            error = sk_add_da(si, &key, SK_ISGROUP(a));
        else
            error = sk_del_da(si, &key, SK_ISGROUP(a));
    }
    break;
}
/*
 * MISSING: add any driver-specific ioctls here.
 */
default:
    error = EINVAL;
}
return (error);
}
/*
 * Add a destination address.
 * Add address to the sw multicast filter table and to
 * our hw device address (if applicable).
 */
/* ARGSUSED */
static int
sk_add_da(

```

```
    struct sk_info *si,
    union mkey *key,
    int ismulti)
{
    struct mfreq mfr;
    /*
     * mfmacthcnt() looks up key in our multicast filter
     * and, if found, just increments its refcnt and
     * returns true.
     */
    if (mfmacthcnt(&si->si_filter, 1, key, 0))
        return (0);
    mfr.mfr_key = key;
    mfr.mfr_value = (mval_t) sk_dahash((char*)key->mk_dhost);
    if (!mfadd(&si->si_filter, key, mfr.mfr_value))
        return (ENOMEM);
    /* MISSING: poke this hash into device's hw address filter */
    return (0);
}
/*
 * Delete an address filter. If key is unassociated, do nothing.
 * Otherwise delete software filter first, then hardware filter.
 */
/* ARGSUSED */
static int
sk_del_da(
    struct sk_info *si,
    union mkey *key,
    int ismulti)
{
    struct mfreq mfr;
    /*
     * Decrement refcnt of this address in our multicast filter
     * and reclaim the entry if refcnt == 0.
     */
    if (mfmacthcnt(&si->si_filter, -1, key, &mfr.mfr_value))
        return (0);
    mfdel(&si->si_filter, key);
    /* MISSING: disable this hash value from the device if necessary */
    return (0);
}
/*
 * compute a hash value for destination addr
 */
static int
```

```
sk_dahash(char *addr)
{
    int hv;
    hv = addr[0] ^ addr[1] ^ addr[2] ^ addr[3] ^ addr[4] ^ addr[5];
    return (hv & 0xff);
}
/*
 * Periodically poll the device for input packets
 * in case an interrupt gets lost or the device
 * somehow gets wedged. Reset if necessary.
 */
static void
sk_watchdog(struct ifnet *ifp)
{
    struct sk_info *si;
#ifdef INTR_KTHREADS
    int s;
#endif
    si = ifptosk(ifp);
    ASSERT(IFNET_ISLOCKED(ifp));
    /* check for a missed interrupt */
    IFNET_UNLOCKNOSPL(ifp);
    sk_intr(si);
    IFNET_LOCKNOSPL(ifp);
    si->si_if.if_timer = SK_DOG;
}
/*
 * Disable the interface.
 */
static void
sk_stop(struct sk_info *si)
{
    struct ifnet *ifp = sktoifp(si);
    ASSERT(IFNET_ISLOCKED(ifp));
    ifp->if_flags &= ~IFF_ALIVE;
    /*
     * Mark an interface down and notify protocols
     * of the transition.
     */
    if_down(ifp);
    sk_reset(si);
}
/*
 * Enable the interface.
 */
```

```
static int
sk_start(
    struct sk_info *si,
    int flags)
{
    struct ifnet *ifp = sktoifp(si);
    int error;
    ASSERT(IFNET_ISLOCKED(ifp));
    error = sk_ifinit(ifp);
    if (error)
        return (error);
    ifp->if_flags = flags | IFF_ALIVE;
    /*
     * Broadcast an ARP packet, asking who has addr
     * on interface ac.
     */
    arpwhohas(&si->si_ac, &si->si_ac.ac_ipaddr);
    return (0);
}
/*
 * private debugging routine.
 */
static void
sk_dump(int unit)
{
    struct sk_info *si;
    struct ifnet *ifp;
    char name[128];
    if (unit == -1)
        unit = 0;
    sprintf(name, "sk%d", unit);
    if ((ifp = ifunit(name)) == NULL) {
        qprintf("sk_dump: %s not found in ifnet list\n", name);
        return;
    }
    si = ifptosk(ifp);
    qprintf("si 0x%x\n", si);
    /* MISSING: qprintf() whatever you want here */
}
```

## **PART NINE**

# **PCI Drivers**

### **Chapter 15: PCI Device Drivers**

Overview of the architecture of the PCI bus attachment and the services offered by the kernel to PCI device drivers.



---

## PCI Device Drivers

The Peripheral Component Interconnect (PCI) bus, initially designed at Intel Corp, is standardized by the PCI Bus Interest Group, a nonprofit consortium of vendors (see “Standards Documents” on page xxviii and “Internet Resources” on page xxvii).

The PCI bus is designed to be a high-performance local bus to connect peripherals to memory and a processor. In many personal computers based on Intel and Motorola processors, the primary system bus is a PCI bus. A wide range of vendors make devices that plug into the PCI bus.

The PCI bus is supported by the O2 workstation as well as other SGI systems. However, the O2 was the first SGI system to support PCI, and IRIX 6.3 was the first release of IRIX to contain PCI bus support. This chapter contains the following topics related to support for the PCI bus:

- “PCI Bus in Silicon Graphics Workstations” on page 376 gives an overview of PCI bus features and implementation.
- “PCI Implementation in O2 Workstations” on page 378 describes the hardware features and restrictions of one PCI implementation.
- “Driver/Kernel Interface for PCI Access” on page 383 discusses the kernel functions that are specifically designed to support PCI device drivers.
- “Example Driver” on page 405 displays the code of a complete PCI driver.

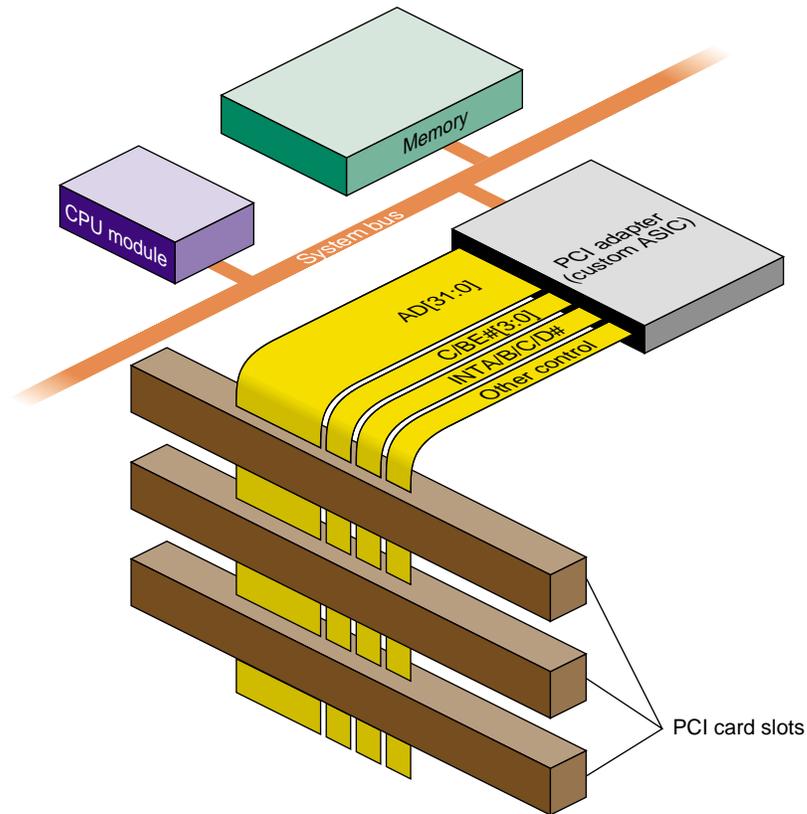
A PCI driver is a kernel-level device driver like other drivers. For information on the architecture of a kernel-level device driver and on how to build and debug one, see Part III, “Kernel-Level Drivers.”

## **PCI Bus in Silicon Graphics Workstations**

This section contains an overview of the main features of PCI hardware attachment, for use as background material for software designers. Hardware designers can obtain a detailed technical paper on PCI hardware through the Silicon Graphics Developer Program. Design issues such as power supply capacities, card dimensions, signal latencies, and arbitration, are covered in that material.

### **PCI Bus and System Bus**

In no Silicon Graphics system is the PCI bus the primary system bus. The primary system bus is always a proprietary bus that connects one or more CPUs with high-performance graphics adapters and main memory. The PCI bus adapter is connected (or “bridged,” in PCI terminology) to the system bus, as shown in Figure 15-1.



**Figure 15-1** PCI Bus In Relation to System Bus

The PCI adapter is a custom circuit with these main functions:

- To act as a PCI bus target when a PCI bus master requests a read or write to memory
- To act as a PCI bus master when a CPU requests a PIO operation
- To manage PCI bus arbitration, allocating bus use to devices as they request it
- To interface PCI interrupt signals to the system bus and the CPU

Different SGI systems have different PCI adapter ASICs. Although all adapters conform to the PCI standard level 2.1, there are significant differences between them in capacities, in optional features such as support for the 64-bit extension, and in performance details such as memory-access latencies.

## Buses, Slots, Cards, and Devices

A system may contain one or more PCI bus adapters. Each bus connects one or more physical *packages*. The PCI standard allows up to 32 physical packages on a bus. A “package” may consist of a card plugged into a slot on the bus. However, a “package” can also consist of an internal chipset mounted directly on the system board, using the PCI bus and occupying one or more virtual slots on the bus. For example, the SCSI adapter in the O2 workstation occupies the first two virtual slots of the PCI bus in that system.

Each physical package can implement from one to eight *functions*. A PCI function is an independent device with its own configuration registers in PCI configuration space, and its own address decoders.

In Silicon Graphics systems, each PCI *function* is integrated into IRIX as a *device*. A PCI device driver manages one or more devices in this sense. A driver does not manage a particular package, or card, or bus slot; it manages one or more logical devices.

## PCI Implementation in O2 Workstations

In the O2 workstation, a proprietary system bus connects the CPU, multimedia devices (audio, video, and graphics) and main memory.

The PCI bus adapter interfaces one PCI bus to this system bus. The PCI bus adapter is a unit on the system bus, on a par with the other devices. The PCI bus adapter competes with the CPU and with multimedia I/O for the use of main memory.

The built-in SCSI adapter, which is located on the main system board, is logically connected to the PCI bus and takes the place of the first two “slots” on the PCI bus, so that the first actual slot is number 2.

## Unsupported PCI Signals

In the O2, the PCI adapter implements a standard, 32-bit PCI bus operating at 33 MHz. The following optional signal lines are not supported.

- The LOCK# signal is ignored; atomic access to memory is not supported.
- The cache-snoop signals SBO# and SDONE are ignored.
- The JTAG signals are not supported.

## 64-bit Address and Data Support

The O2 PCI adapter supports 64-bit data transfers, but not 64-bit addressing. All bus addresses are 32 bits, that is, all PCI bus virtual addresses are in the 4 GB range. The Dual Address Cycle (DAC) command is not supported (or needed).

The 64-bit extension signals AD[63:32], C/BE#[7:4], REQ64# and ACK64# are pulled up as required by the PCI standard.

When the PCI bus adapter operates as a bus master (as it does when implementing a PIO load or store for the CPU), the PCI adapter generates 32-bit data cycles.

When the PCI bus adapter operates as a bus target (as it does when a PCI bus master transfers data using DMA), the PCI adapter does not respond to REQ64#, and hence 64-bit data transfers are accomplished in two, 32-bit, data phases as described in the PCI specification.

### Configuration Register Initialization

When the IRIX kernel probes the PCI bus and finds an active device, it initializes the device configuration registers as follows:

Command Register	The enabling bits for I/O Access, Memory Access, and Master are set to 1. Other bits, such as Memory Write and Invalidate and Fast Back-to-Back are left at 0.
Cache Line Size	0x20 (32, 32-bit words, or 128 bytes).
Latency Timer	0x30 (48 clocks).
Base Address registers	Each register that requests memory or I/O address space is programmed with a starting address. In the O2 system, memory addresses are always greater than 0x8000 0000.

The device driver is free to set any other configuration parameters in the *pfxattach0* entry point (see “Attaching a Device” on page 386).

### Address Spaces Supported

The relationship between the PCI bus address space and the system memory physical address space differs from one system type to another.

The O2 workstation and related systems support a 1 GB physical memory address space (30 bits of physical address used). Any part of physical address space can be mapped into PCI bus address space for purposes of DMA access from a PCI bus master device. The device driver ensures correct mapping through the use of a DMA map object (see “Allocating DMA Maps” on page 390).

Physical memory can be mapped to the PCI bus in normal order or byte-swapped order. Byte-swapping is done on the basis of 64-bit units. When a PCI bus master uses a byte-swapped DMA address as its target, and writes the 64-bit data item 0x0807 0605 0403 0201, the data 0x0102 0304 0506 0708 is delivered to memory.

### **PIO Address Mapping**

For PIO purposes (the CPU loading and storing in device space), memory space defined by each PCI device in its configuration registers is allocated in the upper two gigabytes of the PCI address space, above 0x8000 0000. These addresses are allocated dynamically, based on the contents of the configuration registers of active devices. The I/O address space requested by each PCI device in its configuration registers is also allocated dynamically as the system comes up.

It is possible for a PCI device to request (in the initial state of its Base Address Registers) that its address space be allocated in the first 1 MB of the PCI bus. This request cannot be honored in the O2 workstation. Devices that cannot decode bus addresses above 0x8000 0000 are not supported.

Device drivers get a virtual address to use in PIO access by creating a PIO map (see “Managing PIO Maps for PCI” on page 396).

### **Slot Priority and Bus Arbitration**

Two devices that are built in to the workstation take the positions of PCI bus slots 0 and 1. Actual bus slots begin with slot 2 and go up to the maximum for the system (just the one slot in O2).

The PCI adapter maintains two priority groups. The lower-priority group is arbitrated in round-robin style. The higher-priority group uses fixed priorities based on slot number, with the higher-numbered slot having the higher fixed priority.

The IRIX kernel assigns slots to priority groups dynamically by storing values in an adapter register. There is no kernel interface for changing this priority assignment. The audio and the available PCI slots are in the higher priority group.

## Interrupt Signal Distribution

The PCI adapter can present eight unique interrupt signals to the system CPU. The IRIX kernel uses these interrupt signals to distinguish between the sources of PCI bus interrupts. The system interrupt numbers 0 through 7 are distributed across the PCI bus slots as shown in Table 15-1 (“n.c.” means no connection).

**Table 15-1** PCI Interrupt Distribution to System Interrupt Numbers

PCI Interrupt	Slot 0 (built-in device)	Slot 1 (built-in device)	Slot 2	Slot 3 (When Present)	Slot 4 (When Present)
INTA#	system 0	n.c.	system 2	system 3	system 4
INTB#	n.c.	system 1	system 5	system 7	system 6
INTC#	n.c.	n.c.	system 6	system 5	system 7
INTD#	n.c.	n.c.	system 7	system 6	system 5

Each physical PCI slot has a unique system interrupt number for its INTA# signal. The INTB#, INTC#, and INTD# signals are connected in a spiral pattern to three system interrupt numbers.

## Driver/Kernel Interface for PCI Access

A PCI device driver manages the operation of one or more devices. In this section, “device” has two meanings.

- A device can be a function associated with one set of configuration registers on a PCI card. A PCI card can contain up to eight such functions, but each configuration space is treated as a separate device by IRIX..
- A logical device is a device special file in the */dev* filesystem that refers to the original PCI device. For example, a PCI card that contains a serial port might be associated with */dev/ttyd12*, */dev/ttyf12*, and */dev/ttym12*.

Besides the usual driver entry points for a block or character driver, a PCI driver has to supply the *pxattach()* entry point. This entry point is called to initialize the PCI device and, optionally, to identify any additional logical devices for that PCI device. An entry point named *pxdetach()* is optional.

Besides the usual DDI/DKI functions, the PCI driver calls on kernel functions unique to the PCI bus. These functions are introduced in the following topics. For a summary, see “PCI Function Summary” on page 403.

### Overview of PCI Driver Structure

A PCI driver is a kernel-level device driver that has the general structure described in Chapter 8, “Structure of a Kernel-Level Driver.” It uses the driver/kernel interface described in Chapter 9, “Device Driver/Kernel Interface.” A PCI driver can be loadable or it can be linked with the kernel. In general it is configured into IRIX as described in Chapter 10, “Building and Installing a Driver.”

PCI hardware configuration is more dynamic than the configuration of the VME, EISA or SCSI buses. With other types of bus, the driver learns the hardware configuration when the driver is loaded, and the configuration remains static afterward. IRIX support for the PCI bus is designed to allow support for dynamic reconfiguration in future systems. In principle, a PCI driver has to be designed to allow devices to be attached and detached at any time (no detaching is done in the current release).

The general sequence of operations of a PCI driver is as follows:

1. In the *pfxinith()* entry point, the driver calls a kernel function to register itself as a PCI driver, specifying the kind of device it supports.
2. When the kernel discovers a device of this type, it calls the *pfxattach()* entry point of the driver. The driver creates PIO maps and (optionally) DMA maps to use in addressing the device; initializes the device; and if necessary, registers an interrupt handler.
3. In the normal upper-half entry points such as *pfxopen()*, *pfxread()*, and *pfstrategy()*, the driver operates the device and transfers data.
4. When the device generates an interrupt, the driver's registered interrupt handler is called.
5. If the kernel learns of a bus address error on the PCI bus, it can call an error-handling function registered by the driver to find out which device caused the error.
6. Conceptually, if the kernel learns that the device is being detached, the kernel calls the driver's *pfxdetach()* entry point. The driver notes the device is unusable and stops servicing it through upper-half entry points. (This feature is not implemented in the current release.)

### Driver Flag Bits

As described under "Driver Flag Constant" on page 145, the *pfxdevflag* public name is a byte containing flags for driver characteristics. Since a PCI driver is inevitably a new driver, with no heritage in older versions of IRIX or UNIX, Silicon Graphics, Inc. *strongly recommends* that you design it from the start to be compatible with multiprocessors. The implications of this are discussed under "Planning for Multiprocessor Use" on page 174.

### Initializing and Registering the Driver

A PCI driver must register with the kernel in order to receive notification that devices exist. In the current release this is done in two stages. First the driver calls **pciio\_add\_attach()** to introduce itself to the kernel. See reference page pciio(d3) for the function prototype. The arguments passed in this call are as follows:

- |               |   |
|---------------|---|
| <i>attach</i> | Address of the driver's <i>pfxattach()</i> entry point. This is required.       |
| <i>detach</i> | NULL, or address of the driver's <i>pfxdetach()</i> entry point if implemented. |

*error*            NULL, or address of the driver's error-handling function if any.

*driver\_prefix*   Pointer to the driver's prefix as a character string.

*major*            The major number supported by this driver. The number given in the third field of the descriptive line (see "Descriptive Line" on page 240).

This call associates the driver's *pfattach()* entry point with the driver's prefix string. The *pfdetach()* and error-handler addresses may be passed as NULL when these are not implemented. There is no need to implement *pfdetach()* in IRIX 6.3.

**Note:** This call to *pciio\_add\_attach()* is unique to IRIX 6.3; it is not required or allowed in subsequent releases. You can compile the call conditionally based on the definition of the variable *\_EARLY\_PCI*, which is not defined in later releases—as demonstrated in Example 15-1.

**Tip:** You can create a static list of major numbers as a global variable in the driver descriptive file. See "Variables Section" on page 237 for an example of such an array.

The next step—and the *only* step required in later releases—is to call the *pciio\_driver\_register()* function (see reference page *pciio(d3)* for the syntax). This call specifies the PCI vendor ID and device ID numbers as they appear in the PCI configuration space of any device that this driver supports. The third argument is a character string containing the driver's prefix string (as passed to *pciio\_add\_attach()*). The kernel uses this string to search the switch tables to find the addresses of the driver's *pfattach()* and *pfdetach()* entry points.

Example 15-1 shows a hypothetical example of driver registration.

#### Example 15-1    Driver Registration

```
__int32_t hypo_attach(dev_t); /* forward declaration */
int hypo_init()
{
    int allMyMajors[2];
    ...
    allMyMajors[0] = myMajNum; /* see Example 10-1 */
    allMyMajors[1] = 0;
#ifdef _EARLY_PCI /* need to call pciio_add_attach */
    ret = pciio_add_attach(hypo_attach, /* attach entry point */
                          NULL,NULL, /* detach, error not done */
                          "hypo_", /* prefix string */
                          allMyMajors); /* list of one major# */
#endif
}
```

```
    if (!ret) ...
#endif
    ret = pciio_driver_register(HYPO_VENDOR, HYPO_DEVID, "hypo_", 0);
    if (!ret) ...
}
```

A device driver can register multiple times to handle multiple combinations of vendor ID and device ID.

**Tip:** You should defer the call to `pciio_driver_register()` to the end of the `pfxinit()` routine, when all global data has been initialized. The reason is that, if there is an available device of the specified type, `pfxattach()` might be called immediately, before the `pci_driver_register()` function returns. In a multiprocessor, `pfxattach()` can be called concurrently with the return of `pci_driver_register()` and following code.

A loadable driver, when called at its `pfxunload()` entry point, can unregister before unloading; but that is not required. (See “Unloading” on page 402).

## Attaching a Device

The IRIX support for PCI is designed to allow for future support for hot-swapping and for multinode systems in which devices, slots, buses, and whole nodes come online and offline dynamically. In principle, a PCI device can be attached or detached at any time, meaning that the `pfxattach()` and `pfxdetach()` entry points could be called at any time.

In practice, the only systems supported by IRIX 6.3 for O2 do not permit hot-swapping. Also, the administrator commands that would force a device to detach are not defined as yet. In current systems, `pfxattach()` is only called at boot time and `pfxdetach()` is not called. Nevertheless the driver/kernel interface is designed for future flexibility. For future portability you can design your driver as if that flexibility existed now.

## Matching A Device to A Driver

When the system boots up, the kernel probes the PCI bus configuration space and takes a census of active devices. For each device it notes

- Vendor and device ID numbers
- Requested size of memory space
- Requested size of I/O space

The kernel assigns starting bus addresses for memory and I/O space and sets these addresses in the Base Address Registers (BARs) in the device. Then the kernel looks for a driver that has registered a matching set of vendor and device IDs using `pciio_driver_register()`.

If no matching driver has registered, the device remains inactive. For example, the driver might be a loadable driver that has not been loaded as yet. When the driver is loaded and registers, the kernel will match it to any unattached devices.

When the kernel matches a device to its registered driver, the kernel calls the driver's `pfxattach()` entry point. It passes one argument, a `vertex_hdl_t`, which acts as an opaque handle to a kernel object that describes this device. This handle is used to:

- Store and retrieve the driver's private information about the device
- Request PIO and DMA maps on the device
- Register an interrupt handler for the device

### **Allocating Storage for Device Information**

A driver needs to save information about each device, usually in a structure. Fields in a typical structure might include:

- Locks or semaphores used for mutual exclusion among upper-half entry points and between them and the interrupt handler.
- Addresses of allocated PIO and DMA maps for this device (see "Allocating PIO Maps" on page 388).
- Address of an interrupt connection object for the device (see "Registering an Interrupt Handler" on page 391).
- In a block driver, anchors for a queue of `buf_t` objects being filled or emptied.
- Device status information and flags.

A problem is that at initialization time a driver does not know how many devices it will be asked to manage. For a workstation such as O2 you can expect the number will be small, but you should allow for portability to server-class systems that support dozens of devices. In the past this problem has been handled by allocating an array of a fixed number of information structures, indexed by the device minor number.

This is not a good solution for a PCI driver because PCI configuration is so dynamic. In addition, a loadable driver loses the contents of its global variables when it unloads. The IRIX PCI support gives you a different way.

In a PCI driver, you dynamically allocate memory for an information structure to hold information about the one device being attached. (See “General-Purpose Allocation” on page 190.) You save the address of the structure in the kernel’s hardware vertex, using the **device\_info\_set()** function, which associates an arbitrary pointer with a *vertex\_hdl\_t*.

```
extern void device_info_set(vertex_hdl_t, arbitrary_info_t);
```

The information structure can easily be recovered in any top-half routine; see “Locating Device Information” on page 395.

### Allocating PIO Maps

For almost any device you need at least one PIO map. You use a PIO map to read the memory space or the I/O space of a device. These maps can be allocated when the device is attached, and the addresses of the maps can be stored in the device information structure.

A PIO map is created by **pciio\_piomap\_alloc()**. See reference page `pciio_pio(d3)` for the syntax. This call requires arguments are as follows:

<i>vhdl</i>	The <i>vertex_hdl_t</i> received by the <b>pfattach()</b> routine. Every map must be associated to a specific device at its hardware graph vertex.
<i>desc</i>	Device descriptor structure with one field set (see text following).
<i>space</i>	Constant specifying the space to map: <code>PCIIO_SPACE_CFG</code> (configuration space), or <code>PCIIO_SPACE_WIN(n)</code> .
<i>addr</i>	Offset within the selected <i>space</i> (typically 0).
<i>size</i>	Span of the total area over which this map might be applied.
<i>max</i>	Maximum size of the area that will be mapped at any one time. When the map is always used for the same area, <i>size</i> and <i>max</i> are the same. When the map can be used for smaller segments within a larger area, <i>max</i> is the limit of one segment and <i>size</i> the size of the total extent.
<i>flags</i>	Passed as 0.

A PIO map that you will use to access device configuration registers is based on a space of `PCIIO_SPACE_CFG`. The space selection `PCIIO_SPACE_WIN(n)` means that this map is to be based on Base Address register  $n$ , from 0 through 5, in the PCI configuration space. The device configuration registers specify whether a given base address register defines memory or I/O space. When the space is defined by a 64-bit base address register, use the lower number, the index of the word that contains the configuration bits.

**Tip:** The header file `sys/PCI/pciio.h` declares constants of the form `PCIIO_PIOMAP_CFG` and `PCIIO_PIOMAP_WIN(n)`. You can ignore these; they are not used in any calls. The target space for any kind of map is given with `PCIIO_SPACE_*`.

The device descriptor structure type `dev_desc_t` is declared in `pciio.h`. A descriptor structure is required in this call, but only one field is inspected, `intr_swlevel`. It must be set to one of the interrupt levels of type `pl_t` as declared in `ddi.h`, typically `plhi`, as shown in Example 15-2

**Note:** In subsequent releases the device descriptor is not required and the address can be passed as `NULL`; but for this release it is required.

Example 15-2 shows a function that simplifies the allocation of a PIO map. The space is passed as an argument, as is the size of the space to map. The function makes the simplifying assumption that the map should start at offset 0 in the selected space.

**Example 15-2** Allocation of PCI PIO Map

```
pciio_piomap_t makeMap(vertex_hdl_t dev, int base, size_t size)
{
    struct device_desc_s ddesc;
    ddesc.intr_swlevel = plhi;
    return pciio_piomap_alloc(
        dev,          /* vertex handle */
        &ddesc,        /* dev descriptor w/ in level in it */
        base,         /* space, _CFG or _WIN(n) */
        0,            /* starting offset */
        size, size,   /* size to map */
        0);           /* default endian */
}
```

### Allocating PIO Addresses Directly

In the O2 and some other platforms, PIO addressing is based on fixed hardware resources and a PIO address can be generated without use of a map. You request this using the function `pciio_piotrans_addr()`. In the general case of a PIO address for memory or I/O space, this function call can fail in some systems (as discussed in reference page `pciio_pio(d3)`). When used to obtain a PIO address for configuration space, it generally succeeds in all systems.

### Allocating DMA Maps

For a bus-master device you will need at least one DMA map. A DMA map is created by `pciio_dmamap_alloc()`, which takes a *vertex\_hdl\_t*, a size, and flags regarding the treatment of the mapping. (See reference page `pciio_dma(d3)` for the syntax of this and related functions.)

More map functions are discussed under “Managing PIO Maps for PCI” on page 396 and “Managing DMA Maps for PCI” on page 399.

### Reading the Device Configuration

Typically a PCI driver needs to read the device configuration registers and possibly write to them. In principle, these are PIO operations. However, in the O2 (and possibly other platforms), the special PCI bus cycle called a Configuration cycle is not generated by a simple PIO load or store.

To access the configuration, create a PIO map for the configuration space and extract an address from it. Present this address to `pciio_config_get()` to fetch a word from configuration space, as shown in Example 15-3. (See reference page `pciio_config(d3)` for the syntax of this and related functions.)

The function in Example 15-3 fetches and returns a 32-bit word from configuration space. It obtains a PIO base address for configuration space using `pciio_piotrans_addr()` (see “Allocating PIO Addresses Directly” on page 390). When that call succeeds, as it does in the O2 and other current SGI systems, the address is passed to `pciio_config_get()`.

**Example 15-3** Reading PCI Configuration Space

```

__uint32_t get_cfg_word( vertex_hdl_t vh, int reg)
{
    device_desc_t dd = {0};
    volatile __uint32_t *pio_addr;
    dd.intr_swlevel = plhi;
    pio_addr = pciio_piotrans_addr(vh, dd,
        PCIIO_SPACE_CFG,0,256,0);
    if (pio_addr) /* trans_addr succeeded */
        return pciio_config_get(pio_addr,reg);
    else /* trans_addr failed, simulate hardware failure */
        return (__uint32_t)(-1);
}

```

For a PCI bus master device, the *pfxattach()* function should set the Cache Line Size register to 128 (the size of a cache line in all Silicon Graphics systems).

**Registering an Interrupt Handler**

For devices that can interrupt, a key step during *pfxattach()* is to register an interrupt handler for the device. This is done in a two-step process. First you create an interrupt connection object; then you use that object to specify the interrupt handling function.

The interrupt connection is created with *pciio\_intr\_alloc()*, which takes a *vertex\_hdl\_t* and a flag for the interrupt line that the device uses. (See reference page *pciio\_intr(d3)* for the syntax.)

The interrupt object is used in establishing a handler, and it is needed later to stop taking interrupts (see “Inactivating an Interrupt Handler” on page 401). You probably want to save its address in the device information structure for later use.

After creating the interrupt object, you establish a handler using *pciio\_intr\_connect()*. Its principal arguments are the interrupt object, a handler address, and a value to be passed to the handler when it is called—typically the address of the device information structure you are preparing. If a device will interrupt on line C, interrupt setup could resemble Example 15-4.

**Example 15-4** Setting Up a PCI Interrupt Handler

```
pciio_intr_t intobj;
extern void int_handler(eframe_t *, void *);
int retcode;
intobj = pciio_intr_alloc(
    vhdl, /* as received in attach() */
    0, /* device descriptor is n.a. for pci */
    PCIIO_INTR_LINE_C, /* the line it uses */
    (vertex_hdl_t) 0);
retcode = pciio_intr_connect(
    intobj, /* the interrupt object */
    (intr_func_t) int_handler, /* the handler */
    (intr_arg_t) pDevInfo, /* dev info as input */
    (void*)0 ); /* threads are next release */
if (!retcode) cmn_err(CE_WARN, "oh fiddlesticks");
```

**Note:** The declaration of the interrupt handler function type, *intr\_func\_t*, requires two arguments, the first being an *eframe\_t*. The availability of the *eframe\_t* is unique to IRIX 6.3 for O2. In subsequent releases the interrupt handler receives only one argument, the value passed with the **pciio\_intr\_connect()** call.

The CPU is accepting interrupts when the *pfattach()* entry point is called. If the PCI device is in a state that can produce an interrupt, the interrupt handling function can be called before **pciio\_intr\_connect()** returns. Make sure that all global data used by the interrupt handler has been initialized.

PCI devices can share the four PCI interrupt lines. As a result, in some cases the kernel cannot tell which device caused an interrupt. When there is any doubt, the kernel calls all the interrupt handlers that are registered to that interrupt line. For this reason, your interrupt handler must not assume that its device did cause the interrupt. It should always test to see if an interrupt is really pending, and exit immediately when one is not.

### Return Value from Attach

The return code from *pxattach()* is tested by the kernel. The driver can reject an attachment. When your driver cannot allocate memory, or fails due to another problem, it should:

- Use **cmn\_err()** to document the problem (see “Using *cmn\_err*” on page 251)
- Release any objects such as PIO and DMA maps that were created.
- Release any space allocated to the device such as a device information structure.
- Return an informative return code which might be meaningful in future releases.

More than one driver can register to support the same vendor ID and device ID. When the first driver fails to complete the attachment, the kernel continues on to test the next, until all have refused the attachment or one accepts it. The *pxdetach()* entry point can only be called if the *pxattach()* entry point returns success (0).

### Establishing Logical Devices

Some kinds of physical devices are represented by multiple device special files in */dev*. For example, each serial port appears as at least four devices */dev/tty\**. A tape drive can appear under different names, and a disk device has two device special files for each disk partition, one in */dev/dsk* and one in */dev/rdisk* (raw, or character access). Each logical device represents a slightly different treatment of the same physical device.

The *pxattach()* entry point initializes the real PCI device, but it must also create hardware vertexes to represent the logical devices that should be associated with the same real PCI device. This is done in three steps:

- Create a new hardware vertex connected to the attached vertex, using **hwgraph\_device\_add()**.
- Associate the new vertex with the minor number of the logical device, using **hwgraph\_device\_add\_minor()**.
- Associate the new vertex with the same device information structure, using **device\_info\_set()**.

The `hwgraph_device_add()` function has the following prototype:

```
int hwgraph_device_add(vertex_hdl_t vhdl, /* parent vertex */
                      char *name, /* name of the device */
                      char *prefix, /* driver prefix */
                      vertex_hdl_t *new_vrtx) /* return result */
```

The name argument is not significant in the current release, but it will be significant and visible to users in a future release. It should be one word or numeric characters to label this vertex of the hardware graph, for example “0” (logical unit number) or “nonswap” (feature or access method).

For a simplified example, see Example 15-5. This hypothetical code, which would be part of the `hypo_attach()` entry point, creates two logical devices. The device minor number of the first is 0x01; the second is 0x02—a simplified version of the conventions for minor numbers of tape or serial devices, in which the minor number bits represent device options or features.

**Example 15-5** Creating Logical Devices for a PCI Device

```
vertex_hdl_t subdev;
int ret;
my_dev_info_t *pDev; /* struct stored in PCI vertex */
...
ret = hwgraph_device_add(vhdl, /* attach() input */
                        "left", /* name of minor 01 */
                        "hypo_", /* driver prefix */
                        &subdev); /* output here */

if (ret)...
ret = hwgraph_device_add_minor(subdev, (minor_t)0x01);
if (ret)...
device_info_set(subdev, my_dev_info);
ret = hwgraph_device_add(vhdl, /* attach() input */
                        "right", /* name of minor 02 */
                        "hypo_", /* driver prefix */
                        &subdev); /* output here */

if (ret)...
ret = hwgraph_device_add_minor(subdev, (minor_t)0x02);
if (ret)...
device_info_set(subdev, my_dev_info);
```

## Normal Operation

While handling normal operations on the device, the driver needs to locate device information from the top-half entry points, and needs to translate addresses using maps.

### Locating Device Information

The driver upper-half entry points are called to implement requests from user processes or filesystems that need to open, read, write, map or control the device. These calls can occur at any time; and on a multiprocessor, they can occur multiple times concurrently, on parallel CPUs.

The user process refers to a device through a file descriptor opened to a device special file. The primary argument to any upper-half entry point is the *dev\_t*, a value that distinguishes the device. Traditional drivers extract device numbers from the *dev\_t* (see “The Device Number Types” on page 182).

The first thing any upper-half entry point needs to do is to locate the per-device information structure that was prepared in the *pxattach()* entry point (see “Allocating Storage for Device Information” on page 387). You do this by calling **device\_info\_get()**. However, that function takes a *vertex\_hdl\_t*. You get that from **dev\_to\_vhdl()**. The code, which is repeated over and over in a PCI driver, resembles Example 15-6.

#### Example 15-6 Retrieving Device Information

```
vertex_hdl_t vhd1 = dev_to_vhdl(dev);
my_dev_info_t *pDev = (my_dev_info_t)device_info_get(vhd1);
    if (!pDev) return(ENXIO);
    if (!(pDev->status & USABLE)) return(ENXIO);
```

In the *pxopen()* entry point, the *dev\_t* is received as a reference argument, not by value.

### Verifying Device Usability

In the event that **device\_info\_get()** returns NULL, this device has not been attached, or the *pxattach()* entry point did not allocate an information structure; or the *pxdetach()* entry point has been called. In such cases, the upper-half routine should return ENXIO (no such device). This test is shown in Example 15-6.

Example 15-6 also shows another test. In a future release of IRIX, a driver will be able to implement a *pfxdisable()* entry point, called to make a device temporarily unusable. Even in the current release, your driver might find reasons, such as a persistent device error, to force a device offline. A single flag bit in the device information structure represents this state. Again, a return of ENXIO is appropriate.

### Managing PIO Maps for PCI

The functions that are used to manage PIO maps are summarized in Table 15-2. See reference page *pciio\_pio(d3)* for details.

**Table 15-2** Functions for PIO Maps for PCI

Function	Purpose and Operation
<i>pciio_piomap_addr()</i>	Get a kernel virtual address from a PIO map for a specific offset and length.
<i>pciio_piomap_alloc()</i>	Create a PIO map object, specifying the bus address space, base offset, and length it needs to cover.
<i>pciio_piomap_done()</i>	Make a PIO map inactive until it is next needed (may release hardware resources associated to the map).
<i>pciio_piomap_free()</i>	Release a PIO map object.
<i>pciio_piotrans_addr()</i>	Request immediate translation of a bus address to a kernel virtual address without use of a PIO map. Returns NULL if this system does not support fixed PIO addressing for the requested space.
<i>pciio_config_get()</i>	Fetch a 32-bit value from configuration space using an address returned by <b><i>pciio_piomap_addr()</i></b> .
<i>pciio_config_set()</i>	Store a 32-bit value into configuration space using an address returned by <b><i>pciio_piomap_addr()</i></b> .

Maps are allocated with ***pciio\_piomap\_alloc()***. Its use is covered under “Allocating PIO Maps” on page 388, because you typically will allocate the PIO maps you need while attaching the device.

You use a PIO map by calling ***pciio\_piomap\_addr()***. This function takes a map, an offset within the PCI address space described by the map, and the number of bytes of space that the address will be used to retrieve.

The *pciio\_addr* argument is added to the base offset specified to **pciio\_piomap\_alloc()**, and that in turn is added to the base address assigned by the kernel to this device, to arrive at the bus address needed. The *byte\_count* argument specifies how many bytes beyond the bus address you may access. The returned value is a kernel virtual address that is mapped to the requested PCI bus address for at least that many bytes.

**Tip:** A program variable based on a PIO address should always be declared as “volatile.”

Once you have extracted an address using **pciio\_piomap\_addr()**, the map is active. It remains active until you call either **pciio\_piomap\_done()** or **pciio\_piomap\_free()**. In some systems, it costs nothing to keep a PIO map active. In other systems, an active PIO map may tie up global hardware resources. It is a good idea to call **pciio\_piomap\_done()** when the address is no longer needed.

Some systems also support a one-step translation function, **pciio\_piotrans\_addr()**, as described under “Allocating PIO Addresses Directly” on page 390. However, this function can fail in systems that do not use hard-wired bus maps. The two-step process of allocating a map and then interrogating it is more general.

**pciio\_piotrans\_addr()** always succeeds when getting a PIO address in configuration space. Access to configuration space is done in two steps. First you get a PIO address; then you pass the address to **pciio\_config\_get()** or **pciio\_config\_get()**. An example is shown under “Reading the Device Configuration” on page 390..

### Using Byte-Level PCI Addresses

When you use PIO to fetch or store 32-bit values on 32-bit-aligned PCI addresses, PIO works as you would expect, and a 32-bit value is fetched or returned.

However, when you use PIO to fetch or store less than 32 bits—either a 16-bit value or an 8-bit value—you must use an address that takes account of byte-swapping. The least significant address bits are summarized in Table 15-3.

**Table 15-3** Least Significant Address Bytes for Short PIO

Binary Offset Within 32-bit Memory Word	LSB for 16-Bit Access	LSB for 8-bit Access
0x00	0x02	0x03
0x01	n.a.	0x02
0x02	0x00	0x01
0x03	n.a.	0x00

You can deal with this complication in either of three ways, as follows:

- Always fetch and store 32-bit words. Unpack smaller units in memory.
- Declare the device data as a structure and arrange the order of short fields in the structure so that the least significant address bits work out correctly. For example if the device offers the following logical structure in its PCI memory space:

```
00--> dma base addr, 4 bytes
04--> dma counter, 4 bytes
08--> control reg, 2 bytes
0A--> status reg, 1 byte
0B--> byte fifo, 1 byte
```

Declare this as a C structure as follows:

```
struct {
    unsigned    dma_addr;
    unsigned    dma_count;
    unsigned char byte_fifo;
    unsigned char status;
    unsigned short control;
}
```

- Write C macros for byte-address and halfword-address. The macros would use exclusive-OR to invert two or one (respectively) of the least-significant bits.

## Managing DMA Maps for PCI

The functions that are used to manage simple DMA maps are summarized in Table 15-4. See reference page `pciio_dma(d3)` for syntax.

**Table 15-4** Functions for Simple DMA Maps for PCI

Function	Purpose and Operation
<code>pciio_dmamap_alloc()</code>	Create a DMA map object, specifying the maximum extent of memory the map will have to cover.
<code>pciio_dmamap_addr()</code>	Get the bus virtual address corresponding to a memory address for a specified length.
<code>pciio_dmamap_done()</code>	Make a DMA map inactive (may release hardware resources associated to the map).
<code>pciio_dmamap_free()</code>	Release a DMA map object.
<code>pciio_dmatrans_addr()</code>	Request immediate translation of the address of a contiguous memory buffer to a bus address. Returns NULL unless this system supports fixed DMA addressing

Maps are allocated with `pciio_dmamap_alloc()`. Its use is covered under “Allocating PIO Maps” on page 388, because you typically will allocate the maps you need while attaching the device.

You obtain a map for a single, contiguous span of virtual memory by calling `pciio_dmamap_addr()`. It takes principle arguments of a map, a memory address, and a length. The value returned is a bus address that you can program into a bus master device. When the device accesses that address, it is accessing the specified memory location.

Once you have extracted an address using `pciio_dmamap_addr()`, the map is active. It remains active until you call either `pciio_dmamap_done()` or `pciio_dmamap_free()`. In the O2 workstation it costs nothing to keep a DMA map active. In other systems, an active map may tie up global hardware resources. It is a good idea to call `pciio_dmamap_done()` when the I/O operation is complete.

In systems in which PCI space is hard-wired to specific memory addresses, `pciio_dmamap_alloc()` is a short function and `pciio_dmamap_addr()` is a trivial one. However, these systems also support a one-step translation function, `pciio_dmatrans_addr()`. This function takes a combination of the arguments of `pciio_dmamap_alloc()` and `pciio_dmamap_addr()`, and returns a translated address. In effect, it combines creating a map, using the map, and freeing the map, into a single step.

### Managing Address-Length Lists

In some cases you are not sure whether a memory buffer is contiguous, or perhaps you are sure that it is not. In this case you need to create a list of memory addresses and lengths—one address and length for each segment of memory. Then you need to translate the segments into a list of bus addresses. The list of bus addresses can be programmed into the device one at a time or, if the device supports scatter/gather, you can program all of the list of addresses for transfer in sequence.

Support for these cases is provided by *address-length lists*, an abstract data type that is supported by a family of functions. The IRIX 6.4 contains a complete family of functions for address-length lists. IRIX 6.3 supports a subset necessary to use with DMA maps. The functions related to address-length lists are summarized in Table 15-5. See reference page `alenlist(d4x)` for an overview. For function syntax see `alenlist_ops(d3x)` and `pciio_dma(d3)`.

**Table 15-5** Functions for DMA Using Address-Length Lists

Function	Purpose and Operation
<code>alenlist_destroy()</code>	Release an address-length list.
<code>alenlist_get()</code>	Retrieve the next address and length pair from a list.
<code>buf_to_alenlist()</code>	Create an address-length list to describe the buffer represented by a <code>buf_t</code> object.
<code>kvaddr_to_alenlist()</code>	Create an address-length list to describe a buffer in kernel virtual memory.
<code>pciio_dmamap_list()</code>	Convert an address-length list of memory addresses into an address-length list of corresponding bus addresses.
<code>pciio_dmatrans_list()</code>	Request immediate conversion of an address-length list of memory addresses into an address-length list of corresponding bus addresses. Returns NULL unless this system supports fixed DMA mapping.

The function **buf\_to\_alenlist()** is called in a *pfstrategy()* entry point. It takes a *buf\_t* and returns an address-length list that describes each segment of memory in the buffer that the *buf\_t* describes (see “Structure *buf\_t*” on page 185 and “Entry Point *strategy()*” on page 157). The function **kvaddr\_to\_alenlist()** takes the address and length of any buffer in kernel virtual memory and returns an address-length list to describe that extent of memory.

When you are ready to perform DMA to a buffer, you create an address-length list to describe the buffer, and pass that through **pciio\_dmamap\_list()**. This function returns a new address-length list in which the memory addresses have been replaced by PCI bus addresses.

You step through the contents of the converted address-length list using **alenlist\_get()**, which returns successive pairs of values—an address and a length—from the list. You program each pair of values into the bus master device.

## Detaching A Device

In future releases of IRIX, the *pfdetach()* entry point is called when the kernel decides to detach a PCI device. This can be caused by a hardware failure or by administrator action. In practice, it does not happen at all in IRIX 6.3 for O2. You may provide the entry point, but it is not called.

## Inactivating an Interrupt Handler

The functions for managing interrupt handlers are summarized in Table 15-6. See reference page *pciio\_intr(d3)* for syntax.

**Table 15-6** Functions for Managing PCI Interrupt Handlers

Function	Purpose and Operation
<i>pciio_intr_alloc()</i>	Create an interrupt object that enables interrupts to flow from a specified device.
<i>pciio_intr_connect()</i>	Associate an interrupt object with an interrupt handler function.
<i>pciio_intr_disconnect()</i>	Remove the association between an interrupt object and a handler function.
<i>pciio_intr_free()</i>	Release an interrupt object.

The allocation of an interrupt handler is covered under “Registering an Interrupt Handler” on page 391. When detaching a device, call `pciio_intr_disconnect()` to break the association between the interrupt and the handler function.

### Inactivating Maps and Releasing Objects

There are typically various allocated objects—PIO and DMA maps, interrupt objects—that are addressed from the device information structure stored for the device. All such objects should be released at this time.

### Unloading

When a loadable PCI driver is called at its `pfxunload()` entry point, indicating that the kernel would like to unload it, it must take great pains not to leave any dangling pointers (see “Entry Point unload()” on page 170). A driver should not unload when it has any registered interrupt handlers.

A driver does not have to unregister itself as a PCI driver before unloading. Nor does it have to detach any devices it has attached. However, if any devices are open or memory mapped, the driver should not unload.

If the kernel discovers a device and wants this driver to attach it, the kernel will reload the driver. If the driver has already attached one or more devices, the driver’s information about the state of those devices is safely stored in each hardware vertex. When a process wants to open one of the devices, the driver will be reloaded automatically, and will be able to find the device information again.

## PCI Function Summary

Table 15-7 contains a summary of the PCI-related kernel functions, in alphabetical order.

**Table 15-7** PCI-Related Kernel Functions

Function	Purpose and Operation	Discussed
alenlist_destroy() (alenlist_ops(d3x))	Release an address-length list.	page 399
alenlist_get() (alenlist_ops(d3x))	Retrieve the next address and length pair from a list.	page 399
buf_to_alenlist() (alenlist_ops(d3x))	Create an address-length list to describe the buffer represented by a <i>buf_t</i> object.	page 399
hwgraph_device_add()	Add a device vertex for a logical device.	page 393
hwgraph_device_add_minor()	Associate a logical device vertex with a minor number.	page 393
kvaddr_to_alenlist() (alenlist_ops(d3x))	Create an address-length list to describe a buffer in kernel virtual memory.	page 399
device_info_get()	Retrieve the address of device information from the hardware graph vertex.	page 395
device_info_set()	Store the address of device information in the hardware graph vertex.	page 395
pciio_config_get() (pciio_config(d3))	Fetch a defined register from configuration space using a base address returned by <b>pciio_piomap_addr()</b> .	page 390
pciio_config_set() (pciio_config(d3))	Store a value into one of the defined fields of configuration space using an address returned by <b>pciio_piomap_addr()</b> .	page 390
pciio_dmamap_addr() (pciio_dma(d3))	Get the bus virtual address corresponding to a memory address for a specified length.	page 399
pciio_dmamap_alloc() (pciio_dma(d3))	Create a DMA map object, specifying the maximum extent of memory the map will have to cover.	page 388
pciio_dmamap_done() (pciio_dma(d3))	Make a DMA map inactive (may release hardware resources associated to the map).	page 399

**Table 15-7 (continued)** PCI-Related Kernel Functions

Function	Purpose and Operation	Discussed
<code>pcii_dmamap_free()</code> ( <code>pcii_dma(d3)</code> )	Release a DMA map object.	page 399
<code>pcii_dmamap_list()</code> ( <code>pcii_dma(d3)</code> )	Convert an address-length list of memory addresses into an address-length list of corresponding bus addresses.	page 399
<code>pcii_dmatrans_addr()</code> ( <code>pcii_dma(d3)</code> )	Request immediate translation of the address of a contiguous memory buffer to a bus address. Returns NULL unless this system supports fixed DMA addressing	page 399
<code>pcii_dmatrans_list()</code> ( <code>pcii_dma(d3)</code> )	Request immediate conversion of an address-length list of memory addresses into an address-length list of corresponding bus addresses. Returns NULL unless this system supports fixed DMA mapping.	page 399
<code>pcii_driver_register()</code> ( <code>pcii(d3)</code> )	Notify the kernel that this driver is ready, and tell the vendor and device numbers it supports.	page 384
<code>pcii_driver_unregister()</code> ( <code>pcii(d3)</code> )	Notify the kernel this driver is not available (for example the driver is unloading).	
<code>pcii_intr_alloc()</code> ( <code>pcii_intr(d3)</code> )	Create an interrupt object that enables interrupts to flow from a specified device.	page 391
<code>pcii_intr_connect()</code> ( <code>pcii_intr(d3)</code> )	Associate an interrupt object with an interrupt handler function.	page 391
<code>pcii_intr_disconnect()</code> ( <code>pcii_intr(d3)</code> )	Remove the association between an interrupt object and a handler function.	
<code>pcii_intr_free()</code> ( <code>pcii_intr(d3)</code> )	Release an interrupt object.	
<code>pcii_piomap_addr()</code> ( <code>pcii_pio(d3)</code> )	Get a kernel virtual address from a PIO map for a specific offset and length.	page 396
<code>pcii_piomap_alloc()</code> ( <code>pcii_pio(d3)</code> )	Create a PIO map object, specifying the bus address space, base offset, and length it needs to cover.	page 388

**Table 15-7 (continued)** PCI-Related Kernel Functions

Function	Purpose and Operation	Discussed
pciio_piomap_done() (pciio_pio(d3))	Make a PIO map inactive until it is next needed (may release hardware resources associated to the map).	page 396
pciio_piomap_free() (pciio_pio(d3))	Release a PIO map object.	page 396
pciio_piotrans_addr() (pciio_pio(d3))	Request immediate translation of a bus address to a kernel virtual address without use of a PIO map. Returns NULL unless this system supports fixed PIO addressing.	page 396

## Example Driver

The code in Example 15-7 implements a complete, working, PCI device driver for IRIX 6.3 for O2. This same source code is also available on the SGI Developer Toolbox CDROM, where you can also find the code for the user-level program that tests it.

- Example 15-7 displays the descriptive file for */var/sysgen/master.d*.
- Example 15-8 displays the one-line VECTOR statement for */var/sysgen/system*.
- Example 15-9 displays the header file of device flags and information structure.
- Example 15-10 displays the complete source code.

### Example 15-7 Example PCI Driver for IRIX 6.3—Descriptive File

```
*
*      Barco Chameleon card
*
* Loadable driver: FLAG = fdN
* Non-loadable:   FLAG = c
*FLAG  PREFIX  SOFT  #DEV  DEPENDENCIES
cdN    coco_   73    -

```

\$\$\$

### Example 15-8 Example PCI Driver for IRIX 6.3—Configuration File

```
VECTOR: bustype=PCI module=coco
```

**Example 15-9** Example PCI Driver for IRIX 6.3—Driver Header File

```

/* =====
 *      Input/Output and Byte Swapping
 *      ===== */
#define BYTE_SWAP16(u)  (ushort_t)(((u<<8)&0xff00)|((u>>8)&0x00ff))
#define BYTE_SWAP32(u)  (uint_
t)(((u<<24)|((u<<8)&0xff0000)|((u>>8)&0xff00)|(u>>24))
/*
 * byte swap Input/Output
 */
#if 0
#define Out8(addr, b)   ( *(volatile uchar_t *) (addr) = (b) )
#define Out16(addr, s) ( *(volatile ushort_t *) (addr) = BYTE_SWAP16(s) )
#define Out32(addr, w) ( *(volatile uint_t *) (addr) = BYTE_SWAP32(w) )
#define Inp8(addr)     ( *(volatile uchar_t *) (addr) )
#define Inp16(addr)    BYTE_SWAP16( *(volatile ushort_t *) (addr) )
#define Inp32(addr)    BYTE_SWAP32( *(volatile uint_t *) (addr) )
#endif
/*
 * Input/Output with no byte swap
 */
#define Out8(addr, b)   ( *(volatile uchar_t *) (addr) = (b) )
#define Out16(addr, w) ( *(volatile ushort_t *) (addr) = (w) )
#define Out32(addr, w) ( *(volatile uint_t *) (addr) = (w) );flushbus()
#define Inp8(addr)     ( *(volatile uchar_t *) (addr) )
#define Inp16(addr)    ( *(volatile ushort_t *) (addr) )
#define Inp32(addr)    ( *(volatile uint_t *) (addr) )
/* =====
 *      Sleep/Wakeup/Lock
 *      ===== */
#define COCO_LOCK      splhi
#define COCO_UNLOCK(s) splx(s)
#define SleepEvent(x) psema(x, (PRIBIO|PCATCH) )
#define WakeEvent(x)  vsema(x)
/* =====
 *      Misc. defaults
 *      ===== */
#define DEVICE_ID      0x0001
#define VENDOR_ID      0x11a4
#define DRIVER_PREFIX  "coco_"
#define MAJOR_NUMBER   73
#define AMCC_RAM_SIZE  64
#define CONFIG_RAM_SIZE 16
#define NORMAL_DMA_RAM_SIZE 16

```

```

#define COCO_CONFIG_HDR 68
#define END_OF_CHAIN 0x80000000
#define RW_TIMER 500 /* wait for read/write in clock ticks */
#define SIMRW_TIMER 500 /* wait for sim R/W in clock ticks */
#define CHAIN_FACTOR 10
#define MAPPED_SIZE 17*NBPP
#define COCO_CACHE_SIZE 32
/* =====
 * Configuration Register bits
 * ===== */
#define CONFIG_CCRES 0x00000001 /* reset when zero */
#define CONFIG_FRES 0x00000002 /* Input/Output Fifo reset when 0 */
#define CONFIG_FSDATI 0x00000004 /* Inp. Fifo serial config data */
#define CONFIG_FSDATO 0x00000008 /* Out. Fifo serial config data */
#define CONFIG_FSCLK 0x00000010 /* In/Out Fifo serial config clock */
#define CONFIG_LUTSEL 0x00000020 /* External LUT bank selection (=0) */
/* bits 6-9 is RAM address */
#define CONFIG_DMA_READ_ADDR 0x00000000 /* DMA Read Address */
#define CONFIG_DMA_WRITE_ADDR 0x00000040 /* DMA Write Address */
#define CONFIG_SC_READ_ADDR 0x00000080 /* Scatter-Gather Read Address */
#define CONFIG_SC_WRITE_ADDR 0x00000000 /* Scatter-Gather Write Address */
/* =====
 * Mode Register
 * ===== */
#define COCO_MODE 0x01
#define COCO_SWAP 0x02
#define COCO_SLICE 0x04
#define COCO_DELAY 0x08
#define COCO_FLAG 0x20
#define COCO_WRENA 0x40
/* =====
 * Configuration Register
 * ===== */
#define DMAREG_CCRES 0x00000001L /* Chameleon reset when zero */
#define DMAREG_HICRES 0x00000001L /* Chameleon reset when zero */
#define DMAREG_FRES 0x00000002L /* I/O Fifo reset when zero */
#define DMAREG_FSDATI 0x00000004L /* Input Fifo serial config. data */
#define DMAREG_FSDATO 0x00000008L /* Output Fifo serial config. data */
#define DMAREG_FSCLK 0x00000010L /* I/O serial config clock */
#define DMAREG_LED 0x00000020L /* external LUT bank selection (=0) */
/* bit 6-9: Ram address */
#define DMAREG_RAMRADR 0x00000000L /* DMA read address */
#define DMAREG_RAMWADR 0x00000040L /* DMA write address */
#define DMAREG_RAMPRRD 0x00000080L /* scatter-gather read address */
#define DMAREG_RAMPRWR 0x000000C0L /* scatter-gather write address */

```

```

#define DMAREG_RAMCNT 0x00000100L /* read count copy */
#define DMAREG_RAMWCNT 0x00000140L /* write count copy */
#define DMAREG_FILL 0x00000400L /* enable automatic ram fill for DMA */
#define DMAREG_PTEN 0x00008000L /* Fifo Interface Enable */
#define DMAREG_INTREN 0x00010000L /* DMA Read Interrupt Enable */
#define DMAREG_INTWEN 0x00020000L /* DMA Write Interrupt Enable */
#define DMAREG_INTPREN 0x00040000L /* data-chained DMA read int. enable */
#define DMAREG_INTPWEN 0x00080000L /* data-chained DMA write int. enable */
/* bit 20-21: Device Selection */
#define DMAREG_RAM 0x00000000L /* RAM */
#define DMAREG_RCNT 0x00100000L /* DMA Read Count */
#define DMAREG_WCNT 0x00200000L /* DMA Write Count */
#define DMAREG_STAT 0x00300000L /* Status/Fifo control Register */
#define DMAREG_REN 0x10000000L /* DMA Read Enable */
#define DMAREG_WEN 0x20000000L /* DMA Write Enable */
#define DMAREG_PREN 0x40000000L /* data-chained DMA read enable */
#define DMAREG_PWEN 0x80000000L /* data-chained DMA write enable */
#define DMAREG_NVIFEN 0x00800000L /* Mailbox Interface Enable */
#define DMAREG_DMACVT 0x00400000L /* Unused */
/* Unknown !! */
#define DMAREG_HID0IT 0x00000400L
#define DMAREG_HISCLK 0x00000800L
#define DMAREG_HISDI 0x00001000L
#define DMAREG_HISDO 0x00000080L
#define EOFPROG 0xF0000000L
/* =====
* AMCC Register Offsets
* ===== */
#define AMCC_OP_REG_OMB1 0x00
#define AMCC_OP_REG_OMB2 0x04
#define AMCC_OP_REG_OMB3 0x08
#define AMCC_OP_REG_OMB4 0x0c
#define AMCC_OP_REG_IMB1 0x11
#define AMCC_OP_REG_IMB2 0x14
#define AMCC_OP_REG_IMB3 0x18
#define AMCC_OP_REG_IMB4 0x1c
#define AMCC_OP_REG_FIFO 0x20
#define AMCC_OP_REG_MWAR 0x24
#define AMCC_OP_REG_MWTC 0x28
#define AMCC_OP_REG_MRAR 0x2c
#define AMCC_OP_REG_MRTC 0x30
#define AMCC_OP_REG_MBEF 0x34
#define AMCC_OP_REG_INTCSR 0x38
#define AMCC_OP_REG_MCSR 0x3c
#define AMCC_OP_REG_MCSR_NVDATA (AMCC_OP_REG_MCSR + 2) /* Data in byte 2 */

```

```

#define AMCC_OP_REG_MCSR_NVCMD (AMCC_OP_REG_MCSR + 3) /* Command in byte 3 */
/* Amcc INTCSR interrupt bits */
#define AMCC_INTCSR_WEN 0x00004000
#define AMCC_INTCSR_REN 0x00008000
#define AMCC_INTCSR_INTMB4 0x00001f00
/* enable Output Mbox4, byte 3 only */
#define AMCC_INTCSR_MASK AMCC_INTCSR_INTMB4
#if 0 /* Write/Read Completion Interrupt enable only */
#define AMCC_INTCSR_MASK AMCC_INTCSR_WEN|AMCC_INTCSR_REN
/** "Write/Read Completion Interrupt" with Output Maibox4 */
#define AMCC_INTCSR_MASK AMCC_INTCSR_WEN|AMCC_INTCSR_REN|AMCC_INTCSR_INTMB4
#endif
#define AMCC_INTCSR_RCLR 0x00080000 /* Read Transfer Complete Clear */
#define AMCC_INTCSR_WCLR 0x00040000 /* Write Transfer Complete Clear */
#define AMCC_INTCSR_RST 0x00330000 /* Target/Master Abort and Out Mbox */
/* Amcc MCSR bits */
#define AMCC_REN 0x00004000 /* Read Enable */
#define AMCC_WEN 0x00000400 /* Write Enable */
#define AMCC_RFMAN 0x00002000L
#define AMCC_WFMAN 0x00000200L
#define AMCC_RFPRI 0x00001000L /* Read Priority over Write */
#define AMCC_WFPRI 0x00000100L /* Write Priority over Read */
#define AMCC_RST_FIFOS 0x06000000L /* Reset Fifos */
#define AMCC_RST_ADDON 0x01000000L /* Reset Add-on */
#define AMCC_MCSR_MASK AMCC_RST_FIFOS | AMCC_RST_ADDON
/* Outgoing Mailbox Register 4, byte 3 */
#define AMCC_MB_EOFDMAR 0x01000000 /* 1 = end of read DMA */
#define AMCC_MB_EOFDMAW 0x02000000 /* 1 = end of write DMA */
#define AMCC_MB_EOFPRDMAR 0x04000000 /* 1 = end of prog. read DMA */
#define AMCC_MB_EOFPRDMAW 0x08000000 /* 1 = end of prog. write DMA */
#define AMCC_MB_DIAGN 0x10000000 /* Diagnostic flag */
#define AMCC_MB_COCOMODE 0x20000000 /* Chameleon flag */
#define AMCC_MB_FIFORSTR 0x40000000 /* AMCC input Fifo Reset */
#define AMCC_MB_FIFORSTW 0x80000000 /* AMCC output Fifo Reset */
/* =====
* Device Information
* =====
*
* status: Shows whether the driver is attached/opened, etc.
* dmacfg: The boards Configuration default setting.
* dmabits: DMA operation bit setting (in addition to dmacfg)
* dmatype: Single or Chain DMA type
* dmastat: currennt status of DMA (Idle, Wait, etc.)
* dmacmd: Board's DMA command (Transp or Convert)
* dmawait: Semaphore for wait/wakeup

```

```
*      wp_addr:      Residual phys. address for Write DMA.
*      rp_addr:      Residual phys. address for Read  DMA.
*      rp_size:      Residual size in bytes for Read  DMA.
*      wp_size:      Residual size in bytes for Write DMA.
*      bp:           Current buf_t pointer for Read/Write DMA
*      chain_list:   Address of current chain list buffer
*      addrList:     Current DMA scatter/gather structure (single read/write)
*      r_addrList:   Sim. Read/Write Read scatter/gather structure.
*      w_addrList:   Sim. Read/Write Write scatter/gather structure.
*      page_no:     Pages left to DMA (single read/write)
*      r_page_no:    Pages left for Read to DMA (Sim. read/write).
*      w_page_no:    Pages left for Write to DMA (Sim. read/write).
*      cfg_adr:      Board's Configuration Area Address
*      amcc_adr:     Amcc PCI address
*      conf_adr:     Board's Config Register address
*      norm_adr:     Board's Normal DMA Registers address
*      start_time:   Time DMA started
*      intr_time:    Time board Interrupt completion of DMA
*      call_time:    Time Call came into the driver
*      ret_time:     Time driver returns to the user
*      vhdl:         Vertex handle representing this board.
*      dev_intr:     Our Interrupt Handler structure
*      dev_desc:     Our Device Descriptor structure
*      tid:          Timer ID ..timer waiting for board to interrupt.
*/
typedef struct {
    int          status;
    /* dma stuff */
    uint_t      dmacfg;
    uint_t      dmabits;
    int         dmatype;
    int         dmastat;
    uint_t      dmacmd;
    sema_t      dmawait;
    int         iostat;
    int         dmasize;
    alenaddr_t  wp_addr;
    alenaddr_t  rp_addr;
    size_t      rp_size;
    size_t      wp_size;
    buf_t       *bp;
    caddr_t     chain_list;
    alenlist_t  addrList;
    alenlist_t  r_addrList;
    alenlist_t  w_addrList;
}
```

```
int          page_no;
int          r_page_no;
int          w_page_no;
/* mapped memory */
vhandl_t    *vhandl;
caddr_t     mappedkv;
int         mappedkvlen;
/* addresses */
caddr_t     cfg_adr;
caddr_t     amcc_adr;
caddr_t     conf_adr;
caddr_t     norm_adr;
/* for Dma time measurement */
struct timeval start_time;
struct timeval intr_time;
struct timeval call_time;
struct timeval ret_time;
/* Irix interface structs */
vertex_hdl_t vhdl;
pciio_intr_t dev_intr;
device_desc_t dev_desc;
toid_t      tid;
} card_t;
/* bits for status */
#define CARD_ATTACHED 0x01
#define CARD_OPEN    0x02
/* dmatype values */
#define DMA_PROG      0 /* chained - default value */
#define DMA_SINGLE    1
/* dmastat values */
#define DMA_IDLE      0
#define DMA_LUT_WAIT  1
#define DMA_READ_WAIT 2
#define DMA_WRITE_WAIT 3
#define DMA_RW_WAIT   4
/* iostat values */
#define IO_OK         0
#define IO_ERROR      1
#define IO_TIME       2
/* chained DMA block */
typedef struct {
    paddr_t nextaddr;
    paddr_t addr;
    int     size;
} coco_dmapage_t;
```

**Example 15-10** Example PCI Driver for IRIX 6.3—Driver Source Code

```

/*****
*****      C h a m e l e o n   I r i x   P c i   D r i v e r      *****/
*****
*****/
#include <sys/types.h>
#include "sys/cmn_err.h"
#include "sys/sem.h"
#include <sys/param.h>
#include <sys/errno.h>
#include <sys/syslog.h>
#include <sys/conf.h>
#include <sys/pio.h>
#include "sys/system.h"
#include <sys/time.h>
#include <sys/kmem.h>
#include <sys/ktime.h>
#include <sys/mload.h>
#include <sys/ddi.h>
#include <sys/cred.h>
#include <sys/user.h>
#include <sys/mace.h>
#include <sys/immu.h>
#include <sys/region.h>
#include <sys/alenlist.h>
#include <sys/IP32.h>
#include <sys/PCI/PCI_defs.h>
#include <sys/PCI/pciio.h>
#include "coco.h"
#include "coco_user.h"
char   *coco_mversion = M_VERSION; /* loadable driver requirement */
int     coco_devflag = 0;          /* ddi/dki requirement      */
/* =====
*   Device Driver/PCI entry routines
* ===== */
int coco_unload(void);
int coco_open(dev_t *, int, int, cred_t *);
int coco_close(dev_t, int, int, cred_t *);
int coco_read( dev_t, uio_t *, cred_t *);
int coco_write( dev_t, uio_t *, cred_t *);
int coco_ioctl(dev_t, int, void *, int, cred_t *, int *);
int coco_map ( dev_t, vhandl_t *, off_t, int, int );
int coco_unmap ( dev_t, vhandl_t * );
int coco_init();

```

```
int coco_attach(vertex_hdl_t);
int coco_detach(vertex_hdl_t);
int coco_error(vertex_hdl_t, int );
void coco_intr( eframe_t *, intr_arg_t );
/* =====
 * Supporting internal routines
 * ===== */
static void cocoReset( card_t * );
static void cocoProgFlags(caddr_t, uint_t, ushort_t, ushort_t, ushort_t, ushort_t);
static void cocoSetMode ( card_t *, int, int, int, int );
static void cocoCommand ( card_t *, uint_t, uint_t );
static void cocoBufOut ( card_t *, uint_t *, int );
static void cocoBufIn ( card_t *, uint_t *, int );
static void cocoReadDmaRegs ( card_t *, uint_t * );
static void cocoWriteDmaRegs ( card_t *, uint_t * );
static void cocoSetAddr ( card_t *, uint_t );
static void cocoReadIntRam ( card_t *, uint_t *, int, int );
static void cocoWriteIntRam ( card_t *, uint_t *, int, int );
static void cocoReadExtRam ( card_t *, uint_t *, int );
static void cocoWriteExtRam ( card_t *, uint_t *, int );
static void cocoConvert ( card_t *, uint_t );
static void cocoConvertTest ( card_t *, coco_convert_t *);
static void cocoPrepAmcc ( card_t * );
static void cocoResetAmcc ( card_t *);
static void cocoStartProgDma ( card_t *, iopaddr_t, int, int );
static void cocoStartSingleDma ( card_t *, alenaddr_t, size_t, int );
static void cocoReport ( card_t *, char * );
static void cocoShowChain( char *, coco_dmapage_t * );
static void cocoTimeOut( card_t *);
static void cocoTimeOut2( card_t *);
static void cocoDumpAmcc ( card_t *);
static void cocoReportTime ( caddr_t, card_t *, int );
static void cocoShowAlenlist ( caddr_t, alenlist_t );
static void cocoUnlockUser ( caddr_t, int, int );
static void cocoDiffTime( struct timeval *,struct timeval *,struct timeval *);
static int cocoStrategy ( buf_t * );
static int cocoReadAmccFifo ( card_t * );
static int cocoFifoTest ( card_t *, int );
static int cocoPattern ( int, int );
static int cocoReadMode ( card_t * );
static int cocoReadAddr (card_t *);
static int cocoDmaRegsTest ( card_t * );
static int cocoIntRamTest ( card_t * );
static int cocoExtRamTest ( card_t * );
static int cocoDmaToLuts( card_t *, coco_buf_t *, int );
```

```

static int  cocoReadWrite( card_t *, coco_rw_t * );
static int  cocoStartRWDma ( card_t *, iopaddr_t, int, iopaddr_t, int );
static int  cocoMakeChainRW ( card_t *, coco_dmapage_t **, coco_dmapage_t **);
static int  cocoAlenlistSize ( alenlist_t );
static int  cocoLockUser ( caddr_t, int, int );
static coco_dmapage_t * cocoMakeChain ( card_t *, alenlist_t, int );
static char *cocoIoctlStr(int);
/*
 * temporary declaration of buf_to_alenlist
 * - missing from pciio.h & alenlist.h
 * - also not compatible with IRIX 6.4!
 */
extern alenlist_t buf_to_alenlist(buf_t *);
#define COCO_TEST 0x999
static void  cocoDebug ( coco_rw_t * );
/*
~~~~~
~~~~~
~~~~~      D r i v e r ' s   E n t r y   R o u t i n e s   ~~~~~
~~~~~
~~~~~
*/
/*****
***              c o c o _   i n i t              ***
*****
*
* Name:          coco_init
*
* Purpose:       Called by kernel. Here we declare our PCI interface
*                routines attach, detach and error and register our
*                driver to take care of our card (card is identified by
*                Vendor and Device IDs).
*
* Returns:       None
*
*****/
int
coco_init()
{
    register int ret;
    printf("coco_init()\n");
#ifdef _EARLY_PCI
    /* Identify our attach, detach and error routines */
    ret = pciio_add_attach(coco_attach, coco_detach,
                          (pciio_error_handler_f *)coco_error,

```

```

                                DRIVER_PREFIX, MAJOR_NUMBER);
    if ( ret != 0 ) {
        printf ("coco_init: could not add_attach\n");
        return(0);
    }
#endif
/*   Register and identify the card   */
ret = pciio_driver_register(VENDOR_ID, DEVICE_ID, DRIVER_PREFIX, 0);
if ( ret != 0 ) {
    printf ( "coco_init: could not register\n");
    return(0);
}
}
/*****
***                               c o c o _ a t t a c h                               ***
*****/
*
* Name:          coco_attach
*
* Purpose:      Called by the kernel. Our card is installed and hence
*               we are called. Prepare everything necessary to handle
*               the card. Note that when the card is found, the following
*               is set in the Command field of Config space:
*               - Bus master enabled.
*               - Memory and IO space access enabled.
*               - Cache line is set to 0x20 (32)
*               - Latency Timer is set to 0x30 (48)
*
*               For Chameleon, beside Configuration address space, we need 4
*               Memory address spaces to be mapped:
*               Base_Reg 0 = AMCC Registers and Fifos
*               Base_Reg 3 = Normal DMA registers
*               Base_Reg 4 = Configuration Register.
*
* Returns:      0 for Success, or errno
*
*****/
int
coco_attach(vertex_hdl_t vhdl)
{
    card_t      *cp;
    caddr_t     cfg_adr, mem_ptr, amcc_adr, conf_adr, norm_adr;
    int         ret;
    u_int32_t   vendor_id, device_id, base_reg, cmd_reg;
    device_desc_t dev_desc;

```

```

#ifdef DEBUG
printf ("coco_attach: ---- start -----\n");
#endif
/* =====
 * Configuration Space
 * ===== */
/* Get a pointer to the card's Configuration address space */
cfg_adr = (caddr_t)pciio_piortrans_addr ( vhdl, NULL, PCIIO_SPACE_CFG,
                                         (iopaddr_t)0, COCO_CONFIG_HDR, 0);

if ( cfg_adr == (caddr_t)NULL ) {
    cmn_err ( CE_WARN, "coco_attach: Cannot get to Config space");
    return(EIO);
}
#ifdef DEBUG
printf ("Config. address is 0x%x\n", cfg_adr );
/*
 * Get Vendor_Id, Device_Id and Base_Reg and print them
 * Here we get only Base_Reg 0. Get others if your card uses more.
 * Beside these general fields, get any vendor specific fields
 * and check/print them for debugging purpose.
 */
vendor_id    = pciio_config_get(cfg_adr, PCI_CFG_VENDOR_ID);
device_id    = pciio_config_get(cfg_adr, PCI_CFG_DEVICE_ID);
cmd_reg      = pciio_config_get(cfg_adr, PCI_CFG_COMMAND );
printf ("Coco Vendor_Id = 0x%x, Device_Id = 0x%x, Cmd = 0x%x\n",
        vendor_id, device_id, cmd_reg );
#endif
/* =====
 * Memory Address Space
 * ===== */
/* Get Amcc Register addresses */
amcc_adr = (caddr_t)pciio_piortrans_addr ( vhdl, NULL, PCIIO_SPACE_WIN(0),
                                         (iopaddr_t)0, AMCC_RAM_SIZE, 0);

#ifdef DEBUG
printf ("Amcc address is 0x%x\n", amcc_adr );
#endif
if ( amcc_adr == (caddr_t)NULL ) {
    cmn_err(CE_WARN, "coco_attach: Cannot get to AMCC address space");
    return (EIO);
}
/* Get DMA Config Register address */
conf_adr = (caddr_t)pciio_piortrans_addr ( vhdl, NULL, PCIIO_SPACE_WIN(4),
                                         (iopaddr_t)0, CONFIG_RAM_SIZE, 0);

#ifdef DEBUG
printf ("Config Register Address is 0x%x\n", conf_adr );

```

```

#endif
if ( conf_adr == (caddr_t)NULL ) {
    cmn_err(CE_WARN, "coco_attach: Cannot get to Config Register");
    return (EIO);
}
/* Get Normal DMA Register addresses */
norm_adr = (caddr_t)pciio_piortrans_addr ( vhdl, NULL, PCIIO_SPACE_WIN(3),
                                          (iopaddr_t)0, NORMAL_DMA_RAM_SIZE,
                                          0);

#ifdef DEBUG
printf ("Normal DMA address is 0x%x\n", norm_adr);
#endif
if ( norm_adr == (caddr_t)NULL ) {
    cmn_err(CE_WARN, "coco_attach: Cannot get to Normal DMA address");
    return (EIO);
}
/* allocate an internal structure for this card and save everything */
cp = (card_t *)kmem_zalloc ( sizeof(card_t), KM_NOSLEEP );
if ( cp == (card_t *)NULL ) {
    cmn_err(CE_WARN, "coco_attach: Cannot allocate memory");
    return(ENOMEM);
}
cp->vhdl      = vhdl;
cp->cfg_adr   = cfg_adr;
cp->amcc_adr  = amcc_adr;
cp->conf_adr  = conf_adr;
cp->norm_adr  = norm_adr;
/* =====
 *      Interrupt Handler Registration
 * ===== */
dev_desc = kmem_alloc(sizeof (*dev_desc), KM_SLEEP);
dev_desc->intr_swlevel = COCO_LOCK;
cp->dev_intr = pciio_intr_alloc ( vhdl, dev_desc, PCIIO_INTR_LINE_A, vhdl );
if (cp->dev_intr == (pciio_intr_t)NULL){
    cmn_err(CE_WARN, "coco_attach: Can't pciio_intr_alloc");
    kmem_free (dev_desc, sizeof(*dev_desc) );
    kmem_free ( cp, sizeof(card_t) );
    return(EIO);
}
ret = pciio_intr_connect ( cp->dev_intr, (intr_func_t)coco_intr,
                          (intr_arg_t)cp, 0 );
if ( ret == -1 ) {
    cmn_err(CE_WARN, "coco_attach: Cannot register interrupt handler");
    kmem_free (dev_desc, sizeof(*dev_desc) );
    kmem_free ( cp, sizeof(card_t) );
}

```

```

        return (EIO);
    }
    cp->dev_desc = dev_desc;
    cp->status = CARD_ATTACHED;
    /* allocate memory for mapping */
    cp->mappedkv = kmem_alloc (MAPPED_SIZE,
                              KM_NOSLEEP | KM_PHYSCONTIG | KM_CACHEALIGN);
    if ( cp->mappedkv == (caddr_t)NULL ) {
        cmn_err (CE_NOTE, "coco_attach: Not enough memory for %d mapping",
                 MAPPED_SIZE);
        cmn_err (CE_WARN, "coco_attach: No mapping is allowed");
    }
    else {
        cp->mappedkvlen = MAPPED_SIZE;
        cmn_err (CE_NOTE, "coco_attach: %d bytes available for mapping",
                 cp->mappedkvlen );
    }
    /* save our structure */
    device_info_set ( vhdl, (arbitrary_info_t)cp );
    cmn_err ( CE_NOTE, "coco_attach: driver is ready");
    return(0);
}
/*****
***          c o c o _ d e t a c h          ***
*****/
*
* Name:          coco_detach
*
* Purpose:       Detaches a driver
*
* Returns:       0 Success or errno
*
*****/
int
coco_detach(vertex_hdl_t vhdl)
{
    register card_t      *cp;
    cmn_err(CE_NOTE,
            "coco_detach: Unregister Interrupt handler and free up mem");
    cp = (card_t *)device_info_get ( vhdl );
    /* free up memory for mapping */
    if ( cp->mappedkv )
        kmem_free ( cp->mappedkv, cp->mappedkvlen );
    /* Unregister the Interrupt Handler */
    pciio_intr_disconnect(cp->dev_intr);
}

```

```

    pcio_intr_free(cp->dev_intr);
    kmem_free ( cp->dev_desc, sizeof(*cp->dev_desc) );
    kmem_free ( cp, sizeof(card_t) );
    return(0);
}
/*****
***          c o c o _ i n t r          ***
*****/
*
* Name:      coco_intr
*
* Purpose:   Our interrupt handler.
*
* Returns:   None.
*
*****/
void
coco_intr( eframe_t *ef, intr_arg_t arg )
{
    register card_t  *cp;
    register buf_t   *bp;
    register caddr_t adr_cfg, adr_amcc, adr_norm;
    register uint_t  intcsr, xil_stat, mbox;
    register uint_t  dmatstat, dmatype, dmacfg, dmabits;
    register int     rw, err;
    size_t          p_size;
    alenaddr_t      p_addr;
    /* is it stray interrupt ? */
    cp = (card_t *)arg;
    microtime ( & cp->intr_time)
    adr_amcc = cp->amcc_adr;
    intcsr = Inp32(adr_amcc+AMCC_OP_REG_INTCSR);
    if ( (intcsr & 0x00800000) == 0 ) {
        #ifdef DEBUG
            printf ("cocoIntr: Stray interrupt\n");
        #endif
        return;
    }
    #ifdef DEBUGTIME
        cocoReportTime ( "cocoIntr", cp, 0 );
    #endif
    bp = cp->bp;
    adr_cfg = cp->conf_adr;
    adr_norm = cp->norm_adr;
    dmacfg = cp->dmacfg;

```

```

dmabits = cp->dmabits;
dmastat = cp->dmastat;
dmatype = cp->dmatype;
#ifdef DEBUG
printf ("cocoIntr: started, tid = %d\n", cp->tid );
cocoDumpAmcc(cp);
#endif
/* cancel any outstanding timer */
if ( cp->tid > 0 ) {
    untimeout(cp->tid);
    cp->tid = 0;
}
/* =====
*      Reseting Interrupt and Status
* ===== */
/* disable any Dma and read in Xilinx status */
Out32(adr_cfg, dmacfg | DMAREG_STAT );
xil_stat = Inp32(adr_norm);
mbox     = Inp32(adr_amcc+AMCC_OP_REG_IMB4);
#ifdef DEBUG
printf ("coco_intr: amcc: 0x%x, xilinx: 0x%x, mbox4 = 0x%x\n",
        intcsr, xil_stat, mbox );
#endif
/* Reset Amcc Interrupts and enable interrupts again */
Out32(adr_amcc+AMCC_OP_REG_INTCSR, AMCC_INTCSR_RST | AMCC_INTCSR_RCLR |
      AMCC_INTCSR_WCLR );
Out32(adr_amcc+AMCC_OP_REG_INTCSR, AMCC_INTCSR_MASK);
/* =====
*      End of Dma Read or Write
* ===== */
if ( (mbox & AMCC_MB_EOFDMAR) || (mbox & AMCC_MB_EOFDMAW) ) {
    dmabits &= ~(DMAREG_REN | DMAREG_INTREN);
    dmabits &= ~(DMAREG_WEN | DMAREG_INTWEN );
    #ifdef DEBUG
    if ( mbox & AMCC_MB_EOFDMAR)
        printf ("cocoIntr: End of DMA Read\n");
    else printf ("cocoIntr: End of DMA Write\n");
    #endif
}
/* =====
*      End of Chain Dma Read or Write
* ===== */
if ( (mbox & AMCC_MB_EOFPRDMAR) || (mbox & AMCC_MB_EOFPRDMAW) ) {
    dmabits &= ~(DMAREG_PREN | DMAREG_INTPREN);
    dmabits &= ~(DMAREG_PWEN | DMAREG_INTPWEN);
}

```

```

dmabits &= ~(DMAREG_WEN | DMAREG_REN);
#ifdef DEBUG
if ( mbox & AMCC_MB_EOFPRDMAR )
    printf ("cocoIntr: End Prog DMA Read\n");
else    printf ("cocoIntr: End Prog DMA Write\n");
#endif
}
cp->dmabits = dmabits;
cp->iostat = IO_OK;
switch ( cp->dmastat ) {
    /* =====
    *      Dma to Lut
    * ===== */
case DMA_LUT_WAIT:
    #ifdef DEBUG
    printf ("cocoIntr: Waking up Dma_Lut_Wait\n");
    #endif
    cp->dmastat = DMA_IDLE;
    WakeEvent(&cp->dmawait);
    goto get_out;
    /* =====
    *  Read/Write Dma
    * ===== */
case DMA_READ_WAIT:
case DMA_WRITE_WAIT:
    /* what we do depends on which type of Dma is done */
    switch ( cp->dmatype ) {
        /* chained Dma is done. Simply wake the process up */
        case DMA_PROG:
            #ifdef DEBUG
            printf ("cocoIntr: biodone() read/write\n");
            #endif
            kmem_free ( cp->chain_list,
                        cp->page_no * sizeof(coco_dmapage_t) );
            alenlist_done(cp->addrList);
            cp->addrList = 0;
            cp->dmastat = DMA_IDLE;
            bp->b_resid -= cp->dmasize;
            biodone(cp->bp);
            goto get_out;
        /* single page Dma done. Dma the next page if any */
        case DMA_SINGLE:
            bp->b_resid -= cp->dmasize;
            cp->page_no--;
            if ( cp->page_no <= 0 ) { /* no more pages */

```

```

        #ifdef DEBUG
        printf ("cocoIntr: No more pages to Dma\n");
        printf ("cocoIntr: biodone() read/write\n");
        #endif
        alenlist_done(cp->addrList);
        cp->addrList = 0;
        cp->dmastat = DMA_IDLE;
        biodone(cp->bp);
        goto get_out;
    }
    /* get next page to DMA */
    #ifdef DEBUG
    printf ("cocoIntr: Dma next page ..left = %d\n",
        cp->page_no-1 );
    #endif
    if ( alenlist_get(cp->addrList, NULL, NBPP,
        &p_addr, &p_size) != ALENLIST_SUCCESS ) {
        cmn_err (CE_WARN,
            "cocoIntr: Bad scatter-gather");
        cp->iostat = IO_ERROR;
        #ifdef DEBUG
        printf ("cocoIntr: biodone() read/write\n");
        #endif
        alenlist_done(cp->addrList);
        cp->addrList = 0;
        cp->dmastat = DMA_IDLE;
        bioerror (cp->bp, EIO);
        biodone(cp->bp);
        goto get_out;
    }
    if ( cp->dmastat == DMA_READ_WAIT )
        rw = B_READ;
    else rw = B_WRITE;
    cocoStartSingleDma ( cp, p_addr, p_size, rw );
    goto get_out;
} /*** switch ( cp->dmatype ) ***/
/* =====
 * Simultaneous read/write
 * ===== */
case DMA_RW_WAIT:
    /* what we do depends on which type of Dma is done */
    switch ( cp->dmatype ) {
        /* chained Dma is done. Both read/write is done */
        case DMA_PROG:
            #ifdef DEBUG

```



```
* Name:      coco_unload
*
* Purpose:   Unloads the driver.
*
* Returns:   0 Success or errno
*
*****/
int
coco_unload(void)
{
    printf ("coco_unload: Unloading the driver\n");
    pciio_driver_unregister(DRIVER_PREFIX);
    return(0);
}
/*****
***          c o c o _ o p e n          ***
*****/
*
* Name:      coco_open
*
* Purpose:   Opens the card. This sample driver simply verifies that
*           the card's info structure can be retrieved and checks
*           the card's Base_Register.
*
* Returns:   0 = Success, or errno.
*
*****/
int
coco_open(dev_t *devp, int flag, int otyp, cred_t *cred)
{
    register card_t      *cp;
    register vertex_hdl_t vhdl;
#ifdef DEBUG
    printf ("coco_open: Opening the card\n");
#endif
    /* Get the vertex handle and pointer to card's info */
    vhdl = dev_to_vhdl ( *devp );
    if (vhdl == NULL){
        cmn_err(CE_WARN, "coco_open: dev_to_vhdl returns NULL");
        return(EIO);
    }
    cp = (card_t *)device_info_get ( vhdl );
    /* some error checking first */
    if ( !(cp->status & CARD_ATTACHED) ) {
        cmn_err (CE_WARN, "coco_open: Driver is not attached");
    }
}
```

```

        return (ENODEV);
    }
    if ( cp->status & CARD_OPEN) {
        cmn_err (CE_WARN, "coco_open: Device is busy");
        return (EBUSY);
    }
    cocoReset(cp); /* reset the board */
    cp->status |= CARD_OPEN;
    /* default values */
    cp->dmatype = DMA_SINGLE;
    cp->dmacmd = COCO_TRANSP;
    return(0);
}
/*****
***                c o c o _ c l o s e                ***
*****/
*
* Name:          coco_close
*
* Purpose:       Closes the card. This sample driver's close statement
*                prints out the addresses and Base_Register's value for
*                verification.
*
* Returns:       0 = Success, or errno
*
*****/
int
coco_close(dev_t dev, int flag, int otyp, cred_t *cred)
{
    register card_t      *cp;
    register vertex_hdl_t vhdl;
#ifdef DEBUG
    printf ("coco_close: Closing the card\n");
#endif
    /* get the vertex handle and pointer to card's info */
    vhdl = dev_to_vhdl ( dev );
    cp = (card_t *)device_info_get ( vhdl );
    cp->status &= ~CARD_OPEN;
    return(0);
}
/*****
***                c o c o _ m a p                ***
*****/
*
* Name:          coco_map

```

```

*
* Purpose:   Allocate a piece of continuous memory and map it to user's
*            address space.
*
* Returns:   0 = Success, or errno
*
*****/
int
coco_map ( dev_t dev, vhandle_t *vh, off_t offset, int len, int prot )
{
    register int          ret;
    register caddr_t      kv;
    register vertex_hdl_t vhdl;
    register card_t       *cp;
    /* get the vertex handle and pointer to card's info */
    vhdl = dev_to_vhdl ( dev );
    cp = (card_t *)device_info_get ( vhdl );
    if ( cp->mappedkv == (caddr_t)NULL ) {
        cmn_err (CE_NOTE, "coco_map: Not memory for mapping");
        return (ENOMEM);
    }
    if ( len > cp->mappedkvlen ) {
        cmn_err (CE_NOTE,
            "coco_map: Only %d bytes available for map, requested %d bytes",
                cp->mappedkvlen, len );
        return (ENOMEM);
    }
    ret = v_maphys ( vh, cp->mappedkv, len );
    if ( ret > 0 ) {
        cmn_err (CE_WARN, "coco_map: Could not map, ret = %d", ret );
        return (ret);
    }
    /* save for later */
    cp->vhandle = vh;
    dki_dcache_inval ( cp->mappedkv, cp->mappedkvlen );
    printf ("coco_map: v_maphys returned %d, mapped 0x%x for %d bytes\n",
        ret, cp->mappedkv, len );
    return (0);
}
/*****
***                c o c o _ u n m a p                ***
*****
*
* Name:           coco_unmap
*

```

```

* Purpose:    Unmap the kernel buffer we allocated before.
*
* Returns:    0 = Success, or errno
*
*****/
int
coco_unmap ( dev_t dev, vhandle_t *vh )
{
    return (0);
}
/*****
***          c o c o _ i o c t l          ***
*****/
*
* Name:      coco_ioctl
*
* Purpose:   Handles user Ioctl command. These commands can be
*            found in coco_user.h.
*
* Returns:   0 = Success, or errno
*
*****/
int
coco_ioctl(dev_t dev, int cmd, void *arg, int mode, cred_t *cred,
           int *rvalp )
{
    register card_t      *cp;
    register vertex_hdl_t vhd1;
    register uint_t      *tmp_ibuf;
    register int         tot_bytes, err, i;
    uint_t               tmp_int;
    coco_buf_t           coco_buf;
    coco_rw_t            coco_rw;
    coco_mode_t          coco_mode;
    coco_cmd_t           coco_cmd;
    coco_convert_t       coco_convert;
    uint_t               dmaRegs[DMA_REGS];
#ifdef DEBUG
    printf ("\ncoco_ioctl: ---> command = %s [0x%x]\n",
            cocoIoctlStr(cmd), cmd );
#endif
    /* get the vertex handle and pointer to card's info */
    vhd1 = dev_to_vhdl ( dev );
    cp = (card_t *)device_info_get ( vhd1 );
    /* is device opened ? */

```

```
if (!( cp->status & CARD_OPEN ) ) {
    cmn_err (CE_NOTE, "coco_ioctl: Device is not opened");
    return (EIO);
}
bzero ( &cp->start_time, sizeof(struct timeval) );
bzero ( &cp->intr_time, sizeof(struct timeval) );
bzero ( &cp->call_time, sizeof(struct timeval) );
bzero ( &cp->ret_time, sizeof(struct timeval) );
/* =====
 *   Ioctl commands
 * ===== */
switch ( cmd ) {
    case COCO_TEST:
        /* read in coco_rw_t struct..all we need is there */
        if ( copyin(arg, (char *)&coco_rw, sizeof(coco_rw_t) ) )
            return (EFAULT);
        cocoDebug ( &coco_rw );
        return(0);
    /* =====
     *   Board Identification
     * ===== */
    case COCO_ISPCI:
        *rvalp = 1;
        break;
    /* =====
     *   Reset the board
     * ===== */
    case COCO_RESET:
        cocoReset(cp);
        break;
    /* =====
     *   Set DMA Command
     * ===== */
    case COCO_SETCMD_TRANSP:
        cp->dmacmd = COCO_TRANSP;
        break;
    case COCO_SETCMD_CONVERT:
        cp->dmacmd = COCO_CONVERT;
        break;
    /* =====
     *   Issue a Command
     * ===== */
    case COCO_COMMAND:
        if ( copyin ( (char *)arg, &coco_cmd,
            sizeof(coco_cmd_t) ) )
```

```

        return (EFAULT);
        cocoCommand ( cp, coco_cmd.cmd, coco_cmd.datav );
        break;
/* =====
 *      Set DMA type (Prog or Single)
 * ===== */
case COCO_SET_SINGLE_DMA:
    #ifdef DEBUG
        printf ("Chameleon DMA set to Single Mode\n");
    #endif
    cp->dmatype = DMA_SINGLE;
    break;
case COCO_SET_PROG_DMA:
    #ifdef DEBUG
        printf ("Chameleon DMA set to Prog Mode\n");
    #endif
    cp->dmatype = DMA_PROG;
    break;
/* =====
 *      Set up Chameleon Mode
 * ===== */
case COCO_SET_MODE:
    if ( copyin ( (char *)arg, &coco_mode,
        sizeof(coco_mode_t) ) )
        return (EFAULT);
    cocoSetMode ( cp, coco_mode.mode, coco_mode.swap,
        coco_mode.slice, coco_mode.flag );
    break;
/* =====
 *      Read Chameleon Mode
 * ===== */
case COCO_READ_MODE:
    tmp_int = cocoReadMode (cp);
    if ( copyout((char *)&tmp_int, arg, sizeof(uint_t) ) )
        return (EFAULT);
    break;
/* =====
 *      Raw write Buffer to Fifo      *
 * ===== */
case COCO_RAW_WRITEB_FIFO:
    if ( copyin( (char *)arg, &coco_buf, sizeof(coco_buf_t) ) )
        return (EFAULT);
    if ( coco_buf.buf_size == 0 )
        return (EINVAL);

```

```

/*
 * the data is in coco_buf.buf in user address
 * space. Move it to kernel address space and dump it
 * to the board.
 */
tot_bytes = coco_buf.buf_size * sizeof(uint_t);
tmp_ibuf = (uint_t *)kmem_alloc (tot_bytes, KM_NOSLEEP);
if ( tmp_ibuf == (uint_t *)NULL )
    return (ENOMEM);
if ( copyin((caddr_t)coco_buf.buf, (caddr_t)tmp_ibuf,
    tot_bytes) ) {
    kmem_free (tmp_ibuf, tot_bytes);
    return(EFAULT);
}
cocoBufOut (cp, tmp_ibuf, coco_buf.buf_size);
kmem_free (tmp_ibuf, tot_bytes );
break;
/* =====
 * Raw Read Buffer from Fifo      *
 * ===== */
case COCO_RAW_READB_FIFO:
    if ( copyin( (char *)arg, &coco_buf, sizeof(coco_buf_t) ) )
        return (EFAULT);
    if ( coco_buf.buf_size == 0 )
        return (EINVAL);
/*
 * the buffer to be filled is coco_buf.buf and
 * it can contain coco_buf.buf_size 32-bit values.
 * Read that many into our own temp buffer and move them
 * to user-address space.
 */
tot_bytes = coco_buf.buf_size * sizeof(uint_t);
tmp_ibuf = (uint_t *)kmem_alloc (tot_bytes, KM_NOSLEEP);
if ( tmp_ibuf == (uint_t *)NULL )
    return (ENOMEM);
cocoBufIn(cp, tmp_ibuf, coco_buf.buf_size);
if ( copyout ( (caddr_t)tmp_ibuf, (caddr_t)coco_buf.buf,
    tot_bytes) ) {
    kmem_free ( tmp_ibuf, tot_bytes );
    return (EFAULT);
}
kmem_free ( tmp_ibuf, tot_bytes );
break;
/* =====
 * Read 32-bit from Fifo

```

```
* ===== */
case COCO_RAW_READ_FIFO:
    tmp_int = cocoReadAmccFifo(cp);
    if ( copyout((char *)&tmp_int, arg, sizeof(uint_t) ) )
        return(EFAULT);
    break;
/* =====
*   Write 32-bit value to Fifo
* ===== */
case COCO_RAW_WRITE_FIFO:
    if ( copyin(arg, (char *)&tmp_int, sizeof(uint_t) ) )
        return(EFAULT);
    cocoCommand(cp, COCO_TRANSP, tmp_int );
    break;
/* =====
*   Test Fifo data path
* ===== */
case COCO_FIFO_TEST:
    return (cocoFifoTest(cp, 1) );
/* =====
*   Read DMA Registers
* ===== */
case COCO_RAW_READ_DMA:
    cocoReadDmaRegs(cp, &dmaRegs[0]);
    if ( copyout((char *)&dmaRegs[0], arg,
        sizeof(dmaRegs) ) )
        return (EFAULT);
    break;
/* =====
*   Write DMA Registers
* ===== */
case COCO_RAW_WRITE_DMA:
    if ( copyin(arg, (char *)&dmaRegs[0],
        sizeof(dmaRegs) ) )
        return (EFAULT);
    cocoWriteDmaRegs(cp, &dmaRegs[0] );
    break;
/* =====
*   DMA Registers Access Test
* ===== */
case COCO_DMAREGS_TEST:
    return ( cocoDmaRegsTest(cp) );
/* =====
*   Read Address Register
* ===== */
```

```
case COCO_READ_ADDR:
    tmp_int = cocoReadAddr(cp);
    if ( copyout((char *)&tmp_int, arg, sizeof(uint_t) ) )
        return (EFAULT);
    break;
/* =====
 *   Write Address Register
 * ===== */
case COCO_SET_ADDR:
    if ( copyin(arg, (char *)&tmp_int, sizeof(uint_t) ) )
        return(EFAULT);
    cocoSetAddr(cp, tmp_int);
    break;
/* =====
 *   Internal RAM Test
 * ===== */
case COCO_INIRAM_TEST:
    return ( cocoIntRamTest(cp) );
/* =====
 *   External RAM Test
 * ===== */
case COCO_EXTRAM_TEST:
    return ( cocoExtRamTest(cp) );

/* =====
 *   Read Internal LUTs
 * ===== */
case COCO_READ_RAMIL:
case COCO_READ_RAMIH:
case COCO_READ_RAMO:
    if ( copyin( (char *)arg, (char *)&coco_buf,
        sizeof(coco_buf_t) ) )
        return (EFAULT);
    if ( coco_buf.buf_size == 0 )
        return (EINVAL);
    if ( coco_buf.buf_size > INT_RAM_SIZE )
        return (EINVAL);
    /* allocate a buffer to hold Ram's contents */
    tot_bytes = coco_buf.buf_size * sizeof(uint_t);
    tmp_ibuf = (uint_t *)kmem_alloc (tot_bytes, KM_NOSLEEP);
    if ( tmp_ibuf == (uint_t *)NULL )
        return(EFAULT);
    /* read in Ram's contents into the buffer */
    cocoReadIntRam(cp, tmp_ibuf, cmd, coco_buf.buf_size );
    /* copy the kernel buffer to user's buffer */
```

```

        if ( copyout((caddr_t)tmp_ibuf, (caddr_t)coco_buf.buf,
                    tot_bytes) ) {
            kmem_free ( (caddr_t)tmp_ibuf, tot_bytes );
            return (EFAULT);
        }
        kmem_free ( (caddr_t)tmp_ibuf, tot_bytes );
        break;
/* =====
 *      Write Internal LUTs
 * ===== */
case COCO_FILL_RAMIL:
case COCO_FILL_RAMIH:
case COCO_FILL_RAMO:
    if ( copyin( (char *)arg, &coco_buf, sizeof(coco_buf_t) )
        )
        return (EFAULT);
    if ( coco_buf.buf_size == 0 )
        return (EINVAL);
    if ( coco_buf.buf_size > 2 * INT_RAM_SIZE )
        return (EINVAL);
    /* allocate a buffer to hold user's data */
    tot_bytes = coco_buf.buf_size * sizeof(uint_t);
    tmp_ibuf = (uint_t *)kmem_alloc (tot_bytes, KM_NOSLEEP);
    if ( tmp_ibuf == (uint_t *)NULL )
        return(EFAULT);
    /* copy user's buffer to our own */
    if ( copyin((caddr_t)coco_buf.buf, (caddr_t)tmp_ibuf,
              tot_bytes) ) {
        kmem_free ( tmp_ibuf, tot_bytes );
        return(EFAULT);
    }
    /* fill the Ram with data just moved over */
    cocoWriteIntRam ( cp, tmp_ibuf, cmd,
                    coco_buf.buf_size );
    kmem_free ( tmp_ibuf, tot_bytes );
    break;
/* =====
 *      Read External LUT
 * ===== */
case COCO_READ_RAML:
    if ( copyin( (char *)arg, &coco_buf, sizeof(coco_buf_t) )
        )
        return (EFAULT);
    if ( coco_buf.buf_size == 0 )
        return (EINVAL);
    if ( coco_buf.buf_size > EXT_RAM_SIZE )
        return (EINVAL);

```

```

/* allocate a buffer to hold External Ram's contents */
tot_bytes = coco_buf.buf_size * sizeof(uint_t);
tmp_ibuf = (uint_t *)kmem_alloc (tot_bytes, KM_NOSLEEP);
if ( tmp_ibuf == (uint_t *)NULL )
    return(EFAULT);
/* fill the buffer with Ext. Ram's contents */
cocoReadExtRam(cp, tmp_ibuf, coco_buf.buf_size);
/* copy the data to user's buffer */
if ( copyout((caddr_t)tmp_ibuf, (caddr_t)coco_buf.buf,
    tot_bytes) ) {
    kmem_free ( tmp_ibuf, tot_bytes );
    return (EFAULT);
}
kmem_free ( tmp_ibuf, tot_bytes );
break;
/* =====
*      Write External LUT
* ===== */
case COCO_FILL_RAML:
    if ( copyin( (char *)arg, &coco_buf, sizeof(coco_buf_t) )
        return (EFAULT);
    if ( coco_buf.buf_size == 0)
        return (EINVAL);
    if ( coco_buf.buf_size > 2 * EXT_RAM_SIZE )
        return (EINVAL);
/* allocate a buffer to hold user's data */
tot_bytes = coco_buf.buf_size * sizeof(uint_t);
tmp_ibuf = (uint_t *)kmem_alloc (tot_bytes, KM_NOSLEEP);
if ( tmp_ibuf == (uint_t *)NULL )
    return(EFAULT);
/* copy user's buffer to our own */
if ( copyin((caddr_t)coco_buf.buf, (caddr_t)tmp_ibuf,
    tot_bytes) ) {
    kmem_free ( tmp_ibuf, tot_bytes );
    return(EFAULT);
}
/* fill the Ram with data just moved over */
cocoWriteExtRam ( cp, tmp_ibuf, coco_buf.buf_size);
kmem_free ( tmp_ibuf, tot_bytes );
break;
/* =====
*      Chameleon Single pixel Conversion
* ===== */
case COCO_CONVERT_PIXLE:
    /* read in the coco_convert struct from user's space */

```

```

        if ( copyin(arg, (char *)&tmp_int, sizeof(uint_t) ) )
            return (EFAULT);
        cocoConvert ( cp, tmp_int );
        break;
/* =====
 *      Chameleon Single pixel Conversion Test
 * ===== */
case COCO_CONVERT_TEST:
    /* read in the coco_convert struct from user's space */
    if ( copyin(arg, (char *)&coco_convert,
        sizeof(coco_convert_t) ) )
        return (EFAULT);
    cocoConvertTest ( cp, &coco_convert );
    /* copy the structure back */
    if ( copyout((char *)&coco_convert, arg,
        sizeof(coco_convert_t) ) )
        return(EFAULT);
    break;
/* =====
 *      DMA fill of LUTs (Internals and External)
 * ===== */
case COCO_BLOCK_FILL_RAMIL:
case COCO_BLOCK_FILL_RAMIH:
case COCO_BLOCK_FILL_RAMO:
case COCO_BLOCK_FILL_RAML:
    /* read in coco_buf_t struct..all we need is there */
    if ( copyin(arg, (char *)&coco_buf, sizeof(coco_buf_t) ) )
        return (EFAULT);
    /* empty buffer ? */
    if ( coco_buf.buf_size == 0 )
        return (EINVAL);
    /* DMA the data and report back the result */
    microtime ( &cp->call_time );
    err = cocoDmaToLuts(cp, &coco_buf, cmd );
    microtime ( &cp->ret_time );
    #ifdef DEBUGTIME
    cocoReportTime ( "cocoDmaToLut", cp, 1 );
    #endif
    return (err);
/* =====
 *      Simultaneous Read/Write DMA
 * ===== */
case COCO_RW_BUF:
    /* read in coco_rw_t struct..all we need is there */
    if ( copyin(arg, (char *)&coco_rw, sizeof(coco_rw_t) ) )

```

```

        return (EFAULT);
    if ( coco_rw.buf_size <= 0 ) {
        cmn_err (CE_NOTE,
            "cocoIoctl: Bad Sim.read/write buff size of %d\n",
            coco_rw.buf_size );
        return(EINVAL);
    }
    microtime ( &cp->call_time );
    err = cocoReadWrite(cp, &coco_rw );
    microtime ( &cp->ret_time );
    #ifdef DEBUGTIME
    cocoReportTime ( "cocoReadWrite", cp, 1 );
    #endif
    if ( err != 0 )
        cmn_err (CE_NOTE,
            "cocoIoctl: cocoReadWrite reported error %d\n", err );
    return (err);
} /** end switch */
return(0);
}
/*****
***          c o c o _ r e a d          ***
*****/
*
* Name:          coco_read
*
* Purpose:      Read entry routine. DMAs data from the board to the
*               user's address space.
*
* Returns:      0 = Success, or errno
*
*****/
int
coco_read( dev_t dev, uio_t *uiop, cred_t *crp)
{
    register card_t *cp;
    register vertex_hdl_t vhdl;
    register int ret;
    /* get to the board's info structure */
    vhdl = dev_to_vhdl ( dev );
    cp = (card_t *)device_info_get ( vhdl );
    cp->bp = (buf_t *)NULL;
    cp->addrList = 0;
    cp->dmastat = 0;
    cp->dmabits = 0;

```

```

        bzero ( &cp->start_time, sizeof(struct timeval) );
        bzero ( &cp->intr_time, sizeof(struct timeval) );
        bzero ( &cp->call_time, sizeof(struct timeval) );
        bzero ( &cp->ret_time, sizeof(struct timeval) );
#ifdef DEBUG
        printf ("coco_read: Reading %d bytes\n", uiop->uio_resid );
#endif
        /* do the transfer */
        microtime ( &cp->call_time );
        ret = uiophysio ( cocoStrategy, NULL, dev, B_READ, uiop );
        microtime ( &cp->ret_time );
        if ( cp->addrList ) {
            alenlist_done(cp->addrList);
            cp->addrList = 0;
        }
#ifdef DEBUG
        printf ("coco_read: uiophysio returned %d\n", ret );
#endif
#ifdef DEBUGTIME
        cocoReportTime ( "cocoRead", cp, 1);
#endif
        return ( ret );
    }
}
/*****
***          c o c o _ w r i t e          ***
*****/
*
* Name:          coco_write
*
* Purpose:      Write entry routine. DMAs data from user's address space
*               to the board.
*
* Returns:      0 = Success, or errno
*
*****/
int
coco_write( dev_t dev, uio_t *uiop, cred_t *crp)
{
    register card_t      *cp;
    register vertex_hdl_t vhdl;
    register int         ret;
    /* get to the board's info structure */
    vhdl = dev_to_vhdl ( dev );
    cp = (card_t *)device_info_get ( vhdl );
    cp->bp          = (buf_t *)NULL;
}

```

```

cp->dmastat = 0;
cp->dmabits = 0;
cp->addrList = 0;
bzero ( &cp->start_time, sizeof(struct timeval) );
bzero ( &cp->intr_time, sizeof(struct timeval) );
bzero ( &cp->call_time, sizeof(struct timeval) );
bzero ( &cp->ret_time, sizeof(struct timeval) );
#ifdef DEBUG
printf ("coco_write: Writing %d bytes\n", uiop->uio_resid );
#endif
/* do the transfer */
microtime ( &cp->call_time );
ret = uiophysio ( cocoStrategy, NULL, dev, B_WRITE, uiop );
microtime ( &cp->ret_time );
if ( cp->addrList ) {
    alenlist_done(cp->addrList);
    cp->addrList = 0;
}
#ifdef DEBUG
printf ("coco_write: uiophysio returned %d\n", ret );
#endif
#ifdef DEBUGTIME
cocoReportTime ( "cocoWrite", cp, 1);
#endif
return ( ret );
}
/*****
***                c o c o _ e r r o r                ***
****
*
* Name:          coco_error
*
* Purpose:       Traps PCI bus error.
*
* Returns:       0 = Success, or errno
*
*****/
int
coco_error(vertex_hdl_t vhdl, int error)
{
    printf("coco_error %d\n", error );
    return(0);
}
/*
~~~~~

```

```

~~~~~
~~~~~      Supporting Routines      ~~~~~
~~~~~
*/
/*****
***          cocoStrategy          ***
*****/
*
* Name:      cocoStrategy
*
* Purpose:   Strategy routine. It actually handles read/write
*            and starts the DMA. It is called by uiophysio() kernel
*            routine.
*
* Returns:   0 = Success, or errno
*
*****/
static int
cocoStrategy ( buf_t *bp )
{
    register card_t    *cp;
    register vertex_hdl_t vhdl;
    register coco_dmapage_t *dmaPg;
    register int      s, i, ret, rw, tot_bytes;
    register uint_t   fill_bits;
    alenlist_t        addrList2;
    size_t            p_size;
    alenaddr_t        p_addr;
    iopaddr_t         p_dmaPg;
    if ( !BP_ISMAPPED(bp) ) {
        cmn_err (CE_NOTE, "cocoStrategy: Unmapped buf_t used\n");
        bioerror ( bp, EIO);
        biodone (bp);
        return(EIO);
    }
    /* get to the board's info structure */
    vhdl = dev_to_vhdl ( bp->b_edev );
    cp = (card_t *)device_info_get ( vhdl );
    /* clear error and save bp */
    bioerror(bp, 0);
    bp->b_resid = bp->b_bcount;
    cp->bp = bp;
    if ( bp->b_flags & B_READ )
        rw = B_READ; /* board -> mem */

```

```

else    rw = B_WRITE; /* mem  -> board */
/* create scatter-gather list */
addrList2 = buf_to_alenlist ( bp );
cp->addrList = pciio_dmatrans_list ( cp->vhdl, cp->dev_desc,
                                   addrList2,
                                   PCIIO_DMAMAP_BIGEND);
                                   /* PCIIO_DMAMAP_LITTLEEND); */

alenlist_done (addrList2);
if ( cp->addrList == (alenaddr_t)NULL ) {
    cmn_err (CE_NOTE, "cocoStrategy: Cannot create alenlist");
    bioerror ( bp, EIO);

    biodone (bp);
    return(EIO);
}
#endif
cp->page_no = alenlist_size ( cp->addrList );
#endif
cp->page_no = cocoAlenlistSize ( cp->addrList );
#ifdef DEBUG
printf ("cocoStrategy: %s %d bytes [%d pages] in %s Dma\n",
        rw == B_READ ? "Reading":"Writing",
        bp->b_bcount, cp->page_no,
        cp->dmatype == DMA_PROG ? "Chain":"Single" );
#endif
/* set the command to be executed during Dma */
cp->dmabits = cp->dmacmd;
cocoPrepAmcc ( cp );
/* =====
 *   Chained DMA
 * ===== */
if ( cp->dmatype == DMA_PROG ) {
    dmaPg = cocoMakeChain ( cp, cp->addrList, cp->page_no );
    if ( dmaPg == (coco_dmapage_t *)NULL ) {
        printf ("cocoStrategy: Error creating chain list\n");
        alenlist_done (cp->addrList);
        bioerror (bp, EIO);
        biodone (bp);
        return(EIO);
    }
    tot_bytes = cp->page_no * sizeof(coco_dmapage_t);
    p_dmaPg = pciio_dmatrans_addr ( cp->vhdl, cp->dev_desc,
                                   kvtophys(dmaPg), tot_bytes,
                                   PCIIO_DMAMAP_BIGEND );
    /* write back cache for chain list */
    dki_dcache_wbinval ( dmaPg, tot_bytes );
}

```

```

s = COCO_LOCK();
/* start the Chained Dma */
cocoStartProgDma ( cp, p_dmaPg, bp->b_bcount, rw );
if ( rw == B_READ )
    cp->dmastat = DMA_READ_WAIT;
else cp->dmastat = DMA_WRITE_WAIT;
cp->chain_list = (caddr_t)dmaPg;
#ifdef DEBUG
printf ("cocoStrategy: Waiting for chain interrupt\n");
#endif
/* so we dont wait forever for an interrupt */
cp->tid = itimeout(cocoTimeOut, cp, RW_TIMER,
                 pltimeout, 0, 0, 0);
COCO_UNLOCK(s);
return(0);
}
/* =====
*      Single page DMA
* ===== */
/* get a page to DMA */
if ( alenlist_get(cp->addrList, NULL, NBPP,
                 &p_addr, &p_size) != ALENLIST_SUCCESS ) {
    cmn_err (CE_WARN, "cocoDma: Not enough Memory");
    alenlist_done(cp->addrList);
    bioerror ( bp, EIO);
    biodone(bp);
    return(EIO);
}
/* single dma, memory -> board */

s = COCO_LOCK();
cocoStartSingleDma ( cp, p_addr, p_size, rw );
if ( rw == B_READ )
    cp->dmastat = DMA_READ_WAIT;
else cp->dmastat = DMA_WRITE_WAIT;
cp->chain_list = NULL;
#ifdef DEBUG
printf ("cocoStrategy: Waiting for Single interrupt\n");
#endif
/* so we dont wait forever for an interrupt */
cp->tid = itimeout(cocoTimeOut, cp, RW_TIMER,
                 pltimeout, 0, 0, 0);
#ifdef DEBUG
printf ("cocoStrategy: Waiting for Single interrupt, tid = %d\n",
        cp->tid );

```

```

        #endif
        COCO_UNLOCK(s);
        return(0);
    }
/*****
***          c o c o R e a d W r i t e          ***
****
*
* Name:          cocoReadWrite
*
* Purpose:       Starts simultaneous Read and Write DMA to user's space.
*
* Returns:       0 for success, or errno
*
*****/
static int
cocoReadWrite ( card_t *cp, coco_rw_t *rw )
{
    register caddr_t  dum;
    register caddr_t  w_kvaddr, r_kvaddr;
    register uint_t   *cache_line;
    register int      r_page_no, w_page_no, s, err, i, cache_bytes;
    register int      tot_r_cache, tot_w_cache, tot_bytes;
    register int      coco_adjust_chain, tot_wrong;
    coco_dmapage_t    *w_dmaPg, *r_dmaPg;
    alenlist_t        addrList2;
    iopaddr_t         pr_dmaPg, pw_dmaPg;
    cocoResetAmcc(cp);
    cocoPrepAmcc(cp);
    tot_bytes = rw->buf_size * sizeof(uint_t);
    coco_adjust_chain = rw->adjust_chain;
    /* =====
    *      Scatter-Gather list for Write buffer (mem -> board)
    *      ===== */
    /* lock user's pages in memory for DMA */
    if ( cocoLockUser ( (caddr_t)rw->w_buf, tot_bytes, B_WRITE) != 0 ) {
        cmn_err (CE_WARN, "cocoReadWrite: Cannot lock user's pages");
        return (EFAULT);
    }
    /*
    *   write back and invalidate the data for mem->board
    *   adjust the address for current data cache size
    */
#ifdef R10000
    cache_line = (uint_t *)rw->w_buf;
#endif
}

```

```

cache_line = (uint_t *) ( ((uint_t)cache_line / COCO_CACHE_SIZE) *
                           COCO_CACHE_SIZE );
cache_bytes = tot_bytes + ( (uint_t)rw->w_buf - (uint_t)cache_line);
dki_dcache_wbinval ( (caddr_t)cache_line, cache_bytes );
#else
dki_dcache_wbinval ( rw->w_buf, tot_bytes );
if ( cp->mappedkv )
    dki_dcache_wbinval ( cp->mappedkv, cp->mappedkvlen );
#endif
/* create scatter-gather list of user's buffer */
addrList2 = uvaddr_to_alenlist( (alenlist_t)NULL, (caddr_t)rw->w_buf,
                               (size_t)tot_bytes );

/* did we make it ? */
if ( addrList2 == (alenlist_t)NULL ) {
    cmn_err (CE_WARN, "cocoreadWrite: cannot create Wrt alenlist");
    cocoUnlockUser ( (caddr_t)rw->w_buf, tot_bytes, B_WRITE );
    return (EIO);
}

cp->w_addrList = pciio_dmatrans_list ( cp->vhdl, cp->dev_desc,
                                     addrList2, PCIIO_DMAMAP_BIGEND);
#if 0
w_page_no = alenlist_size ( cp->w_addrList );
#endif
w_page_no = cocoAlenlistSize ( cp->w_addrList );
cp->w_page_no = w_page_no;
alenlist_done ( addrList2 );
#ifdef DEBUG4
printf ("cocoReadWrite: Write %d bytes [%d pages]\n",
        tot_bytes, w_page_no );
#endif
/* =====
 * Scatter-Gather list for Read buffer ( board -> mem )
 * ===== */
/* lock user's pages in memory for DMA */
if ( cocoLockUser ( (caddr_t)rw->r_buf, tot_bytes, B_READ) != 0 ) {
    cmn_err (CE_WARN, "cocoReadWrite: Cannot lock user's pages");
    cocoUnlockUser ( (caddr_t)rw->w_buf, tot_bytes, B_WRITE );
    return (EFAULT);
}
/*
 * Invalidate the cache for board->mem
 * adjust the address for current data cache line size
 */
#endif R10000

```

```

cache_line = (uint_t *)rw->r_buf;
cache_line = (uint_t *) ( ((uint_t)cache_line / COCO_CACHE_SIZE) *
                           COCO_CACHE_SIZE );
cache_bytes = tot_bytes + ( (uint_t)rw->r_buf - (uint_t)cache_line);
dki_dcachel_inval ( (caddr_t)cache_line, cache_bytes);
#else
dki_dcachel_inval ( rw->r_buf, ((tot_bytes/NBPP)+1)*NBPP );
#endif
/* create scatter-gather list */
addrList2 = uvaddr_to_alenlist( (alenlist_t)NULL, (caddr_t)rw->r_buf,
                               (size_t)tot_bytes);

/* did we make it ? */
if ( addrList2 == (alenlist_t)NULL ) {
    cmn_err (CE_WARN, "cocoreadWrite: cannot create Rd alenlist");
    cocoUnlockUser ( (caddr_t)rw->w_buf, tot_bytes, B_WRITE );
    cocoUnlockUser ( (caddr_t)rw->r_buf, tot_bytes, B_READ);
    return (EIO);
}
cp->r_addrList = pciio_dmatrans_list ( cp->vhdl, cp->dev_desc,
                                     addrList2, PCIIO_DMAMAP_BIGEND);

#if 0
r_page_no = alenlist_size ( cp->r_addrList );
#endif
r_page_no = cocoAlenlistSize ( cp->r_addrList );
cp->r_page_no = r_page_no;
alenlist_done ( addrList2 );
#ifdef DEBUG4
printf ("cocoReadWrite: Read %d bytes [%d pages]\n",
        tot_bytes, r_page_no );
#endif
/* assume Single Dma type */
w_dmaPg = (coco_dmapage_t *)NULL;
r_dmaPg = (coco_dmapage_t *)NULL;
/* =====
 * Chain DMA - Prepare chain list for read and write
 * ===== */
if ( cp->dmatype == DMA_PROG ) {
    /* need to adjust read/write byte counts */
    if ( (coco_adjust_chain == 1) &&
          (rw->r_buf != rw->w_buf) ) {
        err = cocoMakeChainRW ( cp, &w_dmaPg, &r_dmaPg );
        if ( err ) {
            cmn_err (CE_WARN,
                    "cocoReadWrite: Could not make Chain List, err = %d",
                    err );
        }
    }
}

```

```

        cocoUnlockUser ( (caddr_t)rw->r_buf, tot_bytes,
                        B_READ);
        cocoUnlockUser ( (caddr_t)rw->w_buf, tot_bytes,
                        B_WRITE);
        alenlist_done(cp->r_addrList);
        alenlist_done(cp->w_addrList);
        return (err);
    }
}
/* no need to adjust read/write byte counts */
else {
    /* ~~~~~
     *   Prepare Chain List for Write
     * ~~~~~ */
#ifdef DEBUG
    printf (
        "cocoReadWrite: Preparing Chain for Write, %d pages\n",
            w_page_no );
#endif
    w_dmaPg = cocoMakeChain ( cp, cp->w_addrList, w_page_no );
    if ( w_dmaPg == (coco_dmapage_t *)NULL ) {
        printf ("cocoReadWrite: Error creating chain list\n");
        cocoUnlockUser ( (caddr_t)rw->r_buf, tot_bytes,
                        B_READ);
        cocoUnlockUser ( (caddr_t)rw->w_buf, tot_bytes,
                        B_WRITE);
        alenlist_done(cp->r_addrList);
        alenlist_done(cp->w_addrList);
        return(EIO);
    }
    /* ~~~~~
     *   Prepare Chain List for Read
     * ~~~~~ */
#ifdef DEBUG
    printf (
        "cocoReadWrite: Preparing Chain for Read, %d pages\n",
            r_page_no );
#endif
    r_dmaPg = cocoMakeChain ( cp, cp->r_addrList, r_page_no );
    if ( r_dmaPg == (coco_dmapage_t *)NULL ) {
        printf ("cocoReadWrite: Error creating chain list\n");
        cocoUnlockUser ( (caddr_t)rw->r_buf, tot_bytes,
                        B_READ);
        cocoUnlockUser ( (caddr_t)rw->w_buf, tot_bytes,
                        B_WRITE);
    }
}

```

```

        alenlist_done(cp->r_addrList);
        alenlist_done(cp->w_addrList);
        kmem_free ( w_dmaPg,
                    w_page_no * sizeof(coco_dmapage_t));
        return(EIO);
    }
}
/*~~~~~
 *   Map to PCI and cache management
 *~~~~~*/
/* map Write chain list */
if ( coco_adjust_chain == 1 )
    tot_w_cache = w_page_no * CHAIN_FACTOR *
                 sizeof(coco_dmapage_t);
else tot_w_cache = w_page_no * sizeof(coco_dmapage_t);
pw_dmaPg = pciio_dmatrans_addr ( cp->vhdl, cp->dev_desc,
                                kvtophys(w_dmaPg),
                                tot_w_cache,
                                PCIIO_DMAMAP_BIGEND );

#ifdef DEBUG
cocoShowChain( "Write Chain list", w_dmaPg );
#endif
/* write back the cache for this chain list */
dki_dcache_wbinval ( w_dmaPg, tot_w_cache );
/* map Read chain list */
if ( coco_adjust_chain == 1 )
    tot_r_cache = r_page_no * CHAIN_FACTOR *
                 sizeof(coco_dmapage_t);
else tot_r_cache = r_page_no * sizeof(coco_dmapage_t);
pr_dmaPg = pciio_dmatrans_addr ( cp->vhdl, cp->dev_desc,
                                kvtophys(r_dmaPg),
                                tot_r_cache,
                                PCIIO_DMAMAP_BIGEND );

#ifdef DEBUG
cocoShowChain( "Read Chain List", r_dmaPg );
#endif
/* write back the cache for this chain list */
dki_dcache_wbinval ( r_dmaPg, tot_r_cache );
}
/* initialize our DMA event semaphore */
initnsema ( &cp->dmawait, 0, "coco" );
#ifdef DEBUG
printf ("cocoReadWrite: Read/Write %d bytes\n", tot_bytes );
#endif
cp->iostat = IO_OK;

```

```

cp->dmabits = cp->dmacmd;
cp->wp_addr = (alenaddr_t)0;
cp->rp_addr = (alenaddr_t)0;
cp->wp_size = 0;
cp->rp_addr = 0;
s = COCO_LOCK();
err = cocoStartRWDma( cp, pw_dmaPg, tot_bytes, pr_dmaPg, tot_bytes);
if ( err == 0 ) {
    err = SleepEvent(&cp->dmawait);
    if ( err ) {
        err = EINTR;
        cmn_err (CE_NOTE, "cocoReadWrite: Interrupted");
    }
    else {
        if ( cp->iostat == IO_OK ) {
            #ifdef DEBUG
                printf ("cocoReadWrite: woken up\n");
            #endif
        }

        if ( cp->iostat == IO_TIME ) {
            cmn_err (CE_NOTE, "cocoReadWrite: Timed out");
            err = ETIME;
        }
        if ( cp->iostat == IO_ERROR ) {
            cmn_err (CE_NOTE, "cocoReadWrite: IO Error");
            err = EIO;
        }
    }
}
else {
    cmn_err (CE_WARN,
            "cocoReadWrite: Could not start Sim read/write");
}
/* we are done */
cp->dmastat = DMA_IDLE;
cp->dmabits = 0;
/*
 * Invalidate the cache for board->mem
 * adjust the address for current data cache line size
 */
#ifdef R10000
cache_line = (uint_t *)rw->r_buf;
cache_line = (uint_t *) ( ((uint_t)cache_line / COCO_CACHE_SIZE) *
                          COCO_CACHE_SIZE );
#endif

```

```

        cache_bytes = tot_bytes + ( (uint_t)rw->r_buf - (uint_t)cache_line);
        dki_dcache_inval ( (caddr_t)cache_line, cache_bytes);
#else
        dki_dcache_inval ( rw->r_buf, ((tot_bytes/NBPP)+1)*NBPP );
#endif
        cocoUnlockUser ( (caddr_t)rw->w_buf, tot_bytes, B_WRITE );
        cocoUnlockUser ( (caddr_t)rw->r_buf, tot_bytes, B_READ);
        alenlist_done(cp->r_addrList);
        alenlist_done(cp->w_addrList);
        if ( r_dmaPg ) kmem_free ( r_dmaPg, tot_r_cache );
        if ( w_dmaPg ) kmem_free ( w_dmaPg, tot_w_cache);
        COCO_UNLOCK(s);
        return(err);
}
/*****
***                c o c o S t a r t R W D m a                ***
****
*
* Name:          cocoStartRWDma
*
* Purpose:       Program the boards for Simultaneous Read/Write DMA.
*                The read and write channel's scatter-gather have been
*                prepared before and are in cp->r_addrList and cp->w_addrList
*                For Chained DMA, the chain list for read and write channels
*                are prepared before and passed to us as parameters.
*
* Returns:       0 = Success, or errno
*
*****/
static int
cocoStartRWDma ( card_t *cp, iopaddr_t pw_dmaPg, int tot_write,
                iopaddr_t pr_dmaPg, int tot_read )
{
    register caddr_t  adr_cfg, adr_norm, adr_amcc;
    register int      i, err, tot_words, tot_bytes;
    register uint_t   dmabits, dmacfg, dmacmd;
    size_t            rp_size, wp_size;
    alenaddr_t        rp_addr, wp_addr;
    adr_cfg = cp->conf_adr;
    adr_norm = cp->norm_adr;
    adr_amcc = cp->amcc_adr;
    dmacfg = cp->dmacfg;
    dmabits = cp->dmabits;
    dmacmd = cp->dmacmd;
    Out32(adr_cfg, dmacfg ); /* clear eof markers */

```

```

/* =====
 *   Chained DMA
 * ===== */
if ( cp->dmatype == DMA_PROG ) {
    /* ~~~~~~
     *   Programming the board for Read/Write Chain Dma
     * ~~~~~~ */
    /* prepare Amcc counters */
    Out32(adr_amcc+AMCC_OP_REG_MRTC, tot_write );
    Out32(adr_amcc+AMCC_OP_REG_MWTC, tot_read );
    /* Program the Write channel Address (board -> mem) */
    dmabits |= DMAREG_PREN | DMAREG_PWEN;
    dmacfg |= dmabits;
    Out32(adr_cfg, dmacfg);
    /* Write channel address (b->m) */
    Out32(adr_cfg, dmacfg | DMAREG_RAM | DMAREG_RAMPRWR );
    Out32(adr_norm, pr_dmaPg);
    /* Read channel address (m->b) */
    Out32(adr_cfg, dmacfg | DMAREG_RAM | DMAREG_RAMPRRD );
    Out32(adr_norm, pw_dmaPg);
    /* Start Read/Write DMA - only Write(b->m) is enabled */
    dmabits |= DMAREG_INTPWEN | DMAREG_WEN | DMAREG_REN | dmacmd;
    cp->dmasize = tot_write;
    cp->dmabits = dmabits;
    cp->iostat = IO_OK;
    cp->dmastat = DMA_RW_WAIT;
    #ifdef DEBUG
    printf (
"cocoStartRWDma: read chain = 0x%x for %d, write chain = 0x%x for %d\n",
        pr_dmaPg, tot_read, pw_dmaPg, tot_write );
    cocoReport (cp, "cocoStartRWDma: Chain Dma started");
    #endif
    /* start the DMA */
    pciio_flush_buffers ( cp->vhdl );
    microtime ( &cp->start_time );
    Out32(adr_cfg, dmacfg | dmabits );
    /* so we dont wait forever for an interrupt */
    cp->tid = itimeout(cocoTimeOut2, cp, SIMRW_TIMER,
        pltimeout, 0, 0, 0);

    return(0);
}
/* =====
 *   Single page DMA
 * ===== */
/* get next page for Read channel if nothing left from past */

```

```
    if ( cp->wp_addr == (alenaddr_t)0 ) {
        if ( cp->w_page_no <= 0 ) {
            #ifdef DEBUG
                printf (
"cocostartRWDma: Premature Write, w_page_no = %d, r_page_no = %d, wp_addr =
0x%x, rp_addr = 0x%x\n",
                cp->w_page_no, cp->r_page_no, cp->wp_addr, cp->rp_addr );
            #endif
            cmn_err (CE_WARN, "cocostartRWDma: Premature end of Write");
            return(EIO);
        }
        if ( alenlist_get(cp->w_addrList, NULL, NBPP,
            &wp_addr, &wp_size) != ALENLIST_SUCCESS ) {
            cmn_err (CE_WARN, "cocostartRWDma: bad scater-gather");
            return(EIO);
        }
        cp->w_page_no--;
    }
    /* some bytes left from past */
    else {
        wp_addr = cp->wp_addr;
        wp_size = cp->wp_size;
        #ifdef DEBUG
            printf (
                "cocostartRWDma: %d old bytes at 0x%x for read chan\n",
                    wp_size, wp_addr );
        #endif
    }
    /* get next page for Write channel if nothing left from past */
    if ( cp->rp_addr == (alenaddr_t)0 ) {
        if ( cp->r_page_no <= 0 ) {
            #ifdef DEBUG
                printf (
"cocostartRWDma: Premature Read, r_page_no = %d, w_page_no = %d, wp_addr = 0x%x,
rp_addr = 0x%x\n",
                cp->r_page_no, cp->w_page_no, cp->wp_addr, cp->rp_addr );
            #endif
            cmn_err (CE_WARN, "cocostartRWDma: Premature end of Read");
            return(EIO);
        }
        if ( alenlist_get(cp->r_addrList, NULL, NBPP,
            &rp_addr, &rp_size) != ALENLIST_SUCCESS ) {
            cmn_err (CE_WARN, "cocoReadWrite: bad scater-gather");
            return(EIO);
        }
    }
}
```

```

        cp->r_page_no--;
    }
    /* some bytes left from past */
    else {
        rp_addr = cp->rp_addr;
        rp_size = cp->rp_size;
#ifdef DEBUG
        printf (
            "cocoStartRWDma: %d old bytes at 0x%x for Write chan\n",
                rp_size, rp_addr );
#endif
    }
    /* adjust - we should write as much as we can read */
    cp->wp_addr = (alenaddr_t)0;
    cp->rp_addr = (alenaddr_t)0;
    cp->wp_size = 0;
    cp->rp_size = 0;
#ifdef DEBUG
    printf (
        "cocoStartRWDma: Original sizes, rp_size = %d at 0x%x, wp_size = %d at 0x%x\n",
            rp_size, rp_addr, wp_size, wp_addr );
#endif
    /* Write more than read ? */
    if ( wp_size == rp_size )
        tot_bytes = wp_size;
    if ( wp_size > rp_size ) {
        tot_bytes = rp_size;
        cp->wp_addr = (alenaddr_t)((int)wp_addr + tot_bytes);
        cp->wp_size = wp_size - tot_bytes;
    }
    /* read more than write ? */
    if ( rp_size > wp_size ) {
        tot_bytes = wp_size;
        cp->rp_addr = (alenaddr_t)((int)rp_addr + tot_bytes);
        cp->rp_size = rp_size - tot_bytes;
    }
    tot_words = (int)tot_bytes/sizeof(uint_t);
#ifdef DEBUG
    printf (
        "cocoStartRWDma: Single Dma %d [%d words], rp_addr = 0x%x, wp_addr = 0x%x\n",
            tot_bytes, tot_words, rp_addr, wp_addr );
#endif
    tot_words--; /* counts down to 0xffff in hardware */
    /* ~~~~~~
    *   Program the board for Read and Write

```

```

    * ~~~~~*/
    /* prepare Amcc counters */
    /* Out32(adr_amcc+AMCC_OP_REG_MRTC, tot_bytes); */
    Out32(adr_amcc+AMCC_OP_REG_MRTC, tot_bytes);
    Out32(adr_amcc+AMCC_OP_REG_MWTC, tot_bytes);
    /* program for Write channel ( board -> mem ) */
    Out32(adr_cfg, dmacfg | DMAREG_WCNT);
    Out32(adr_norm, (uint_t)tot_words);
    Out32(adr_cfg, dmacfg | DMAREG_RAM | DMAREG_RAMWADR); /* b->m */
    Out32(adr_norm, (uint_t)rp_addr);
    /* program for Read channel ( mem -> board ) */
    Out32(adr_cfg, dmacfg | DMAREG_RCNT);
    Out32(adr_norm, (uint_t)tot_words);
    Out32(adr_cfg, dmacfg | DMAREG_RAM | DMAREG_RAMRADR); /* b->m */
    Out32(adr_norm, (uint_t)wp_addr);
    /* start read/write Dma - only write channel (b->m) is needed */
    dmabits |= DMAREG_INTWEN | DMAREG_REN | DMAREG_WEN | dmacmd;
    cp->dmasize = wp_size;
    cp->dmabits = dmabits;
    cp->dmastat = DMA_RW_WAIT;
    cp->iostat = IO_OK;
    /* start the DMA */
    pciio_flush_buffers ( cp->vhdl );
    microtime ( &cp->start_time );
    Out32(adr_cfg, dmacfg | dmabits | dmacmd );
    /* so we dont wait forever for an interrupt */
    cp->tid = itimeout(cocoTimeOut2, cp, SIMRW_TIMER,
                    plttimeout, 0, 0, 0);
    return(0);
}
/*****
***          c o c o R e s e t          ***
*****/
*
* Name:          cocoReset
*
* Purpose:      This routine initializes the board by resetting the add-on
*               logic, AMCC fifos, Chameleon chip and FIFOs.
*               After reset, FIFO flags are programmed, default DMA controller
*               config. register is written and Chameleon chip is enabled
*
* Returns:      None
*
*****/
static void

```

```

cocoReset( card_t *cp )
{
    uint_t  stat;
    register caddr_t  adr_amcc;
    register caddr_t  adr_cfg;
    register uint_t  tmp;
    adr_amcc = cp->amcc_adr;
    adr_cfg  = cp->conf_adr;
    /* disable interrupts */
    Out32(adr_amcc+AMCC_OP_REG_INTCSR, AMCC_INTCSR_RST | AMCC_INTCSR_RCLR |
        AMCC_INTCSR_WCLR );

    /* reset AMCC fifos and Xilinx */
    Out32(adr_amcc+AMCC_OP_REG_MCSR, AMCC_RST_ADDON );
    Out32(adr_amcc+AMCC_OP_REG_MCSR, AMCC_RST_FIFOS );
    /* reset Chameleon and FIFOs and bring FIFOs to programming flags */
    Out32(adr_cfg, DMAREG_CCRES | DMAREG_FRES | DMAREG_FSCLK );
    Out32(adr_cfg, DMAREG_FSCLK);
    Out32(adr_cfg, DMAREG_FSCLK);
    /* program FIFO flags */
    cocoProgFlags ( adr_cfg, DMAREG_FRES,100,450,100,450);
    /* bring FIFOs in functional mode and Chameleon out of reset */
    Out32(adr_cfg, DMAREG_FRES | DMAREG_FSCLK);
    Out32(adr_cfg, DMAREG_CCRES | DMAREG_FSCLK);
    Out32(adr_cfg, DMAREG_CCRES | DMAREG_FSCLK | DMAREG_PTEN);
    /* set default config. reg */
    cp->dmacfg = DMAREG_NVIFEN | DMAREG_PTEN | DMAREG_CCRES | DMAREG_FSCLK;
    /* initialize Chameleon */
    cocoSetMode(cp, 0,0,0,1);
    /* enable Amcc Interrupt */
    Out32(adr_amcc+AMCC_OP_REG_INTCSR, AMCC_INTCSR_MASK );
}
/*****
***              c o c o P r o g F l a g s              ***
*****/
*
* Name:          cocoProgFlags
*
* Purpose:       Program offsets of ASIC fifo flags by writing 18 bit
*                serial registers of FIFO containg 9 bit AF and 9 bit AE
*                flag offset
*
* Returns:       None
*
*****/
static void

```

```
cocoProgFlags ( caddr_t adr_cfg, uint_t d, ushort_t iae, ushort_t iaf,
                ushort_t oae, ushort_t oaf )
{
    ushort_t andfl;
    ushort_t shft;
    ushort_t bit0;
    ushort_t bit1;
    uint_t   bit;
    andfl=0x01;
    shft=0;
#ifdef DEBUG
    printf ("cocoProgFlags: iae = 0x%x, oae = 0x%x, oaf = 0x%x\n",
           iae, oae, oaf );
#endif
    do {
        bit0 = (iaf & andfl);
        bit1 = (oaf & andfl);
        bit = d | 0x00000010;
        if (bit0 == 0)
            bit |= 0x00000004;
        if (bit1 == 0)
            bit |= 0x00000008;
        Out32(adr_cfg, bit);
        bit &= 0xFFFFFFFF;
        Out32(adr_cfg, bit);
        andfl = andfl << 1;
        shft++;
    }
    while ( shft <= 8 );
    andfl = 0x01;
    shft = 0;
    do {
        bit0 = (iae & andfl);
        bit1 = (oae & andfl);
        bit = d | 0x00000010;
        if (bit0==0)
            bit |= 0x00000004;
        if (bit1==0)
            bit |= 0x00000008;
        Out32(adr_cfg, bit);
        bit = bit & 0xFFFFFFFF;
        Out32(adr_cfg, bit);
        andfl = andfl << 1;
        shft++;
    }
}
```

```

        while (shft<=8);
    }
/*****
***          c o c o S e t   M o d e          ***
****
*
* Name:      cocoSetMode
*
* Purpose:   Set Chameleon Mode Register
*
* Returns:   None
*
*****/
static void
cocoSetMode ( card_t *cp, int mode, int swap, int slice, int flag )
{
    uint_t cmd;
    cmd = COCO_WRENA | COCO_DELAY;
    if (mode)
        cmd |= COCO_MODE;

    if (swap)
        cmd |= COCO_SWAP;
    if (slice)
        cmd |= COCO_SLICE;
    if (flag)
        cmd |= COCO_FLAG;

#ifdef DEBUG
    printf ("cocoSetMode: setting mode to 0x%x\n", cmd );
#endif
    cocoCommand (cp, COCO_SETMODE, cmd);
}
/*****
***          c o c o C o m m a n d          ***
****
*
* Name:      cocoCommand
*
* Purpose:   Sends a command to Chameleon. Writes the command to
*            the controller and writes the data to Fifo.
*
* Returns:   None
*
*****/

```

```

*****/
static void
cocoCommand ( card_t *cp, uint_t cmd, uint_t data )
{
    register uint_t  stat;
#ifdef 0
    stat = Inp32(cp->amcc_adr+AMCC_OP_REG_MCSR);
    while ( stat & 0x1 ) {
        printf ("Fifo full..waiting\n");
        delay(1);
        stat = Inp32(cp->amcc_adr+AMCC_OP_REG_MCSR);
    }
#endif
    /* set Chameleon connad in Config. Register */
    Out32(cp->conf_adr, cmd | cp->dmacfg );
    /* writes the data for the command to AMCC Fifo */
    Out32(cp->amcc_adr + AMCC_OP_REG_FIFO, data);
}
/*****
***          c o c o B u f O u t          ***
*****
*
* Name:          cocoBufOut
*
* Purpose:      Dumps a buffer containig 32-bit values to Chameleon Fifo.
*               Data is given to Chameleon in Transparent mode.
*
* Returns:      None
*
*****/
static void
cocoBufOut ( card_t *cp, uint_t *buf, int size )
{
    register int    i, j;
    register uint_t cmd, stat;
#ifdef DEBUG
    printf ("cocoBufOut: outputing %d words to Fifo\n", size );
#endif
    /* set to Transparent Mode */
    cmd = cp->dmacfg | COCO_TRANSP;
    Out32(cp->conf_adr, cmd );
    /* Dump the data */
    for ( i = 0; i < size; i++) {
        j = 0;
        stat = Inp32(cp->amcc_adr+AMCC_OP_REG_MCSR);

```

```

        if ( stat & 0x03 ) { /* fifo has room */
            Out32(cp->amcc_adr+AMCC_OP_REG_FIFO, buf[i] );
            continue;
        }
        printf (
"cocoBufOut: Fifo is full, stat = 0x%x, discarding %d bytes\n",
                stat, size - i );
        break;
    }
}
/*****
***          c o c o B u f I n          ***
*****/
*
* Name:          cocoBufIn
*
* Purpose:       Reads Fifo entries into a buf.
*
* Returns:       None
*
*****/
static void
cocoBufIn ( card_t *cp, uint_t *buf, int size )
{
    register int    i;
    for ( i = 0; i < size; i++ ) {
        buf[i] = cocoReadAmccFifo(cp);
    }
}
/*****
***          c o c o R e a d A m c c F i f o          ***
*****/
*
* Name:          cocoReadAmccFifo
*
* Purpose:       Reads a 32-bit value word from AMCC Internal Fifo
*                Returns zero if Fifo is empty
*
* Returns:       Fifo value or zero if fifo is empty
*
*****/
static int
cocoReadAmccFifo ( card_t *cp )
{

```

```

        register uint_t  stat;
        register int    i;
        for ( i = 0; i < 3; i++ )
            stat = Inp32(cp->amcc_adr+AMCC_OP_REG_MCSR);
        /* if ( stat != 0x000000e6 ) { */
        if ( stat != 0x1026 ) {
            stat = Inp32(cp->amcc_adr + AMCC_OP_REG_FIFO );
            return(stat);
        }
        return(0);
    }
}
/*****
***          c o c o F i f o T e s t          ***
****
*
* Name:          cocoFifoTest
*
* Purpose:      Fills Fifo with a pattern, read it back one by one
*               and compare. Reports back the results.
*
* Returns:      0 = Success, 1 = Failed.
*
*****/
static int
cocoFifoTest ( card_t *cp, int pat )
{
    register int  i, j, err;
    uint_t  res, expect;
    #ifdef DEBUG
    printf ("cocoFifoTest:  testing Fifo, pattern: %d\n", pat );
    #endif
    err = 0;  /* assume success */
    /* test for 50 times .. */
    for( j = 0; j < 50; j++) {

        /* fill Fifo with pattern */
        for(i=0;i<750;i++)
            cocoCommand (cp, COCO_TRANSP, cocoPattern(pat,i) );
        /* read it back and compare */
        for( i = 0;i < 750; i++) {
            res = cocoReadAmccFifo( cp );  /* what we read */
            expect = cocoPattern( pat, i); /* what we expect */
            if (res != expect) {
                printf(
                    "cocoFifoTest: Fifo entry %d, expected 0x%x, read 0x%x (round %d)\n",

```

```
        i, expect, res, j);
        err = 1;
    }
}
return(err);
}
/*****
***          c o c o P a t t e r n          ***
****
*
* Name:          cocoPattern
*
* Purpose:       Pattern generator for testing purpose.
*
* Returns:       The generated pattern
*
*****/
static int
cocoPattern ( int pat, int cnt )
{
    register int res;
    /* for now, we always return pattern 1 */
    pat=1;

    switch (pat) {
        case 1:
            res = cnt;
            break;
        case 2:
            res = ~cnt;
            break;
        case 3:
            res = 0xFFFFFFFFL;
            break;
        case 4:
            res = 0x00000000L;
            break;
        case 5:
            if (cnt % 2 != 1)
                res = 0xaaaaaaaaL;
            else res = 0x55555555L;
            break;
        case 6:
            cnt = (cnt & 0x01FL);
    }
}
```

```

        res = (1 << cnt);
        break;
    case 7:
        cnt = (cnt & 0x01FL);
        res = (1 << cnt);
        res = ~res;
        break;
    case 8:
        if ( cnt % 2 == 1)
            res = 0x0f0f0f0fL;
        else res = 0xf0f0f0f0L;
        break;
    case 9:
        cnt = (cnt & 0xFF);
        res = cnt + ( cnt << 8) + (cnt << 16) + (cnt << 24);
        break;
    case 10:
        cnt = (cnt & 0xFFFF);
        res = cnt + ( cnt << 16);
        break;
    }
    return(res);
}
}
/*****
***          c o c o R e a d M o d e          ***
*****/
*
* Name:          cocoReadMode
*
* Purpose:       Reads Chameleon Mode register
*
* Returns:       The generated pattern
*
*****/
static int
cocoReadMode ( card_t *cp )
{
    register uint_t  mode;
    /* set Chameleon command in Config. Register */
    cocoCommand ( cp, COCO_READMODE, 0x0 );
    mode = cocoReadAmccFifo(cp);
#ifdef DEBUG
    printf ("cocoReadMode:  Mode = 0x%x\n", mode );
#endif
    return (mode);
}

```

```

}
/*****
***          c o c o R e a d D m a R e g s          ***
****
*
* Name:      cocoReadDmaRegs
*
* Purpose:   Reads Xilinx DMA Registers
*
* Returns:   None.
*
*****/
static void
cocoReadDmaRegs ( card_t *cp, uint_t *dmaRegs )
{
    register int  i, k;
    /* read DMA registers into dmaRegs table */
    i = 0;
    for ( k = 13; k <= 16; k++) {
        Out32(cp->conf_adr, cp->dmacfg | (k << 6) );
        dmaRegs[i++] = Inp32(cp->norm_adr);
    }
}
/*****
***          c o c o W r i t e D m a R e g s          ***
****
*
* Name:      cocoWriteDmaRegs
*
* Purpose:   Write values in dmaRegs[] to Xilinx DMA Registers
*
* Returns:   None.
*
*****/
static void
cocoWriteDmaRegs ( card_t *cp, uint_t *dmaRegs )
{
    register int  i, k;
    /* read DMA registers into dmaRegs table */
    i = 0;
    for ( k = 13; k <= 16; k++) {
        Out32(cp->conf_adr, cp->dmacfg | (k << 6) );
        Out32(cp->norm_adr, dmaRegs[i++] );
    }
}

```

```

/*****
***          c o c o R e a d A d d r          ***
****
*
* Name:          cocoReadAddr
*
* Purpose:       Reads Address Register
*
* Returns:       Address Register value
*
*****/
static int
cocoReadAddr (card_t *cp)
{
    register uint_t  addr_reg;
    cocoCommand ( cp, COCO_READADDR, 0x0 );
    addr_reg = cocoReadAmccFifo(cp);

    #ifdef DEBUG
    printf ("cocoreadAddr:  Address Reg = 0x%x\n", addr_reg );
    #endif
    return(addr_reg);
}
/*****
***          c o c o S e t A d d r          ***
****
*
* Name:          cocoSetAddr
*
* Purpose:       Set Address Register to given value
*
* Returns:       None
*
*****/
static void
cocoSetAddr ( card_t *cp, uint_t val )
{
    cocoCommand ( cp, COCO_SETADDR, val );
}
/*****
***          c o c o D m a R e g s T e s t          ***
****
*
* Name:          cocoDmaRegsTest
*
*****/

```

```

* Purpose:   Tests access to Xilinx DMA registers.
*           Write and read to DMA registers, compare values and
*           report back the results.
*
* Returns:   Test result ( 0 = Success, 1 = Failed)
*
*****/
static int
cocoDmaRegsTest ( card_t *cp )
{
    register int i, j, k, err, res, expct, dmacfg, pat;
    register caddr_t  adr_cfg, adr_norm;
    cmn_err (CE_NOTE, "testing Chameleon DMA registers access");
    adr_cfg = cp->conf_adr;
    adr_norm = cp->norm_adr;
    dmacfg  = cp->dmacfg;
    err     = 0;
    pat     = 1;
    /* try for 500 times */
    for ( i = 0; i < 500; i++ ) {

        /* write to DMA Regsiers */
        for ( k = 13; k <= 16; k++ ) {
            Out32(adr_cfg, dmacfg | (k << 6) );
            Out32(adr_norm, cocoPattern(pat, i+k) );
        }

        /* read back registers and compare */
        for ( k = 13; k <= 16; k++ ) {
            Out32(adr_cfg, dmacfg | (k << 6) );
            res = Inp32(adr_norm);
            expct = cocoPattern(pat, i+k);
            /* not equal ..report it ! */
            if ( res != expct ) {
                cmn_err ( CE_NOTE,
                    "Dma reg %d, expected 0x%x, read 0x%x, round %x",
                    k, expct, res, i );
                err = 1;
            }
        }
    }
    cmn_err (CE_NOTE, "Chameleon DMA registers access test %s",
        (err == 0 ? "Succeeded": "Failed") );
    return(err);
}

```

```

/*****
***          c o c o I n t R a m T e s t          ***
****
*
* Name:          cocoIntRamTest
*
* Purpose:       Test CHameleon Internal LUTs.
*               Fills Internal LUT with a pattern and reads them back
*               for comparison and reports the results.
*
* Returns:       Test result ( 0 = Success, 1 = Failed)
*
*****/
static int
cocoIntRamTest ( card_t *cp )
{
    register int i, err, pat, num;
    register uint_t res;
    cmn_err (CE_NOTE, "Testing Chameleon Internal Ram access\n");
    num = 1;
    err = 0;

    /* =====
     *   Fill LUTs
     * ===== */
    for ( i = 0; i < INT_RAM_SIZE; i++ ) {
        pat = cocoPattern(num, i);
        /* set address of internal LUT location */
        cocoCommand (cp, COCO_SETADDR, i);
        /* fill Internal LUTs with pattern */
        cocoCommand (cp, COCO_FILLRAMIL, pat );
        cocoCommand (cp, COCO_FILLRAMIH, pat );
        cocoCommand (cp, COCO_FILLRAMO, pat );
    }
    /* =====
     *   Read LUTs and compare
     * ===== */
    for ( i = 0; i < INT_RAM_SIZE; i++ ) {

        pat = cocoPattern(num, i) & 0x03ffff;
        /*
         * set address of internal RAMIL location
         * read value from that location and compare
         */
        cocoCommand (cp, COCO_SETADDR, i);
    }
}

```

```

cocoCommand ( cp, COCO_READRAMIL, 0x0 );
res = cocoReadAmccFifo (cp);

/* is what we read ok ? */
if ( (res & 0x03ffff) != pat ) {
    cmn_err (CE_NOTE,
             "Chameleon RAMIL at 0x%x, expected 0x%x, read 0x%x\n",
             i, pat, res );
    err = 1;
}
/*
 * set address of internal RAMIH location
 * read value from that location and compare
 */
cocoCommand (cp, COCO_SETADDR, i);
cocoCommand ( cp, COCO_READRAMIH, 0x0 );
res = cocoReadAmccFifo (cp);
/* is what we read ok ? */
if ( (res & 0x03ffff) != pat ) {
    cmn_err (CE_NOTE,
             "Chameleon RAMIH at 0x%x, expected 0x%x, read 0x%x\n",
             i, pat, res );
    err = 1;
}
/*
 * set address of internal RAMO location
 * read value from that location and compare
 */
cocoCommand (cp, COCO_SETADDR, i);
cocoCommand ( cp, COCO_READRAMO, 0x0 );
res = cocoReadAmccFifo (cp);
/* is what we read ok ? */
if ( (res & 0x03ffff) != pat ) {
    cmn_err (CE_NOTE,
             "Chameleon RAMO at 0x%x, expected 0x%x, read 0x%x\n",
             i, pat, res );
    err = 1;
}
}
cmn_err (CE_NOTE, "Chameleon Internal RAM test %s\n",
        ( err == 0 ? "Succeeded":"Failed" ) );
return(err);
}
/*****
***                               c o c o E x t R a m T e s t                               ***
*****/

```

```
*****
*
* Name:      cocoExtRamTest
*
* Purpose:   Test CHameleon External LUTs.
*           Fills External LUT with a pattern and reads them back
*           for comparison and reports the results.
*
* Returns:   Test result ( 0 = Success, 1 = Failed)
*
*****/
static int
cocoExtRamTest ( card_t *cp )
{
    register int  i, err, pat, num;
    register uint_t  res;
    cmn_err (CE_NOTE, "Testing Chameleon External Ram access\n");
    num = 1;
    err = 0;

    /* =====
     *   Fill RAML
     * ===== */
    for ( i = 0; i < EXT_RAM_SIZE; i++ ) {
        /* set address of internal LUT location */
        cocoCommand (cp, COCO_SETADDR, i);
        /* fill External LUTs with pattern */
        cocoCommand (cp, COCO_FILLRAML, cocoPattern(num, i) );
    }
    /* =====
     *   Read External LUT and compare
     * ===== */
    for ( i = 0; i < EXT_RAM_SIZE; i++ ) {

        pat = cocoPattern(num, i);
        /*
         * set address of internal RAMIL location
         * read value from that location and compare
         */
        cocoCommand (cp, COCO_SETADDR, i);
        cocoCommand ( cp, COCO_READRAML, 0x0 );
        res = cocoReadAmccFifo (cp);

        /* is what we read ok ? */
        if ( res != pat ) {
```

```

        cmn_err (CE_NOTE,
                "Chameleon RAML at 0x%x, expected 0x%x, read 0x%x\n",
                i, pat, res );
        err = 1;
    }
}
cmn_err (CE_NOTE, "Chameleon External RAM test %s\n",
        ( err == 0 ? "Succeeded":"Failed" ) );
return(err);
}
/*****
***                c o c o R e a d I n t R a m                ***
****
*
* Name:            cocoReadIntRam
*
* Purpose:        Read Chameleon's Internal LUTs into given buffer.
*
* Returns:        None.
*
*****/
static void
cocoReadIntRam ( card_t *cp, uint_t *buf, int lut, int size )
{
    register uint_t cmd;
    register int    i;
    /* which LUT we should read ? */
    switch ( lut ) {
        case COCO_READ_RAMIL:  cmd = COCO_READRAMIL;  break;
        case COCO_READ_RAMIH:  cmd = COCO_READRAMIH;  break;
        case COCO_READ_RAMO:   cmd = COCO_READRAMO;   break;
    }
    /* fill buffer with contents of Internal Ram */
    for ( i = 0; i < size; i++ ) {
        cocoCommand (cp, COCO_SETADDR, i);
        cocoCommand (cp, cmd, 0x0 );
        buf[i] = cocoReadAmccFifo (cp);
    }
}
/*****
***                c o c o W r i t e I n t R a m                ***
****
*
* Name:            cocoWriteIntRam
*
*****/

```

```

* Purpose:   Write Chameleon's Internal LUTs into given buffer.
*
* Returns:   None.
*
*****/
static void
cocoWriteIntRam ( card_t *cp, uint_t *buf, int lut, int size )
{
    register uint_t  cmd;
    register int     i;
    /* which LUT we should read ? */
    switch ( lut ) {
        case COCO_FILL_RAMIL:  cmd = COCO_FILLRAMIL;  break;
        case COCO_FILL_RAMIH:  cmd = COCO_FILLRAMIH;  break;
        case COCO_FILL_RAMO:   cmd = COCO_FILLRAMO;   break;
    }
    /* fill buffer with contents of Internal Ram */
    for ( i = 0; i < size; i+=2 ) {
        cocoCommand (cp, COCO_SETADDR, buf[i]);
        cocoCommand (cp, cmd, buf[i+1] );
    }
}
/*****
***               c o c o R e a d E x t R a m               ***
*****/
*
* Name:       cocoReadExtRam
*
* Purpose:    Read Chameleon's External LUT into given buffer.
*
* Returns:    None.
*
*****/
static void
cocoReadExtRam ( card_t *cp, uint_t *buf, int size )
{
    register uint_t  cmd;
    register int     i;
    /* fill External Ram with contents of the buffer */
    for ( i = 0; i < size; i++ ) {
        cocoCommand (cp, COCO_SETADDR, i);
        cocoCommand (cp, COCO_READRAML, 0x0 );
        buf[i] = cocoReadAmccFifo (cp);
    }
}

```

```

/*****
***          c o c o W r i t e E x t R a m          ***
*****/
*
* Name:      cocoWriteExtRam
*
* Purpose:   Write Chameleon's External LUT into given buffer.
*
* Returns:   None.
*
*****/
static void
cocoWriteExtRam ( card_t *cp, uint_t *buf, int size )
{
    register int    i;
    /* fill buffer with contents of Internal Ram */
    for ( i = 0; i < size; i+=2 ) {
        cocoCommand (cp, COCO_SETADDR, buf[i]);
        cocoCommand (cp, COCO_FILLRAM, buf[i+1] );
    }
}
/*****
***          c o c o C o n v e r t          ***
*****/
*
* Name:      cocoConvert
*
* Purpose:   Converts a single value
*
* Returns:   None
*
*****/
static void
cocoConvert ( card_t *cp, uint_t val )
{
    cocoCommand ( cp, COCO_CONVERT, val );
}
/*****
***          c o c o C o n v e r t T e s t          ***
*****/
*
* Name:      cocoConvertTest
*
* Purpose:   Converts a single value, read the converted value back
*            and compare to known value.

```

```

*
* Returns:    None
*
*****/
static void
cocoConvertTest ( card_t *cp, coco_convert_t *cv )
{
    register uint_t  res;

    /* convert the pixle value */
    cv->result = 0;          /* assume success */
    cocoCommand ( cp, COCO_CONVERT, cv->in );
    /* read the converted value and compare */
    res = cocoReadAmccFifo (cp);
    if ( res != cv->out )
        cv->result = 1;
}
/*****
***                c o c o D m a T o L u t s                ***
*****/
*
* Name:        cocoDmaToLuts
*
* Purpose:     DMAs data in user's buffer to one of Internal LUTs or
*              the External LUTs (identified by cmd param). The DMA is
*              done either in SINGLE or PROG (chained) mode depending
*              on current set up (cp->dmatype).
*
* Returns:     0 = Success, or errno
*
*****/
static int
cocoDmaToLuts( card_t *cp, coco_buf_t *cb, int cmd )
{
    register caddr_t  kvaddr, amcc_adr;
    register coco_dmapage_t  *dmaPg;
    register int      len, page_no, s, i, dmatype, err, tot_bytes;
    register int      cache_bytes;
    register uint_t   dmabits, *ib, olddmacmd, dmacmd;
    register uint_t   *cache_line;
    alenlist_t        addrList2;
    alenlist_t        addrList;
    size_t            p_size;
    alenaddr_t        p_addr;
    iopaddr_t         p_dmaPg;

```

```

amcc_adr = cp->amcc_adr;
dmacmd   = cp->dmacmd;
/* =====
 *      Scatter-Gather list preparation
 * ===== */
/* get the user's address and size */
len = cb->buf_size * sizeof(uint_t);

/* lock user pages into memory for DMA */
if ( cocoLockUser ( (caddr_t)cb->buf, len, B_WRITE) != 0 ) {
    cmn_err (CE_WARN, "cocoDmaToLuts: Cannot lock user pages");
    return (EFAULT);
}

/*
 * write back and invalidate the data for mem->board
 * adjust the address for current data cache size
 */
cache_line = (uint_t *)cb->buf;
cache_line = (uint_t *) ( ((uint_t)cache_line / COCO_CACHE_SIZE) *
                          COCO_CACHE_SIZE );
cache_bytes = len + ( (uint_t)cb->buf - (uint_t)cache_line);
dki_dcache_wbinval ( (caddr_t)cache_line, cache_bytes );
/* create scatter-gather list of user's buffer */
addrList2 = uvaddr_to_alenlist( (alenlist_t)NULL, (caddr_t)cb->buf,
                               (size_t)len);
if ( addrList2 == (alenlist_t)NULL ) {
    cmn_err (CE_WARN, "cocoDmaToLuts: cannot create alenlist");
    cocoUnlockUser ( (caddr_t)cb->buf, len, B_WRITE);
    return (EIO);
}
addrList = pciio_dmatrans_list ( cp->vhdl, cp->dev_desc,
                               /* addrList2, 0 ); */
                               addrList2, PCIIO_DMAMAP_BIGEND);
#if 0
page_no = alenlist_size ( addrList ); /* total pages to DMA */
#endif
page_no = cocoAlenlistSize ( addrList );
cp->page_no = page_no;
alenlist_done (addrList2);
/* initialize our DMA event semaphore */
cp->dmabits = 0;
olddmacmd = dmacmd;
dmatype = cp->dmatype;
initnsema ( &cp->dmawait, 0, "coco" );

```

```
cocoResetAmcc ( cp );
cocoPrepAmcc ( cp );
/* =====
 *   Set proper DMA flags
 * ===== */
dmacmd = DMAREG_FILL;
/* set proper LUT fill bit */
switch ( cmd ) {
    case COCO_BLOCK_FILL_RAMIL:
        dmacmd |= COCO_FILLRAMIL;
        #ifdef DEBUG
            printf (
                "cocoDmaToLuts: Writing %d bytes to RAMIL, %d pages\n",
                    len, page_no );
        #endif
        break;
    case COCO_BLOCK_FILL_RAMIH:
        dmacmd |= COCO_FILLRAMIH;
        #ifdef DEBUG
            printf (
                "cocoDmaToLuts: Writing %d bytes to RAMIH, %d pages\n",
                    len, page_no );
        #endif
        break;
    case COCO_BLOCK_FILL_RAMO:
        dmacmd |= COCO_FILLRAMO;
        #ifdef DEBUG
            printf (
                "cocoDmaToLuts: Writing %d bytes to RAMIO, %d pages\n",
                    len, page_no );
        #endif
        break;
    case COCO_BLOCK_FILL_RAML:
        dmacmd |= COCO_FILLRAML;
        #ifdef DEBUG
            printf (
                "cocoDmaToLuts: Writing %d bytes to RAML, %d pages\n",
                    len, page_no );
        #endif
        break;
}
cp->dmacmd = dmacmd;
/* =====
 *   Chained DMA
 * ===== */
```

```

if ( dmatype == DMA_PROG ) {
    #ifdef DEBUG
    printf ("cocoDmaToLuts: Chain Dma %d pages\n", page_no );
    #endif
    dmaPg = cocoMakeChain ( cp, addrList, page_no );
    if ( dmaPg == (coco_dmapage_t *)NULL ) {
        printf ("cocoDmaLut: Error creating chain list\n");
        cocoUnlockUser ( (caddr_t)cb->buf, len, B_WRITE );
        alenlist_done(addrList);
        return(EIO);
    }
    tot_bytes = page_no * sizeof(coco_dmapage_t);
    p_dmaPg = pciio_dmatrans_addr ( cp->vhdl, cp->dev_desc,
        kvtophys(dmaPg), tot_bytes,
        PCIIO_DMAMAP_BIGEND );

    #ifdef DEBUG
    cocoShowChain( "Lut Chain List", dmaPg );
    #endif
    /* write back cache for the chain list */
    dki_dcache_wbinval(dmaPg, tot_bytes );
    s = COCO_LOCK();
    err = 0;
    /* dma from memory -> board (selected lut buffer) */
    cocoStartProgDma ( cp, p_dmaPg, len, B_WRITE );
    cp->dmastat = DMA_LUT_WAIT;
    #ifdef DEBUG
    printf ("cocoDmaToLut: Waiting DMA Interrupt\n");
    #endif
    if ( SleepEvent(&cp->dmawait) != 0 ) {
        #ifdef DEBUG
        printf ("cocoDmaToLut: Interrupted\n");
        #endif
        err = EINTR;
    }
    else {
        #ifdef DEBUG
        printf ("cocoDmaToLuts: Woken up..all done\n");
        #endif

        err = 0;
    }
    /* we are done */
    cp->dmastat = DMA_IDLE;
    cp->dmabits = 0;
    cp->dmacmd = olddmacmd;
}

```

```

        COCO_UNLOCK(s);
        cocoUnlockUser ( (caddr_t)cb->buf, len, B_WRITE );
        alenlist_done(addrList);
        kmem_free ( dmaPg, page_no * sizeof(coco_dmapage_t) );
        return(err);
    }
    /* =====
    *      Single page DMA
    * ===== */
    for ( i = 0; i < page_no; i++ ) {
        /* get a page to DMA */
        if ( alenlist_get(addrList, NULL, NBPP,
            &p_addr, &p_size) != ALENLIST_SUCCESS ) {
            cmn_err (CE_WARN, "cocoDma: Bad scatter-gather");
            cocoUnlockUser ( (caddr_t)cb->buf, len, B_WRITE );
            alenlist_done(addrList);
            return(ENOMEM);
        }
        #ifdef DEBUG
        printf (
            "cocoDmaToLuts: Loop ...DMA page %d of %d (%d bytes), p_addr = 0x%x\n",
                i+1, page_no, p_size, p_addr );
        #endif
        s = COCO_LOCK();
        err = 0;
        cp->dmastat = DMA_LUT_WAIT;
        cocoStartSingleDma ( cp, p_addr, p_size, B_WRITE );
        #ifdef DEBUG
        printf ("cocoDmaToLut: Loop, waiting for Interrupt\n");
        #endif
        if ( SleepEvent(&cp->dmawait) != 0 ) {
            #ifdef DEBUG
            printf ("cocoDmaToLut: Interrupted...\n");
            #endif
            err = EINTR;
            COCO_UNLOCK(s);
            break;
        }
        else {
            COCO_UNLOCK(s);
            #ifdef DEBUG
            printf ("cocoDmaToLut: Loop ..Woken up\n");
            #endif
        }
    }
}

```

```

    }
    #ifdef DEBUG
    printf ("cocoDmaToLut: DMA is done\n");
    #endif
    /* we are done */
    cp->dmastat = DMA_IDLE;
    cp->dmabits = 0;
    cp->dmacmd = olddmacmd;
    cocoUnlockUser ( (caddr_t)cb->buf, len, B_WRITE );
    alenlist_done(addrList);
    return(err);
}

/*****
***          c o c o L o c k U s e r          ***
*****/
*
* Name:          cocoLockUser
*
* Purpose:       Given a user's address space pointer, it locks all user's
*                pages in memory in preparation for DMA.
*
* Returns:       0 = Success or errno
*
*****/
static int
cocoLockUser ( caddr_t user_buf, int len, int direction )
{
    /* lock pages in memory */
    if ( userdma( (caddr_t)user_buf, len, direction) == 0 )
        return (EFAULT);
    return (0);
}

/*****
***          c o c o U n l o c k U s e r          ***
*****/
*
* Name:          cocoUnlockUser
*
* Purpose:       Given a user's address, unlock all pages for that address
*
* Returns:       None.
*
*****/
static void

```

```

cocoUnlockUser ( caddr_t user_buf, int len, int direction )
{
    undma ( user_buf, len, direction );
}
/*****
***          c o c o P r e p A m c c          ***
*****/
*
* Name:          cocoPrepAmcc
*
* Purpose:       Prepares AMCC chip for DMA.
*
* Returns:       None.
*
*****/
static void
cocoPrepAmcc ( card_t *cp )
{
    register caddr_t  adr_amcc;
    register uint_t  mcsr, intcsr;
    register uint_t  tmp;
    adr_amcc = cp->amcc_adr;
    /* =====
    *          Prepare AMCC chip
    * ===== */
    /*
    * Prepare AMCC as follow:
    *   - Read Maibox registers to make sure they are empty
    *   - Set Amcc DMA count to 0
    *   - Enable Amcc Read/Write interrupts
    */
    tmp = Inp32(adr_amcc+AMCC_OP_REG_IMB4);          /* Read In mbox */
    tmp = Inp32(adr_amcc+AMCC_OP_REG_OMB4);          /* read Out mbox */
    Out32(adr_amcc+AMCC_OP_REG_MRTC, 0 );
    Out32(adr_amcc+AMCC_OP_REG_MWTC, 0 );
    Out32(adr_amcc+AMCC_OP_REG_INTCSR, AMCC_INTCSR_MASK );
}
/*****
***          c o c o R e s e t A m c c          ***
*****/
*
* Name:          cocoResetAmcc
*
* Purpose:       Resets Addon and Amcc Fifos
*

```

```

* Returns:    None.
*
*****/
static void
cocoResetAmcc ( card_t *cp )
{
    register caddr_t  adr_amcc;
    adr_amcc = cp->amcc_adr;
    /* Reset Amcc Fifos and Addon interface */
    Out32(adr_amcc+AMCC_OP_REG_MCSR, AMCC_RST_ADDON );
    Out32(adr_amcc+AMCC_OP_REG_MCSR, AMCC_RST_FIFOS );
}
/*****
***                c o c o S t a r t P r o g D m a                ***
*****/
*
* Name:        cocoStartProgDma
*
* Purpose:     Programs the board for Chained DMA and starts the Dma
*
* Returns:     None.
*
*****/
static void
cocoStartProgDma ( card_t *cp, iopaddr_t p_dmaPg, int tot_bytes, int rw )
{
    register caddr_t  adr_cfg, adr_norm, adr_amcc;
    register uint_t  adr_bits, enable_bits, dmacfg, dmacmd, dmabits;
    adr_cfg = cp->conf_adr;
    adr_norm = cp->norm_adr;
    adr_amcc = cp->amcc_adr;
    dmacfg = cp->dmacfg;
    dmacmd = cp->dmacmd;
    dmabits = cp->dmabits;
    cp->dmasize = tot_bytes;
#ifdef DEBUG
    cocoReport ( cp, "before cocoStartProgDma");
#endif
    /* clear eof marks */
    Out32(adr_cfg, dmacfg );
    /* =====
    * board -> memory
    * ===== */
    if ( rw == B_READ ) {
        /* set Amcc counters */

```

```

        Out32(adr_amcc+AMCC_OP_REG_MWTC, tot_bytes );
        Out32(adr_amcc+AMCC_OP_REG_MRTC, 0 );
        adr_bits    = DMAREG_RAMPRWR;
        enable_bits = DMAREG_INTPWEN | DMAREG_WEN;
        dmabits |= DMAREG_PWEN;
    }
    /* =====
    * memory -> board
    * ===== */
    else {
        /* set Amcc counters */
        Out32(adr_amcc+AMCC_OP_REG_MRTC, tot_bytes );
        Out32(adr_amcc+AMCC_OP_REG_MWTC, 0 );
        adr_bits    = DMAREG_RAMPRRD;
        enable_bits = DMAREG_INTPREN | DMAREG_REN;
        dmabits |= DMAREG_PREN;
    }
    /* enable chaining */
    dmacfg |= dmabits;
    Out32(adr_cfg, dmacfg );
    /* set address of chained list */
    Out32(adr_cfg, dmacfg | DMAREG_RAM | adr_bits );
    Out32(adr_norm, p_dmaPg);

    /* start the DMA */
    dmabits |= enable_bits | dmacmd;
    cp->dmabits = dmabits;
#ifdef DEBUG
    cocoReport ( cp, "after cocoStartProgDma");
#endif
    pciio_flush_buffers ( cp->vhdl );
    microtime ( &cp->start_time );
    Out32(adr_cfg, dmacfg | dmabits );
}
/*****
***          c o c o S t a r t S i n g l e D m a          ***
*****
*
* Name:          cocoStartSingleDma
*
* Purpose:      Programs the board for Single page DMA (read or write)
*
* Returns:      None.
*
*****/

```

```

static void
cocoStartSingleDma ( card_t *cp, alenaddr_t p_addr, size_t p_size, int rw )
{
    register caddr_t  adr_cfg, adr_norm, adr_amcc;
    register uint_t   adr_bits, enable_bits, dmacfg, dmacmd, dmabits, temp;
    register int      tot_words;
    adr_cfg = cp->conf_adr;
    adr_norm = cp->norm_adr;
    adr_amcc = cp->amcc_adr;
    dmacfg = cp->dmacfg;
    dmacmd = cp->dmacmd;
    dmabits = cp->dmabits;
    enable_bits = 0;
    cp->dmasize = p_size;
#ifdef DEBUG
    printf ("cocoStartSingleDma: Dma'ing %d bytes %s, p_addr = 0x%x\n",
           p_size, rw==B_READ ? "from board":"to board", p_addr );
#endif
    Out32(adr_cfg, dmacfg ); /* clear eof marks */
    /*
     * Select read/write count in DMA controller
     * and set the DMA size in 32-bit values
     */
    tot_words = (int)p_size/sizeof(uint_t);
    if ( rw == B_READ ) { /* board -> memory */
        Out32(adr_amcc+AMCC_OP_REG_MWTC, (uint_t)p_size );
        Out32(adr_amcc+AMCC_OP_REG_MRTC, 0xffffffff );
        /* Out32(adr_amcc+AMCC_OP_REG_MRTC, (uint_t)p_size ); */
        dmabits |= DMAREG_WEN;
        Out32(adr_cfg, dmacfg | DMAREG_WCNT);
    }
    else {
        Out32(adr_amcc+AMCC_OP_REG_MRTC, (uint_t)p_size );
        Out32(adr_amcc+AMCC_OP_REG_MWTC, 0xffffffff );
        /* Out32(adr_amcc+AMCC_OP_REG_MWTC, (uint_t)p_size ); */
        dmabits |= DMAREG_REN;
        Out32(adr_cfg, dmacfg | DMAREG_RCNT);
    }
    /* set transfer counts in words - -1 since counts down to 0xffff */
    tot_words--;
    Out32(adr_norm, (uint_t)tot_words);
    /* select Read/Write Address and set the DMA address */
    if ( rw == B_READ ) {
        Out32(adr_cfg, dmacfg | DMAREG_RAM | DMAREG_RAMWADR); /* b->m */
    }
}

```

```

else {
    Out32(adr_cfg, dmacfg | DMAREG_RAM | DMAREG_RAMRADR); /* m->b */
}
Out32(adr_norm, (uint_t)p_addr);
/* enable the right interrupt */
if ( rw == B_READ )
    dmabits |= DMAREG_INTWEN; /* board -> memory */
else dmabits |= DMAREG_INTREN; /* memory -> board */
dmabits |= cp->dmacmd;
cp->dmabits = dmabits;
#ifdef DEBUG
cocoReport ( cp, "cocoStartSingleDma exit");
#endif
pciio_flush_buffers ( cp->vhdl );
microtime ( &cp->start_time);
Out32(adr_cfg, dmacfg | dmabits );
}
/*****
***          c o c o M a k e C h a i n          ***
*****/
*
* Name:          cocoMakeCHain
*
* Purpose:       given an alenlist, creates chained list for Prog DMA.
*
* Returns:       NULL for error or address of chained list
*
*****/
static coco_dmapage_t *
cocoMakeChain ( card_t *cp, alenlist_t addrList, int page_no )
{
    register coco_dmapage_t * dmaPg;
    register int          i, tot_words;
    size_t                p_size;
    alenaddr_t            p_addr;
    iopaddr_t             pa;

    dmaPg = (coco_dmapage_t *)kmem_alloc (page_no * sizeof(coco_dmapage_t),
                                         KM_NOSLEEP | KM_PHYSCONTIG | KM_CACHEALIGN);
    if ( dmaPg == (coco_dmapage_t *)NULL ) {
        printf ("cocoMakeChain: Not enough mem for chained list\n");
        return( (coco_dmapage_t *)NULL);
    }
    /* fill the chained list with address-size values */
    for ( i = 0; i < page_no; i++ ) {

```

```

        if ( alenlist_get(addrList, NULL, NBPP,
            &p_addr, &p_size) != ALENLIST_SUCCESS ) {
            cmn_err (CE_WARN, "cocoMakeChain: Bad alenlist\n");
            return( (coco_dmapage_t *)NULL);
        }
        pa = pciio_dmatrans_addr ( cp->vhdl, cp->dev_desc,
            kvtophys(&dmapg[i+1]),
            sizeof(coco_dmapage_t),
            PCIIO_DMAMAP_BIGEND );

        dmapg[i].nextaddr = (paddr_t)pa;
        dmapg[i].addr = (paddr_t)p_addr;
        tot_words = (int)p_size/sizeof(uint_t);
        dmapg[i].size = tot_words - 1;
    }
    dmapg[i-1].size |= END_OF_CHAIN;
    return ( dmapg );
}
/*****
***              c o c o M a k e C h a i n R W              ***
*****/
*
* Name:          cocoMakeCHainRW
*
* Purpose:      Creates chain list for simultaneous read and write
*               but makes sure we read and write the same amount of bytes
*               in each entry
*
* Returns:      0 = Success, or errno
*
*****/
static int
cocoMakeChainRW ( card_t *cp,
                 coco_dmapage_t **w_dmapg,
                 coco_dmapage_t **r_dmapg )
{
    register coco_dmapage_t *wp, *rp;
    register alenaddr_t wp_resadr, rp_resadr;
    register iopaddr_t pa;
    register size_t wp_left, rp_left;
    register int i, tot_rchain, tot_wchain, w_page, r_page;
    register int tot_bytes, tot_words;
    size_t wp_size, rp_size;
    alenaddr_t wp_addr, rp_addr;
    w_page = cp->w_page_no;
    r_page = cp->r_page_no;

```

```

tot_rchain = r_page * CHAIN_FACTOR;
tot_wchain = w_page * CHAIN_FACTOR;
wp_resadr = 0;
rp_resadr = 0;
wp_left = 0;
rp_left = 0;
*w_dmaPg = (coco_dmapage_t *)NULL;
*r_dmaPg = (coco_dmapage_t *)NULL;
/* allocate chain list for Write */
wp = (coco_dmapage_t *)kmem_alloc (tot_wchain * sizeof(coco_dmapage_t),
                                KM_NOSLEEP | KM_PHYSCONTIG | KM_CACHEALIGN);
if ( wp == (coco_dmapage_t *)NULL ) {
    cmn_err (CE_WARN, "cocoMakeChainRW: Not enough memory");
    return (ENOMEM);
}
/* allocate chain list for Read */
rp = (coco_dmapage_t *)kmem_alloc (tot_rchain * sizeof(coco_dmapage_t),
                                KM_NOSLEEP | KM_PHYSCONTIG | KM_CACHEALIGN);
if ( rp == (coco_dmapage_t *)NULL ) {
    kmem_free (wp, tot_wchain * sizeof(coco_dmapage_t) );
    cmn_err (CE_WARN, "cocoMakeChainRW: Not enough memory");
    return (ENOMEM);
}
/* make sure we are at the top of the list */
alenlist_cursor_init ( cp->r_addrList, NULL, NULL );
alenlist_cursor_init ( cp->w_addrList, NULL, NULL );

#ifdef DEBUG
printf ("cocoMakeChainRW: ----- Started -----\n");
cocoShowAlenlist ( "Write Alenlist", cp->w_addrList );
cocoShowAlenlist ( "Read Alenlist", cp->r_addrList );
#endif
for ( i = 0; i++ ) {
    /* =====
    *      Write address and count
    * ===== */
    /* if no Write bytes left, get a new page */
    if ( wp_resadr == (alenaddr_t)NULL ) {
        /* we must have data for write */
        if ( w_page <= 0 ) {
            cmn_err (CE_WARN,
                "cocoMakeChainRW: Premature end of Write, w_page = %d, r_page = %d",
                    w_page, r_page );
            #ifdef DEBUG4
            printf ("%d Write pages, %d Read pages\n",

```

```

        tot_wchain / CHAIN_FACTOR,
        tot_rchain / CHAIN_FACTOR );
#ifdef DEBUG3
rp[i-1].size |= END_OF_CHAIN;
wp[i-1].size |= END_OF_CHAIN;
cocoShowChain ("Write Chain", wp );
cocoShowChain ("Read Chain", rp );
#endif
cocoShowAlenlist ( "Write Alenlist", cp->w_addrList );
cocoShowAlenlist ( "Read Alenlist", cp->r_addrList );
#endif
kmem_free (wp, tot_wchain * sizeof(coco_dmapage_t) );
kmem_free (rp, tot_rchain * sizeof(coco_dmapage_t) );
return(EIO);
}
if ( alenlist_get(cp->w_addrList, NULL, NBPP,
&wp_addr, &wp_size) != ALENLIST_SUCCESS ) {
cmn_err (CE_WARN,
"cocoMakeChainRW: Bad Write alenlist\n");
kmem_free (wp, tot_wchain * sizeof(coco_dmapage_t) );
kmem_free (rp, tot_rchain * sizeof(coco_dmapage_t) );
return(EIO);
}
w_page--;
}
/* some old Write bytes left - use those */
else {
#ifdef DEBUG
printf ("using %d old Write bytes left at 0x%x\n",
wp_left, wp_resadr );
#endif
wp_addr = wp_resadr;
wp_size = wp_left;
}
/* =====
*      Read address and count
* ===== */
/* if no Read bytes left, get a new page */
if ( rp_resadr == (alenaddr_t)NULL ) {
/* we must have data for read */
if ( r_page <= 0 ) {
cmn_err (CE_WARN,
"cocoMakeChainRW: Premature end of Read, r_page = %d, w_page = %d",
r_page, w_page );
#ifdef DEBUG4

```

```

printf ("%d Write pages, %d Read pages\n",
        tot_wchain / CHAIN_FACTOR,
        tot_rchain / CHAIN_FACTOR );
#ifdef DEBUG3
wp[i-1].size |= END_OF_CHAIN;
rp[i-1].size |= END_OF_CHAIN;
cocoShowChain ("Read Chain", rp );
cocoShowChain ("Write Chain", wp );
#endif
cocoShowAlenlist ( "Write Alenlist", cp->w_addrList );
cocoShowAlenlist ( "Read Alenlist", cp->r_addrList );
#endif
kmem_free (wp, tot_wchain * sizeof(coco_dmapage_t) );
kmem_free (rp, tot_rchain * sizeof(coco_dmapage_t) );
return(EIO);
}
if ( alenlist_get(cp->r_addrList, NULL, NBPP,
&rp_addr, &rp_size) != ALENLIST_SUCCESS ) {
cmn_err (CE_WARN,
"cocoMakeChainRW: Bad Read alenlist\n");
kmem_free (wp, tot_wchain * sizeof(coco_dmapage_t) );
kmem_free (rp, tot_rchain * sizeof(coco_dmapage_t) );
return(EIO);
}
r_page--;
}
/* some old Read bytes left - use those */
else {
#ifdef DEBUG
printf ("using %d old Read bytes left at 0x%x\n",
rp_left, rp_resadr );
#endif
rp_addr = rp_resadr;
rp_size = rp_left;
}
/* =====
* Adjust Read and Write counts
* ===== */
wp_resadr = 0;
rp_resadr = 0;
wp_left = 0;
rp_left = 0;
/* Writing and Reading the same amount ? no adjustment */
if ( wp_size == rp_size ) {
#ifdef DEBUG

```

```

        printf (
"Same byte count (%d) at read 0x%x and write 0x%x ..no adjustment\n",
        wp_size, rp_addr, wp_addr );
        #endif
        tot_bytes = wp_size;
    }
    /* Writing more than reading ? Adjust Write */
    if ( wp_size > rp_size ) {
        tot_bytes = rp_size;
        wp_resadr = (alenaddr_t)((int)wp_addr + tot_bytes);
        wp_left = wp_size - tot_bytes;
        #ifdef DEBUG
        printf ("Write adjusted ..%d bytes left at 0x%x\n",
            wp_left, wp_resadr );
        #endif
    }

    /* Reading more than write ? Adjust Read */
    if ( rp_size > wp_size ) {
        tot_bytes = wp_size;
        rp_resadr = (alenaddr_t)((int)rp_addr + tot_bytes);
        rp_left = rp_size - tot_bytes;
        #ifdef DEBUG
        printf ("Read adjusted ..%d bytes left at 0x%x\n",
            rp_left, rp_resadr );
        #endif
    }
    /* =====
    *      Make Read and Write Chain Entries
    * ===== */
    tot_words = (int)tot_bytes/sizeof(uint_t);
    tot_words--;
    /* Make Write Chain entry */
    pa = pciio_dmatrans_addr ( cp->vhdl, cp->dev_desc,
        kvtophys(&wp[i+1]),
        sizeof(coco_dmapage_t),
        PCIIO_DMAMAP_BIGEND );
    wp[i].nextaddr = (paddr_t)pa;
    wp[i].addr = (paddr_t)wp_addr;
    wp[i].size = tot_words;
    /* Make Write Chain entry */
    pa = pciio_dmatrans_addr ( cp->vhdl, cp->dev_desc,
        kvtophys(&rp[i+1]),
        sizeof(coco_dmapage_t),
        PCIIO_DMAMAP_BIGEND );

```

```

rp[i].nextaddr = (paddr_t)pa;
rp[i].addr = (paddr_t)rp_addr;
rp[i].size = tot_words;
/* end of Write chain list ? */
if ( (w_page <=0) && (wp_resadr == (alenaddr_t)NULL) )
    wp[i].size |= END_OF_CHAIN;
/* end of Read chain list ? */
if ( (r_page <=0) && (rp_resadr == (alenaddr_t)NULL) )
    rp[i].size |= END_OF_CHAIN;
/* end of loop ? */
if ( (rp[i].size & END_OF_CHAIN) &&
      (wp[i].size & END_OF_CHAIN) )
    break;
/* ran out of memory ? */
if ( (i >= tot_rchain) || (i >= tot_wchain) ) {
    cmn_err (CE_WARN,
            "cocoMakeChainRW: Ran out of Chain entries");
    kmem_free (wp, tot_wchain * sizeof(coco_dmapage_t) );
    kmem_free (rp, tot_rchain * sizeof(coco_dmapage_t) );
    return(EIO);
}

} /*** for ( i = 0;; i++ ) ***/
*w_dmaPg = wp;
*r_dmaPg = rp;
return(0);
}

/*****
***          c o c o T i m e O u t          ***
*****/
*
* Name:          cocoTimeout
*
* Purpose:       We timedout waiting for a read/write interrupt.
*
* Returns:       None.
*
*****/
static void
cocoTimeout ( card_t *cp )
{
    register uint_t intcsr;
    cmn_err (CE_NOTE, "cocoTimeout: Read/Write Timed out");
    alenlist_done ( cp->addrList );
}

```

```

        cp->addrList = 0;
#ifdef DEBUG
        cocoDumpAmcc(cp);
#endif
        bioerror ( cp->bp, ETIME);
        biodone (cp->bp);
    }
/*****
***          c o c o T i m e O u t 2          ***
*****/
*
* Name:          cocoTimeOut2
*
* Purpose:       We timedout waiting for a simultaneous read/write intr.
*
* Returns:       None.
*
*****/
static void
cocoTimeOut2 ( card_t *cp )
{
    cmn_err (CE_NOTE, "cocoTimeOut2: Sim. Read/Write Timed out");
#ifdef DEBUG
    cocoDumpAmcc(cp);
#endif
    cp->iostat = IO_TIME;
    WakeEvent(&cp->dmawait);
}
/*
~~~~~
~~~~~          D e b u g i n g   R o u t i n e s          ~~~~~
~~~~~
~~~~~
*/
/*****
***          c o c o I o c t l S t r          ***
*****/
*
* Name:          cocoIoctlStr
*
* Purpose:       return string version of an Ioctl command
*
* Returns:       pointer to string
*

```

```

*****/
static char *
cocoIoctlStr ( int cmd )
{
    switch ( cmd ) {
        case COCO_COMMAND: return ("Coco_Command");
        case COCO_RAW_READ_FIFO: return ("Raw_Read_Fifo");
        case COCO_RAW_WRITE_FIFO: return ("Raw_Write_Fifo");
        case COCO_RAW_READB_FIFO: return ("Raw_ReadB_Fifo");
        case COCO_RAW_WRITEB_FIFO: return ("Raw_WriteB_Fifo");
        case COCO_FIFO_TEST: return ("Fifo_Test");
        case COCO_SET_MODE: return ("Set_Mode");
        case COCO_READ_MODE: return ("Read_Mode");
        case COCO_RAW_READ_DMA: return ("Raw_Read_Dma");
        case COCO_RAW_WRITE_DMA: return ("Raw_Write_Dma");
        case COCO_READ_ADDR: return ("Read_Addr");
        case COCO_SET_ADDR: return ("Set_Addr");
        case COCO_DMAREGS_TEST: return ("DmaRegs_Test");
        case COCO_INTRAM_TEST: return ("IntRam_Test");
        case COCO_EXTRAM_TEST: return ("ExtRam_Test");
        case COCO_READ_RAMIL: return ("Read_RamIL");
        case COCO_READ_RAMIH: return ("Read_RamIH");
        case COCO_READ_RAMO: return ("Read_RamO");
        case COCO_READ_RAML: return ("Read_RamL");
        case COCO_FILL_RAMIL: return ("Fill_RamIL");
        case COCO_FILL_RAMIH: return ("Fill_RamIH");
        case COCO_FILL_RAMO: return ("Fill_RamO");
        case COCO_FILL_RAML: return ("Fill_RamL");
        case COCO_CONVERT_PIXLE: return ("Convert_Pixle");
        case COCO_CONVERT_TEST: return ("Convert_Test");
        case COCO_SET_SINGLE_DMA: return ("Set_Single_Dma");
        case COCO_SET_PROG_DMA: return ("Set_Prog_Dma");
        case COCO_BLOCK_FILL_RAMIL: return ("Block_Fill_RamIL");
        case COCO_BLOCK_FILL_RAMIH: return ("Block_Fill_RamIH");
        case COCO_BLOCK_FILL_RAML: return ("Block_Fill_RamL");
        case COCO_BLOCK_FILL_RAMO: return ("Block_Fill_RamO");
        case COCO_SETCMD_TRANSP: return ("SetCmd_Transp");
        case COCO_SETCMD_CONVERT: return ("SetCmd_Convert");
        case COCO_RESET: return ("Reset");
        case COCO_RW_BUF: return ("RW_Buff");
        case COCO_ISPCI: return ("Is_PCI");
        defaults: return ("Unknown");
    }
}
/*****

```

```

***                c o c o R e p o r t                ***
*****
*
* Name:           cocoReport
*
* Purpose:        Prints out the bit setting in dmacfg and other status
*
* Returns:        None.
*
*****/
static void
cocoReport ( card_t *cp, char *s )
{
    register uint_t  dmacfg, dmabits, dmacmd, ram;
    register int     dmastat, dmatype;
    dmacmd = cp->dmacmd;
    dmabits = cp->dmabits;
    dmastat = cp->dmastat;
    dmatype = cp->dmatype;
    printf ("%s: ", s );
    printf ("dmabits: " );
    if ( dmabits & DMAREG_WEN ) printf ("|Wen");
    if ( dmabits & DMAREG_REN ) printf ("|Ren");
    if ( dmabits & DMAREG_PREN) printf ("|Pren");
    if ( dmabits & DMAREG_PWEN) printf ("|Pwen");
    if ( dmabits & DMAREG_FILL) printf ("|Reg_Fill");
    if ( dmabits & DMAREG_INTWEN) printf ("|IntWen");
    if ( dmabits & DMAREG_INTREN) printf ("|IntRen");
    if ( dmabits & DMAREG_INTPWEN ) printf ("|IntPwen");
    if ( dmabits & DMAREG_INTPREN ) printf ("|IntPren");
    ram = dmabits & 0x0f000000;
    if ( ram == COCO_FILLRAML) printf ("|Fill_RamL");
    if ( ram == COCO_FILLRAMIL) printf ("|Fill_RamIL");
    if ( ram == COCO_FILLRAMIH) printf ("|Fill_RamIH");
    if ( ram == COCO_FILLRAMO) printf ("|Fill_RamO");
    if ( ram == COCO_TRANSP)   printf ("|Transp");
    if ( ram == COCO_CONVERT) printf ("|Convrt");
    printf (" dmastat: ", dmastat);
    if ( dmastat == DMA_IDLE )      printf ("Idle");
    if ( dmastat == DMA_LUT_WAIT ) printf ("Dma_Lut_Wait");
    if ( dmastat == DMA_READ_WAIT ) printf ("Dma_Read_Wait");
    if ( dmastat == DMA_WRITE_WAIT ) printf ("Dma_Write_Wait");
    if ( dmastat == DMA_RW_WAIT )  printf ("Dma_RW_Wait");
    printf ("\n");
}

```

```

/*****
***          c o c o S h o w C h a i n          ***
*****
*
* Name:      cocoShowChain
*
* Purpose:   Displays contents of a chain list
*
* Returns:   None.
*
*****/
static void
cocoShowChain( char *title, coco_dmaPage_t *dmaPg )
{
    register int i, tot_bytes, tot_words;
    register long all_bytes;
    printf ("--- %s ---\n", title );
    all_bytes = 0;
    for (i = 0; i++ ) {
        tot_words = dmaPg[i].size;
        tot_bytes = ( (tot_words &~END_OF_CHAIN) +1) * sizeof(uint_t);
        all_bytes += tot_bytes;
        printf ("addr = 0x%x, size = %d next addr = 0x%x %s\n",
            dmaPg[i].addr, tot_bytes, dmaPg[i].nextaddr,
            tot_words & END_OF_CHAIN ? "[ End ]:" : " " );
        if ( tot_words & END_OF_CHAIN )
            break;
    }
    printf ("--- Total of %d entries, %d bytes ---\n\n", i, all_bytes );
}
/*****
***          c o c o D u m p A m c c          ***
*****
*
* Name:      cocoDumpAmcc
*
* Purpose:   Dumps Amcc registers for debugging purpose.
*
* Returns:   None.
*
*****/
static void
cocoDumpAmcc ( card_t *cp )
{
    register caddr_t amcc_adr, norm_adr, cfg_adr;

```

```

register uint_t  dmacfg, xil_stat, war, wcnt, rar, rcnt, mbeif;
amcc_adr = cp->amcc_adr;
norm_adr = cp->norm_adr;
cfg_adr  = cp->conf_adr;
dmacfg   = cp->dmacfg;
/* disable any Dma and read in Xilinx status */
Out32(cfg_adr, dmacfg | DMAREG_STAT );
xil_stat = Inp32(norm_adr);
war      = Inp32(amcc_adr+AMCC_OP_REG_MWAR);
wcnt     = Inp32(amcc_adr+AMCC_OP_REG_MWTC);
rar      = Inp32(amcc_adr+AMCC_OP_REG_MRAR);
rcnt     = Inp32(amcc_adr+AMCC_OP_REG_MRTC);
mbeif    = Inp32(amcc_adr+AMCC_OP_REG_MBEF);

wcnt &= 0x01ffffff;
rcnt &= 0x01ffffff;
printf ("***  cocoDumpAmcc ***\n");
printf ("WAR = 0x%x      WTC = %d [ 0x%x ]\n", war, wcnt, wcnt );
printf ("RAR = 0x%x      RTC = %d [ 0x%x ]\n", rar, rcnt, rcnt );
printf ("MBEF = 0x%x  Xilinx = 0x%x\n", mbeif, xil_stat );
}
/*****
***              c o c o S h o w A l e n l i s t              ***
*****
*
* Name:          cocoShowAlenlist
*
* Purpose:       Displays the contents of a given Alenlist
*
* Returns:       None.
*
*****/
static void
cocoShowAlenlist ( caddr_t title, alenlist_t al )
{
    size_t          size;
    long            tot_bytes;
    alenaddr_t      addr;
    register int    count, i;
    /* reset the cursor for the alenlist */
    alenlist_cursor_init ( al, NULL, NULL );
    printf ("cocoShowAlenlist: --- %s ---\n", title );
    tot_bytes = 0;
    count     = 0;
    for ( ;; ) {

```

```

        if ( alenlist_get(al, NULL, NBPP, &addr, &size) !=
            ALENLIST_SUCCESS ) {
            break;
        }
        printf ("addr = 0x%x, size = %d ...[%d]\n",addr, size, count );
        tot_bytes += size;
        count++;
    }
    printf ("--- Total of %d bytes [ %d entries ]---\n\n",
        tot_bytes, count );
    /* reset the cursor now */
    alenlist_cursor_init ( al, NULL, NULL );
}
/*****
***          c o c o A l e n l i s t S i z e          ***
*****/
*
* Name:          cocoAlenlistSize
*
* Purpose:      Returns number of pairs in a given alenlist.
*
* Returns:      Number of address/size entries
*
*****/
static int
cocoAlenlistSize ( alenlist_t al )
{
    register int      count;
    size_t           size;
    alenaddr_t       addr;
    alenlist_cursor_init ( al, NULL, NULL );
    count = 0;
    for ( ;; ) {
        if ( alenlist_get(al, NULL, NBPP, &addr, &size) !=
            ALENLIST_SUCCESS ) {
            break;
        }
        count++;
    }
    alenlist_cursor_init ( al, NULL, NULL );
    return (count);
}
#ifdef DEBUGTIME
/*****
***          c o c o R e p o r t T i m e          ***
*****/

```

```

*****
*
* Name:      cocoReportTime
*
* Purpose:   Displays start, intr and end time of Dma
*
* Returns:   None.
*
*****/
static void
cocoReportTime ( caddr_t title, card_t *cp, int final )
{
    register long   s_sec, s_milsec, s_micsec;
    register long   i_sec, i_milsec, i_micsec;
    register long   e_sec, e_milsec, e_micsec;
    register struct timeval *tv;
    struct timeval  dtv;
    if ( final == 0 ) {
        tv = &cp->start_time;
        s_sec   = tv->tv_sec & 0x000000ff;
        s_milsec = tv->tv_usec / 1000;
        s_micsec = tv->tv_usec % 1000;
        tv = &cp->intr_time;
        i_sec   = tv->tv_sec & 0x000000ff;
        i_milsec = tv->tv_usec / 1000;
        i_micsec = tv->tv_usec % 1000;
        cocoDiffTime( &cp->start_time, &cp->intr_time, &dtv );
        tv = &dtv;
        e_sec   = tv->tv_sec & 0x000000ff;
        e_milsec = tv->tv_usec / 1000;
        e_micsec = tv->tv_usec % 1000;
        printf (
"%s: %s, %d bytes, ..Start: %d.%d.%d, ..Intr: %d.%d.%d, ..Diff: %d.%d.%d\n",
            title, cp->dmatype == DMA_PROG ? "Chain":"Single", cp->dmasize,
            s_sec, s_milsec, s_micsec,
            i_sec, i_milsec, i_micsec,
            e_sec, e_milsec, e_micsec );

        return;
    }
    tv = &cp->call_time;
    s_sec   = tv->tv_sec & 0x000000ff;
    s_milsec = tv->tv_usec / 1000;
    s_micsec = tv->tv_usec % 1000;
    tv = &cp->ret_time;
}

```

```

        i_sec    = tv->tv_sec & 0x000000ff;
        i_milsec = tv->tv_usec / 1000;
        i_micsec = tv->tv_usec % 1000;
        cocoDiffTime ( &cp->call_time, &cp->ret_time, &dtv );
        tv = &dtv;
        e_sec    = tv->tv_sec & 0x000000ff;
        e_milsec = tv->tv_usec / 1000;
        e_micsec = tv->tv_usec % 1000;
        printf ("%s: Call at %d.%d.%d, ...Ret at %d.%d.%d, ..Diff: %d.%d.%d\n",
                title,
                s_sec, s_milsec, s_micsec,
                i_sec, i_milsec, i_micsec,
                e_sec, e_milsec, e_micsec );
    }
}
/*****
***                c o c o D i f f T i m e                ***
****
*
* Name:          cocoDiffTime
*
* Purpose:      Given two timeval struct, it calculates the difference
*
* Returns:     None.
*
*****/
static void
cocoDiffTime ( struct timeval *st, struct timeval *et, struct timeval *dt )
{
    register long  s_sec, s_milsec, s_micsec;
    register long  e_sec, e_milsec, e_micsec;
    register long  s_totmic, e_totmic, diff_mic;
    s_sec    = st->tv_sec & 0x000000ff;
    s_totmic = (s_sec * 1000000) + st->tv_usec;
    e_sec    = et->tv_sec & 0x000000ff;
    e_totmic = (e_sec * 1000000) + et->tv_usec;
    diff_mic = e_totmic - s_totmic;
    dt->tv_sec = diff_mic / 1000000;
    dt->tv_usec = diff_mic % 1000000;
}
#endif
#define GOOD  (PG_M | PG_G | PG_SV | PG_VR)
static void
cocoDebug ( coco_rw_t *rw )
{
    register int  tot_bytes, tot_wrong, i;

```

```

register uint_t      *kvi;
/* =====
 *      Scatter-Gather list for Write buffer (mem -> board)
 * ===== */
tot_bytes = rw->buf_size * sizeof(uint_t);
/* lock user's pages in memory for DMA */
if ( cocoLockUser ( (caddr_t)rw->w_buf, tot_bytes, B_WRITE) != 0 ) {
    cmn_err (CE_WARN, "cocoDebug: Cannot lock user's pages");
    return;
}
kvi = (uint_t *)maputokv ( (caddr_t)rw->w_buf, tot_bytes,
                          /* PG_UNCACHED | GOOD ); */
                          pte_cachebits() | GOOD );
/* check memory before the dki_dcache */
if ( kvi ) {
    tot_wrong = 0;
    for ( i = 0; i < rw->buf_size; i++ ) {
        if ( kvi[i] != 0x0 ) {
            if ( tot_wrong %100 == 0 ) {
                printf (
                    "cocoDebug: kvi[%d] = 0x%x at 0x%x kv and 0x%x user\n",
                    i, kvi[i], kvi+i, rw->w_buf+i );
            }
            tot_wrong++;
        }
    }
    if ( tot_wrong )
        printf (
            "cocoDebug BEFORE: --- %d pixles were wrong in w_buf (0x%x)\n\n",
            tot_wrong, rw->w_buf );
    else printf ("cocoDebug BEFORE: Ok\n\n");
}
/* write back and invalidate the data for mem->board */
dki_dcache_wbinval ( (caddr_t)rw->w_buf , tot_bytes );
/* check after dki_dcache */
if ( kvi ) {
    tot_wrong = 0;
    for ( i = 0; i < rw->buf_size; i++ ) {
        if ( kvi[i] != 0x0 ) {
            if ( tot_wrong %100 == 0 ) {
                printf (
                    "cocoDebug: kvi[%d] = 0x%x at 0x%x kv and 0x%x user\n",
                    i, kvi[i], kvi+i, rw->w_buf+i );
            }
            tot_wrong++;
        }
    }
}

```

```
        }
    }
    if ( tot_wrong )
        printf (
            "cocoDebug AFTER: --- %d pixles were wrong in w_buf (0x%x)\n\n",
                tot_wrong, rw->w_buf );
    else printf ("cocoDebug AFTER: Ok\n\n");
}
unmaputokv ( (caddr_t)kvi, tot_bytes );
cocoUnlockUser ( (caddr_t)rw->w_buf, tot_bytes, B_WRITE );
}
```

**PART NINE**

## **STREAMS Drivers**

### **Chapter 16: STREAMS Drivers**

How STREAMS drivers are integrated into the IRIX system.



---

## STREAMS Drivers

The IRIX implementation of STREAMS drivers is intended to be compatible with the multiprocessor implementation of STREAMS in UNIX version SVR4.2.

STREAMS programming in SVR4.2 is documented in *STREAMS Modules and Drivers, UNIX SVR4.2*. That book contains detailed discussion and many examples of STREAMS programming.

References in this chapter to *STREAMS Modules and Drivers* are to the edition copyright 1992 by UNIX System Laboratories, published by UNIX Press/Prentice-Hall, and bearing ISBN 0-13-066879. If you are using an earlier edition, you should upgrade it. If you have a later edition, you may have to interpret references carefully.

This chapter contains the following major sections:

- “Driver Exported Names” on page 500 summarizes the public names and functions that a STREAMS driver must export.
- “Building and Debugging” on page 504 describes the ways that building a STREAMS driver are like and unlike other kernel-level drivers.
- “Special Considerations for Multiprocessing” on page 505 describes the methods you must use to work with the multi-threaded STREAMS monitor.
- “Special Considerations for IRIX” on page 507 details the points at which IRIX differs from the SVR4 STREAMS environment.
- “Summary of Standard STREAMS Functions” on page 512 lists the available kernel functions used by STREAMS drivers.
- “STREAMS Modules for X Input Devices” on page 514 describes the use of configuration files for special input devices used by the X display manager.

## Driver Exported Names

A STREAMS driver or module must define certain public names for use by *lboot*, as described in “Summary of Driver Structure” on page 140. Only one of these names, the info structure, is unique to a STREAMS driver or module; all the others are also defined by kernel-level device drivers.

The public names all begin with a prefix (see “Driver Name Prefix” on page 140); the same prefix is specified in the configuration file (see “Describing the Driver in /var/sysgen/master.d” on page 235).

### Streamtab Structure

A STREAMS driver or module must provide a global *streamtab* structure containing pointers to the *qinit* structures for its read and write queues. These structures in turn point to required *module\_info* structures. The name of the streamtab is *pxinfo*.

### Driver Flag Constant

A STREAMS driver or module should provide a driver flag constant containing either 0 or the flag *D\_MP*. (See “Driver Flag Constant” on page 145 and “Flag *D\_MP*” on page 145). The name of the constant is *pxdevflag*.

**Note:** A driver or module that does not export *pxdevflag* is assumed to use SVR3 calling conventions at its *pxopen()* and *pxclose()* entry points. However, this support will be withdrawn in a release of IRIX in the very near term. If you are porting a STREAMS driver or module to IRIX you are urged to make sure it uses SVR4 conventions and exports a *pxdevflag* containing at least 0.

## Initialization Entry Points

A STREAMS driver or module can define an entry point *pfxinit()*, or an entry point *pfxstart()*, or both. These entry points will be called during boot if the driver or module is included in the kernel, or when the driver or module is loaded if it is loadable. The operation of these entry points is the same as for device drivers (see “Initialization Entry Points” on page 147).

Many STREAMS drivers perform all initialization at open time, and have no *pfxinit()* or *pfxstart()* entry points. Many STREAMS modules perform initialization when they receive the `I_PUSH` ioctl message.

## Entry Point *open()*

A STREAMS driver (but not module) must export a *pfxopen()* entry point. The argument list for a STREAMS driver’s open differs from that of a device driver. The prototype for a STREAMS *pfxopen()* entry point is:

```
int
pfxopen(queue_t *q, dev_t *devp, int oflag, int sflag, cred_t *crp);
```

The argument values are

- \*q*            Pointer to the *queue* structure being opened.
- \*devp*        Pointer to a *dev\_t* value from which you can extract both the major and minor device numbers.
- oflag*        Flag bits specifying user mode options on the **open()** call.
- sflag*        Flag bits specifying the type of STREAM open: driver, module or clone.
- \*crp*        Pointer to a *cred\_t* object—an opaque structure for use in authentication.

The *pfxopen()* entry point is a public name. In addition a pointer to it must be defined in the *qinit* structure for the read queue.

### Entry Point `close()`

A STREAMS driver (but not module) must export a `pfxclose()` entry point. The argument list for a STREAMS driver's close differs from that of a device driver. The prototype for a STREAMS `pfxclose()` entry point is:

```
int
pfxclose(queue_t *q, int oflag, cred_t *crp);
```

The argument values are the same as passed to `pfxopen()`. The `pfxclose()` entry point is a public name. In addition a pointer to it must be defined in the `qinit` structure for the read queue.

### Put Functions `wput()` and `rput()`

Every STREAMS driver and module must define a `put()` function to handle messages as they are delivered to a queue.

The prototype of a `put()` function is as follows:

```
int
name(queue_t *q, mblk_t *mp);
```

Because the `put()` function for a given queue is addressed from the associated `qinit` structure, there is no requirement that the `put()` function be a public name, and no requirement that it begin with the prefix string. The `put()` function for the write queue, which handles messages moving “downstream” from the user process toward the driver, is conventionally called the `wput()` function. All write queues need a `wput()` function.

The `put()` function for the read queue, which handles messages moving “upstream” from the driver toward the user process, is conventionally called the `rput()` function. In some cases the `rput()` function is not required, for example in a driver where all upstream messages are generated by an interrupt handler.

Typically, a **put()** function decides what to do by switching on the message type value from `mp->b_datap->db_type`. A **put** routine must do at least one of the following:

- Process the message, if immediate processing is required, consuming the message or transforming it.
- Pass the original or processed message to the next component in the stream by calling the **putnext()** function (see the `putnext(D3)` reference page).
- Queue the message for deferred processing by the service routine with the **putq()** function (see the `putq(D3)` reference page).

When all processing is deferred to the service function, the address of the kernel function **putq()** can be given as a queue's **put()** function.

In a multiprocessor, a **put()** function can be called concurrently with user-level code, and concurrently with another **put()** function for the same or a different queue. A service function for the same or different queue can also be executing concurrently.

### Service Functions **rsrv()** and **wsrv()**

When a STREAMS driver defers message processing by setting the kernel function **putq()** address as the driver's **put()** function, the queue must also define a service function **srv()**.

Because the **srv()** function for a given queue is addressed from the associated *qinit* structure, there is no requirement that the **srv()** function be a public name, and no requirement that it begin with the prefix string.

The prototype of a **svr()** function is as follows:

```
int
name(queue_t *q);
```

The **srv()** function for the write queue, which handles messages moving “downstream” from the user process toward the driver, is conventionally called the **wsrv()** function. The **srv()** function for the read queue, which handles messages moving “upstream” from the driver toward the user process, is conventionally called the **rsrv()** function.

An **srv()** function is called by the STREAMS monitor to deal with queued messages. It is called at a time chosen by the monitor, not necessarily related to any call to the **put()** function for the same queue. In a multiprocessor, only one instance of **srv()** is called per queue at any time. However, one or more instances of the **put()** function could execute concurrently with the **srv()** function—so any data that is used in common by **put()** and **srv()** must be protected with a lock (see “Waiting and Mutual Exclusion” on page 206). User-level code can also execute concurrently with a service function.

The service function is expected to dispose of all queued messages through one of the following actions:

- Consuming and freeing the message.
- Passing the message on to the following queue using **putnext()** (see the **putnext(D3)** reference page).
- Replacing the message on the same queue using **putbq()** for processing later (see the **putbq(D3)** reference page).

The service function implements flow control (which the **put()** function cannot do). Before applying **putnext()**, the service function calls a flow control function such as **canputnext()** to find out if the following queue can accept a message. If the following queue cannot accept a message, the service function replaces the message with **putbq()** and exits.

A STREAMS module or driver that is not multiprocessor-aware (lacks **D\_MP** in its **pf\_xdevflags**) uses one set of functions for flow control (see the **canput(D3)** and **bcanputnext(D3)** reference pages), while one that is multiprocessor-aware uses a different set (see **canputnext(D3)** and **bcanputnext(D3)** ).

## Building and Debugging

A STREAMS driver or module is a kernel module and is compiled using the same compiler options as any driver (see “Compiling and Linking” on page 230).

You configure each STREAMS driver or module as part of the IRIX kernel by:

- Placing the executable module in */var/sysgen/boot*
- Writing a descriptive file and placing it in */var/sysgen/master.d* (see “Describing the Driver in */var/sysgen/master.d*” on page 235)
- Placing a USE or INCLUDE line in */var/sysgen/system* (see “Configuring a Kernel” on page 238)

When a STREAMS driver or module is loadable, you specify the appropriate options in the descriptive file (see “Master File for Loadable Drivers” on page 240). You can configure a STREAMS driver or module to be autoregistered and loaded automatically (see “Registration” on page 241). Alternatively, you can require a STREAMS driver or module to be loaded manually using the *ml* command (see “Loading” on page 241).

When you have configured a debugging kernel (see “Preparing the System for Debugging” on page 245), the symbols of a STREAMS driver or module are available for display. You can set breakpoints using *symmon* (see “Using *symmon*” on page 254). You can display symbols using *symmon* or *idbg* (see “Using *idbg*” on page 264). In particular, *idbg* has built-in support for displaying the contents of structures used by a STREAMS module or driver (see “Commands to Display STREAMS Structures” on page 270).

## Special Considerations for Multiprocessing

In IRIX releases prior to 6.2, the STREAMS monitor was single-threaded, so that only one **put()** or **srv()** function in the entire system could execute at any time. That one **put()** or **srv()** function might execute concurrently with user-level code, but no two STREAMS functions could execute concurrently.

Beginning with IRIX 6.2, the STREAMS monitor is multi-threaded. Depending on the version of IRIX and on the number of CPUs in the system, the following functions can run concurrently in any combination: one **srv()** function for each queue; any number of **put()** functions for each queue; and one or more user processes. For general discussion of the consequences, see “Planning for Multiprocessor Use” on page 174.

In the multithreaded monitor, when a module or driver calls **putq()** or **qenable()**, the service function for the enabled queue can begin to execute at any time. It can begin execution before **putq()** or **qenable()** call has returned, and can run concurrently with the module or driver that enabled the queue.

The STREAMS monitor runs concurrently with interrupt handling. For this reason, the interrupt handler of a STREAMS driver must take an extra step before it performs any STREAMS-related processing such as **allocb()**, **putq()**, or **qenable()**. The IRIX-unique functions provided for this purpose are summarized in Table 16-1.

**Table 16-1** Multiprocessing STREAMS Functions

Name	Can Sleep?	Summary
<code>streams_interrupt(D3)</code>	N	Synchronize interrupt-level function with STREAMS mechanism.
<code>STREAMS_TIMEOUT(D3)</code>	N	Synchronize timeout with STREAMS mechanism.

Suppose that the interrupt handler of a STREAMS driver needs to add a message to the read queue with **putq()**. It cannot simply call that function, since the STREAMS monitor might be using the queue at the same time in another CPU. The driver must define a function in which the **putq()** call is written. The name of this function and the pointer to the queue are passed to **streams\_interrupt()**. As soon as possible, **streams\_interrupt()** gets control of the queue and executes the passed function.

A callback function scheduled using **itimeout()** and similar functions (see “Waiting for Time to Pass” on page 216) must also be synchronized with the STREAMS monitor.

Suppose that a STREAMS driver or module needs to schedule a function to execute at a later time. (In a nonSTREAMS driver the function would be scheduled with **itimeout()**.) In the time-delayed function is a call to **qenable()**. That call cannot be executed freely whenever the interval expires, because the state of the STREAMS monitor is not known at that time.

The `STREAMS_TIMEOUT` macros provide a solution. Like **itimeout()**, it schedules a function to be executed at a later time. However, it defers calling the function until the function is synchronized with the STREAMS monitor, so that it can execute calls such as **qenable()**.

## Special Considerations for IRIX

While IRIX is largely compatible with UNIX SVR4.2, there are points of difference in the implementation of IRIX that have to be reflected in the design of a STREAMS driver or module. This topic lists points at which the contents of *STREAMS Modules and Drivers, UNIX SVR4.2* is not a correct description of IRIX and STREAMS use within IRIX.

### Extension of Poll and Select

Under IRIX, the `poll()` system function is not limited to testing STREAMS, but can be applied to file descriptors of all types (see the `poll(2)` and `select(2)` reference pages). In addition the `select()` function can be applied to STREAMS file descriptors. You may want to note this under the heading “STREAMS System Calls” in Chapter 2 of *STREAMS Modules and Drivers, UNIX SVR4.2*.

### Support for Pipes

IRIX supports two kinds of pipes with different semantics, as described in the `pipe(2)` reference page. The default type of pipe is compatible with UNIX SVR3, and does not conform to the description in Chapter 2 of *STREAMS Modules and Drivers, UNIX SVR4.2* under the heading “Creating a STREAMS-based Pipe.”

The SVR4 pipe semantics are enabled on a system-wide basis by using the `sysctl` command to set the tuning parameter `svr3pipe` to 0. First test the configuration as shown in Example 16-1.

#### Example 16-1 Testing Pipe Configuration

```
# sysctl | grep svr3pipe
svr3pipe = 1 (0x1)
```

## Service Scheduling

At two points in *STREAMS Modules and Drivers, UNIX SVR4.2* (Under “Service Procedure” in Chapter 4 and under “Message Processing” in Chapter 5), the book explicitly says that in a uniprocessor, enabled service functions are always executed before returning to user-level processing. This promise is not supported by IRIX. In both uniprocessors and multiprocessors, user-level processes can potentially execute after a service function is enabled and before it executes.

## Supplied STREAMS Modules

*STREAMS Modules and Drivers, UNIX SVR4.2*, Chapter 4, refers to some example STREAMS drivers named CHARPROC, CANONPROC, and ASCEBC. These examples are not supplied with IRIX.

The following STREAMS-based modules are supplied with IRIX. You can read their reference pages in volume 7:

- alp(7)            Algorithm pool management module.
- clone(7)        Clone-open driver; see “Support for CLONE Drivers” on page 510.
- connld(7)       Line discipline for unique stream connections.
- kbd(7)          Generalized string translation module.
- log(7)          Interface to STREAMS error logging and event tracing.
- sad(7)          STREAMS Administrative Driver.
- streamio(7)    STREAMS ioctl commands.
- timod(7)        Transport Interface cooperating STREAMS module.
- tirdwr(7)      Transport Interface read/write interface STREAMS module.
- tsd(7)          TELNET server protocol STREAMS device.

## No #idefs

Chapter 4 of *STREAMS Modules and Drivers, UNIX SVR4.2* refers in a note to the use of the #idef and a transition period for SVR3-compatible drivers. None of this material is relevant to IRIX. IRIX is SVR4-compatible, with no special provision for SVR3 drivers.

## Different I/O Hardware Model

Chapter 5 of *STREAMS Modules and Drivers, UNIX SVR4.2* discusses the use of memory-mapped hardware and of Dual-Access RAM (DARAM). None of these considerations are relevant in a MIPS processor. The MIPS I/O model is discussed in Chapter 1, “Physical and Virtual Memory.”

## Different Network Model

Chapter 10 of *STREAMS Modules and Drivers, UNIX SVR4.2* describes the TPI interface model. This model is supported in IRIX. When an application uses the TLI library functions such as `t_open()`, the library uses IRIX-provided TPI STREAMS modules which implement the protocol described in chapter 10.

Chapter 11 of *STREAMS Modules and Drivers, UNIX SVR4.2* describes the Data Link Provider Interface (DLPI) as implemented using STREAMS facilities.

The IRIX networking support is not STREAMS-based, but rather is based on BSD *ifnet* architecture. This is discussed in Chapter 14, “Network Device Drivers.” The IRIX network support includes DLPI support as an add-on feature to the *ifnet* driver interface. If you are porting a network device driver to IRIX, it is better to convert it to the *ifnet* interface. You can install a DLPI-based network device driver, but only other STREAMS modules could use it—there would be no connection to the rest of the IRIX networking system.

## Support for CLONE Drivers

*STREAMS Modules and Drivers, UNIX SVR4.2* discusses CLONE drivers; that is, STREAMS drivers that generate a new minor device number for each open. Refer to Chapter 3, “The CLONE Driver,” and to Chapter 8, “Cloning.” Clone opens and the clone driver are implemented under IRIX; this section clarifies the discussion in the SVR4 manual.

The essence of cloned access to a STREAMS driver is that the user process is indifferent to the minor device number, and simply wants to open a stream from this driver. A cloned stream is created using the following steps:

1. Recognize that the process calling **open()** is indifferent to the minor device number and simply wants cloned access.
2. Choose an unused minor device number from the set of minor numbers the driver supports.
3. Construct a new device number *dev\_t* value based on the chosen minor number, and assign it to the argument passed to *pfxopen()*.

### Using the CLONE Driver

The IRIX-supplied clone driver automates some of these steps for your driver. In order to use it, prepare a device special file with these characteristics:

- A device name that is related to the actual device name
- The major device number (10 decimal) that specifies the clone driver
- A minor device number equal to the major number of the actual driver

You can view the descriptive file for the clone driver in */var/sysgen/master.d/clone*. This file sets its major number (10) and states that it is not loadable. Although the clone driver is not specifically configured in the */var/sysgen/system/irix.sm* file, it is included in any kernel because it is listed as a dependency in the descriptive file of several other drivers (use *fgrep clone /var/sysgen/master.d/\** to see which drivers depend on it; and see “Listing Dependencies” on page 236). You can specify it as a dependency in the same way, if your driver depends on it.

When a user process opens a device special file with the major number of the clone driver, the kernel naturally calls the clone driver’s open entry point. The clone driver verifies that the minor number passed is the major number of an existing, STREAMS driver. (If it is not, the clone driver returns ENXIO).

The clone driver sets up the *qinit* structure appropriately for the target driver's queue and calls that driver's *pfxopen()* entry point, passing the CLONEOPEN flag in the *sflag* argument (see "Entry Point *open()*" on page 501).

### Recognizing a Clone Request Independently

It is not essential to use the clone driver. You can instead designate a particular minor device number to stand for "clone open." You prepare a device special file with these characteristics:

- A device name related to the actual device name
- The major number of your driver
- Some minor number you define to mean "clone open"

When a user process opens this device special file, the kernel calls the *pfxopen()* entry point of your driver. It does not pass the CLONEOPEN flag in *sflag*, but your driver can recognize a request for a clone open based on the minor device number.

### Responding to a Clone Request

In response to a clone request coming from either of the two methods described, your *pfxopen()* entry point must select an unused minor device number. (If no minor number is available, return EBUSY.)

Text in Chapter 3 of *STREAMS Modules and Drivers, UNIX SVR4.2* seems to suggest that your driver should scan through the kernel's *cdevsw* table to find an unused minor number (see "Kernel Switch Tables" on page 141). Under IRIX, the *cdevsw* table is not accessible to drivers. The reason is that the table layout differs between 32-bit and 64-bit kernels, and can change between releases. Instead, your driver must know the minor numbers that it supports, and must know which ones are currently in use.

**Tip:** You can design your driver so that the number of supported devices is specified in the descriptive file in */var/sysgen/master.d*, and passed in to the driver through that descriptive file (see "Variables Section" on page 237). Your driver can allocate and initialize an array of device information structures in its *pfxinit()* entry point.

Your driver constructs a new *dev\_t* value, specifying its major number and the selected minor number. The *makedevice()* function is used for this (see the *makedevice(D3)* reference page, which has some sample code for use in a clone open). The new *dev\_t* value is stored into the *\*devp* argument passed to *pfxopen()*.

## Summary of Standard STREAMS Functions

The supported kernel functions for STREAMS operations are summarized for reference in Table 16-2. To declare the necessary prototypes and data types, include *sys/types.h* and *sys/stream.h*.

**Table 16-2** Kernel Entry Points

Name	Can Sleep?	Summary
adjmsg(D3)	N	Trim bytes from a message.
allocb(D3)	N	Allocate a message block.
bcanput(D3)	N	Test for flow control in a specified priority band.
bcanputnext(D3)	N	Test for flow control in a specified priority band.
bufcall(D3)	N	Call a function when a buffer becomes available.
canput(D3)	N	Test for room in a message queue.
canputnext(D3)	N	Test for room in a message queue.
copyb(D3)	N	Copy a message block.
copymsg(D3)	N	Copy a message.
datamsg(D3)	N	Test whether a message is a data message.
dupb(D3)	N	Duplicate a message block.
dupmsg(D3)	N	Duplicate a message.
enableok(D3)	N	Allow a queue to be serviced.
esballoc(D3)	N	Allocate a message block using an externally-supplied buffer.
esbcall(D3)	N	Call a function when an externally-supplied buffer can be allocated.
flushband(D3)	N	Flush messages in a specified priority band.
flushq(D3)	N	Flush messages on a queue.
freeb(D3)	N	Free a message block.
freemsg(D3)	N	Free a message.

**Table 16-2 (continued)** Kernel Entry Points

<b>Name</b>	<b>Can Sleep?</b>	<b>Summary</b>
freezestr(D3)	N	Freeze the state of a stream.
getq(D3)	N	Get the next message from a queue.
insq(D3)	N	Insert a message into a queue.
linkb(D3)	N	Concatenate two message blocks.
msgdsiz(D3)	N	Return number of bytes of data in a message.
msgpullup(D3)	N	Concatenate bytes in a message.
noenable(D3)	N	Prevent a queue from being scheduled.
OTHERQ(D3)	N	Get a pointer to queue's partner queue.
pcmsg(D3)	N	Test whether a message is a priority control message.
pullupmsg(D3)	N	Concatenate bytes in a message.
putbq(D3)	N	Place a message at the head of a queue.
putctl(D3)	N	Send a control message to a queue.
putctl1(D3)	N	Send a control message with a one-byte parameter to a queue.
putnext(D3)	N	Send a message to the next queue.
putnextctl(D3)	N	Send a control message to a queue.
putnextctl1(D3)	N	Send a control message with a one-byte parameter to a queue.
putq(D3)	N	Put a message on a queue.
qenable(D3)	N	Schedule a queue's service routine to be run.
qprocsoff(D3)	Y	Enable put and service routines.
qprocson(D3)	Y	Disable put and service routines.
qreply(D3)	N	Send a message in the opposite direction in a stream.
qsize(D3)	N	Find the number of messages on a queue.
RD(D3)	N	Get a pointer to the read queue.

<b>Name</b>	<b>Can Sleep?</b>	<b>Summary</b>
rmvb(D3)	N	Remove a message block from a message.
rmvq(D3)	N	Remove a message from a queue.
SAMESTR(D3)	N	Test if next queue is of the same type.
strqget(D3)	N	Get information about a queue or band of the queue.
strqset(D3)	N	Change information about a queue or band of the queue.
unbufcall(D3)	N	Cancel a pending bufcall request.
unfreezestr(D3)	N	Unfreeze the state of a stream.
unlinkb(D3)	N	Remove a message block from the head of a message.
WR(D3)	N	Get a pointer to the write queue.

## STREAMS Modules for X Input Devices

The Silicon Graphics, Inc. implementation of the X display manager, *Xsgi*, is a customized version of the MIT X11 Sample Server. Besides other enhancements such as integration with Silicon Graphics proprietary graphics subsystems, *Xsgi* implements a generalized input subsystem so that unusual input devices can easily be integrated into the X window system. The input system is based on STREAMS modules.

### The X Input Subsystem

While X mandates that every X server support a keyboard and mouse, there is no standard system interface for accessing such devices on UNIX systems. This means each vendor has its own input subsystem for its X server. SGI's input subsystem not only meets the basic requirement to support a keyboard and mouse but also has the following features:

- A shared memory input queue is supported for high performance
- A wide variety of input devices is supported, including 3D devices such as the Spaceball

- Input devices are supported abstractly; knowledge of specific input devices is isolated to modular kernel-level device drivers
- Hardware cursor tracking is supported in the kernel

These features provide a more functional, responsive input subsystem than that available in the MIT Sample Server.

The programming interface to the input subsystem from the X client API is covered in the *X11 Input Extension Library Specification*, an online book that is distributed with the IRIX Developer's Option.

**Note:** Numerous code examples demonstrating the X input system are available in the X developer component (`x_dev` component) of the IRIX Developer Option. Source for STREAMS modules to integrate a Spaceball, a dial-and-button box, and other devices can be found in subdirectories of `/usr/share/src/X`.

### Shared Memory Input Queue

A shared memory input queue (called a *shmiq* in Silicon Graphics code comments, and pronounced "shmick") is a fast way of receiving input device events by eliminating the filesystem overhead to receive data from input devices. Instead of the X server reading the input devices through file descriptors, a kernel-level driver deposits input events directly into a region of the X server's address space, organized as a ring buffer.

The IRIX *shmiq* device driver is implemented as a STREAMS multiplexor. This allows an arbitrary number of input sources (in the form of STREAMS modules) to be linked to it so all input sources are funneled through the *shmiq*.

In addition to processing input events from input device modules, the *shmiq* driver also processes events from the graphics subsystem, and updates the screen cursor position. This allows smooth cursor movement since cursor positioning is done in kernel code, without *Xsgi* involvement.

## IDEV Interface

X input devices are integrated into the shmiq driver by implementing STREAMS modules that translate raw device input into abstract events which are sent to the shmiq driver (and on to the server). For example, an input device that connects to a serial port can be integrated in the form of a STREAMS module that is pushed onto the stream from that serial device, and translates incoming bytes into event messages.

The shmiq driver expects messages from all input devices to be in the form of IDEV events, as documented in the */usr/include/sys/idev.h* header file; hence this is called the IDEV interface. IDEV device events appear as valuator, button, and pointer state changes.

The IDEV interface defines two-way communications between the input device and *Xsgi*. Besides the uniform set of IDEV input events, the interface defines a standard set of abstract commands that *Xsgi* can send down (using IOCTL messages) to initialize and control input devices. This allows the server to see input devices as abstract input sources and does not require special server code to be written every time a new input device is supported. Instead, device specific knowledge of each device is encapsulated in an IDEV-based STREAMS module linked into the kernel.

## Input Device Naming

*Xsgi* recognizes as input devices, any device special files in the */dev/input* directory. At a minimum this includes */dev/input/keyboard* and */dev/input/mouse*. Other input devices that are to be integrated into the IDEV interface must also appear in */dev/input*.

Typically an X input device is defined as a link from */dev/input* to some other device special file, for example a serial port in the */dev/tty\** series. The filename in */dev/input* determines the name of the STREAMS module that is used to interface that device to the IDEV input system. For example, if the file is */dev/input/calcomp*, the *calcomp* STREAMS module is loaded and pushed onto the stream from the device.

When a single STREAMS module is used to support two or more devices, you can use a hyphen-digit suffix on the filename. For example, the *calcomp* STREAMS module would be used for both */dev/input/calcomp-1* and */dev/input/calcomp-2*.

When a device is initialized (as described in the next section), the STREAMS module is asked to return the X name of the input device. This name can be the same as the name of the device and the module, or it can be different. Typically the device and module names will reflect the hardware type (for example *calcomp*), while the X name reflects the kind of device (for example *tablet*).

## Opening Input Devices

An input device is opened at one of two times: when the X server starts up, and when an X client requests an open.

### Starting Up the Server

When *Xsgi* starts up, it opens each device name in */dev/input* and for each one it:

- Loads a STREAMS module that has the same name as the name of the device special file, and pushes it onto the stream from the device, below the *shmiq* multiplexor.

The STREAMS module may be loadable, and most IDEV modules are loadable.

- Looks for a file in */usr/lib/X11/input/config* having the same name as the module. The device controls in that file are sent down the stream as IOCTL messages.

The format of device controls is discussed under “Device Controls” on page 518.

- Asks the device to describe itself. This is done by sending down an IOCTL message of the type *IDEVDESC*. The module must return the IOCTL message with descriptive data.

The IDEV IOCTL structures are declared in */usr/include/sys/idev.h*. A key element of the device description is the X name of the input device.

- Looks for a file in */usr/lib/X11/input/config* having the X name of the device as returned in the device description. The X init controls in this file are processed by the X server.

The format of X init controls is discussed under “Device Controls” on page 518.

- Unless *autostart* was specified for this device, the device is closed.

### Opening from a Client

An X application can use the **XListInputDevices()** function to get a list of available input devices. Then it can call **XOpenDevice()** to open a selected device, so that input events from that device will be processed by the X server (see the **XListInputDevices(3X)** and **XOpenDevice(3X)** reference pages).

When **XOpenDevice()** is called for an input device that is not already open, it repeats the process done at startup time:

- Loads the STREAMS module and pushes it on the device stream, feeding the **shmiq** multiplexor.
- Sends device controls from a file in */usr/lib/X11/input/config* having the same name as the module.
- Asks the device (module) to describe itself, including the X name of the device.
- Processes X init controls from a file in */usr/lib/X11/input/config* having the X name of the device.

### Device Controls

Device controls are string values that are passed via an **IOCTL** message to the STREAMS module for an input device at the time the device is opened. You can use device controls as a way of configuring the device module at runtime. Device controls are interpreted only by the module.

X init controls have the same syntax as device controls, but are processed by the X server after the device has been initialized.

### Where Controls Are Stored

You can issue X server device controls on the fly by calling **XSGIDeviceControl** from within a program, or by storing them in configuration files in the */usr/lib/X11/input/config* directory. Specific documentation on controls can be found in */usr/lib/X11/input/config/README*.

There are (potentially) two configuration files per device. As noted under “Opening Input Devices” on page 517, the X server looks for device controls in a file with the same name as the STREAMS module that implements the device. After the module returns the X name of the device, the X server looks for X init controls in a file with the X name of the device.

Some devices use the same name for the STREAMS module and for the X device (*tablet*, *mouse*), but some use different names for the two. For example, the STREAMS module for the Spaceball device is *sball*, while the X name is *spaceball*.

The X server intercepts about a dozen **x\_init** controls. For a list of the **x\_init** controls and some of the more common **device\_init** controls, see the file

### Control Syntax

When the X server opens a file to look for device controls, it searches the file for a single set of controls with the following format:

```
device_init {
    name    "value"
    ...
}
```

Each *name* may have at most 15 characters. Each *value* may have at most 23 characters. Each pair of name and value are put in an IOCTL message of *idevOtherControl* type and sent down to the device module for interpretation.

When the X server opens a file to look for X init controls, it searches the file for a single set of controls with the following format:

```
x_init {
    name    "value"
    ...
}
```

The syntax is the same, except for the use of **x\_init** instead of **device\_init**.

The specific *name* and *value* strings that the X server supports are documented in the file */usr/lib/X11/input/config/README*. Any *name* strings that are not recognized by the X server are sent down to the device module, just as if they were device controls.



---

## Silicon Graphics Driver/Kernel API

This appendix summarizes the Silicon Graphics Driver/Kernel Authorized Programming Interface in tabular form. The data structures, entry points, and kernel functions are listed alphabetically with cross-references to the pages where they are discussed. The tables also show which functions and structures are compatible with SVR4 and which are unique to IRIX.

The tables in this appendix are based on the reference pages in volume D. The reference pages in volume D constitute the formal, engineering definition of the Driver/Kernel API. When discussion in this book disagrees with the contents of a reference page, the reference page takes precedence (however, any such disagreement should be reported by email to techpubs@sgi.com).

- “Driver Exported Names” on page 522 tabulates the names of data and functions that a driver must export.
- “Kernel Data Structures and Declarations” on page 523 tabulates the objects used in the interface.
- “Kernel Functions” on page 525 tabulates the IRIX kernel services used by drivers.

Each table in this appendix has a column headed “Versions.” The codes in this column have the following meanings:

SV	Syntactically and semantically portable from SVR4 UNIX, as documented in the <i>UNIX SVR4.2 Device Driver Reference</i> .
SV*	Syntactically portable from UNIX SVR4, but semantics may differ. Read the discussion and reference page carefully when porting.
5.3	Portable from IRIX version 5.3.
5.3*	Portable from IRIX 5.3, but interface has changed in some detail or new ability has been added.
6.2	Introduced in IRIX version 6.2.

## Driver Exported Names

The kernel loader, *lboot*, recognizes certain exported names of static data and functional entry points. These exported names are summarized in Table A-1.

**Table A-1** Driver Exported Names

Name	Summary	Discussed	Versions
attach	Notify driver of device attachment.	page 386	6.3
close(D2)	Notify driver of final close of minor device.	page 153	SV, 5.3
detach	Notify driver of removed device.	page 401	6.3
devflag(D1)	Show driver attributes to <i>lboot</i> .	page 145	SV*, 5.3*
edtfinit(D2)	Initialize driver from VECTOR information.	page 148	5.3
halt(D2)	Notify driver of system shutdown.	page 171	SV, 5.3
info(D1)	Show driver entries to STREAMS interface.	page 500	SV, 5.3
init(D2)	Initialize driver early in system startup.	page 148	SV*, 5.3
intr(D2)	Notify driver of device interrupt.	page 167	SV, 5.3
ioctl(D2)	Call driver to implement ioctl() call.	page 154	SV*, 5.3
map(D2)	Call driver to implement IRIX mmap().	page 163	5.3
mmap(D2)	Call driver to implement mmap().	page 165	SV*, 5.3
open(D2)	Call driver to open a device.	page 150	SV, 5.3
print(D2)	Call block driver to display filesystem error.	page 172	SV, 5.3
put(D2)	Call STREAMS driver to receive message.	page 502	SV, 5.3
read(D2)	Call character driver to read data.	page 155	SV, 5.3
size(D2)	Call block driver to get device capacity.	page 172	SV, 5.3
srv(D2)	Call driver to service queued messages.	page 503	SV, 5.3
start(D2)	Initialize driver late in system startup.	page 149	SV, 5.3
strategy(D2)	Call block driver to read or write data.	page 157	SV*, 5.3
unload(D2)	Call loadable driver prior to unloading it.	page 170	5.3

**Table A-1 (continued)** Driver Exported Names

Name	Summary	Discussed	Versions
unmap(D2)	Call driver to notify it of unmap() call.	page 166	5.3
write(D2)	Call character driver to write data.	page 155	SV, 5.3

The following reference pages have overview information on exported names: intro(D1), intro(D2), and prefix(D1).

**Note:** The following SVR4 exported names are not used in IRIX drivers: `chpoll`, `_load`, and `_unload`. The latter is replaced by `pfxload()` without the leading underscore.

## Kernel Data Structures and Declarations

The driver/kernel interface is based on shared use of certain data types and defined constant values. For general information on these interface objects, see the intro(D4) and intro(D5) reference pages.

The interface objects used by device drivers are summarized in Table A-2. .

**Table A-2** Device Driver Interface Objects

Name	Summary	Discussed	Versions
alenlist (alenlist(d4x))	List of addresses and lengths of memory segments.	page 400	6.3
buf(D4)	Block read/write request structure.	page 185	SV*, 5.3*
eisa_dma_cb(D4)	DMA command block for EISA slave DMA.		5.3
eisa_dma_buf(D4)	DMA command buffer for EISA slave DMA.		5.3
errnos(D5)	Error numbers valid for driver use.		SV*, 5.3
iovec(D4)	Describes an I/O buffer segment to the read or write entry points.	page 184	SV, 5.3
signals(D5)	Lists signal numbers valid for driver use.		SV*, 5.3
uio(D4)	Describes an I/O request to the read or write entry points.	page 184	SV*, 5.3

**Note:** The following data structures used in SVR4 drivers are not used in IRIX: *dma\_buf* and *dma\_cb*. The *eisa\_dma\_buf* and *eisa\_dma\_cb* structures are similar but are used only in EISA drivers.

The interface objects used by STREAMS drivers are summarized in Table A-3

**Table A-3** STREAMS Driver Interface Objects

Name	Summary	Discussed	Versions
copyreq(D4)	Copy request structure.		SV, 5.3
copyresp(D4)	Copy response structure.		SV, 5.3
datab(D4)	Message data block.		SV, 5.3
free_rtn(D4)	Describes a message-free routine.		SV, 5.3
iocblk(D4)	Describes ioctl() data or response.		SV, 5.3
linkblk(D4)	Describes multiplexed link.		SV, 5.3
module_info(D4)	Describes module attributes.		SV, 5.3
msgb(D4)	Describes all or part of a message.		SV, 5.3
qinit(D4)	Points to handlers and parameters for a queue.		SV, 5.3
queue(D4)	Describes a queue of messages.		SV, 5.3
streamtab(D4)	Points to the queues handled by a driver.		SV, 5.3
stroptions(D4)	Lists stream-head options.		SV, 5.3

## Kernel Functions

The IRIX kernel makes available the functions summarized in Table A-4. For PCI drivers, see also the functions listed in Table 15-7 on page 403

**Table A-4** Kernel Functions

Name	Summary	Discussed	Versions
adjmsg(D3)	Trim bytes from a message.		SV, 5.3
alloca(D3)	Allocate a message block.		SV, 5.3
ASSERT(D3)	Debugging macro designed for use in the kernel (compare to assert(3X)).	page 254	5.3
badaddr(D3)	Test physical address for input.	page 198	5.3
badaddr_val(D3)	Test physical address for input and return the input value received.	page 198	6.2
bcanput(D3)	Test for flow control in a specified priority band.		SV, 5.3
bcanputnext(D3)	Test for flow control in a specified priority band.		SV, 5.3
bcmp(D3)	Compare data between kernel locations.	page 196	SV, 5.3
bcopy(D3)	Copy data between locations in the kernel.	page 196	SV, 5.3
biodone(D3)	Mark a <i>buf_t</i> as complete and wake any process waiting for it.	page 219	SV, 5.3
bioerror(D3)	Manipulate error fields within a <i>buf_t</i> .	page 219	SV, 5.3
biowait(D3)	Suspend process pending completion of block I/O.	page 219	SV, 5.3
bp_mapin(D3)	Map buffer pages into kernel virtual address space.	page 202	SV, 5.3
bp_mapout(D3)	Release mapping of buffer pages.	page 202	SV, 5.3
bptophys(D3)	Get physical address of buffer data.	page 201	5.3
brelease(D3)	Return a buffer to the system's free list.	page 192	SV, 5.3
btod(D3)	Return number of 512-byte "sectors" in a byte count (round up).	page 200	5.3

<b>Table A-4 (continued)</b>		Kernel Functions	
<b>Name</b>	<b>Summary</b>	<b>Discussed</b>	<b>Versions</b>
btop(D3)	Return number of I/O pages in a byte count (truncate).	page 200	SV, 5.3
btopr(D3)	Return number of I/O pages in a byte count (round up).	page 200	SV, 5.3
bufcall(D3)	Call a function when a buffer becomes available.		SV, 5.3
bzero(D3)	Clear kernel memory for a specified size.	page 195	SV, 5.3
canput(D3)	Test for room in a message queue.		SV, 5.3
canputnext(D3)	Test for room in a message queue.		SV, 5.3
clrbuf(D3)	Erase the contents of a buffer described by a buf_t.	page 202	SV, 5.3
cmn_err(D3)	Display an error message or panic the system.	page 251	SV*, 5.3
copyb(D3)	Copy a message block.		SV, 5.3
copyin(D3)	Copy data from user address space.	page 195	SV, 5.3
copymsg(D3)	Copy a message.		SV, 5.3
copyout(D3)	Copy data to user address space.	page 195	SV, 5.3
cpsema(D3)	Conditionally decrement a semaphore's state.	page 224	5.3
cvsema(D3)	Conditionally increment a semaphore's state	page 224	5.3
datamsg(D3)	Test whether a message is a data message.		SV, 5.3
delay(D3)	Delay for a specified number of clock ticks.	page 216	SV, 5.3
disable_sysad_parity()	Disable memory parity checking on SysAD bus.		
dki_dcache_inval(D3)	Invalidate the data cache for a given range of virtual addresses.	page 203	5.3

**Table A-4 (continued)** Kernel Functions

<b>Name</b>	<b>Summary</b>	<b>Discussed</b>	<b>Versions</b>
dki_dcache_wb(D3)	Write back the data cache for a given range of virtual addresses.	page 203	5.3
dki_dcache_wbinval(D3)	Write back and invalidate the data cache for a given range of virtual addresses.	page 203	5.3
dma_map(D3)	Load DMA mapping registers for an imminent VME transfer.		5.3
dma_mapbp(D3)	Load DMA mapping registers for an imminent VME transfer.		5.3
dma_mapaddr(D3)	Return the “bus virtual” address for a given VME map and address.		5.3
dma_mapalloc(D3)	Allocate a VME DMA map.		5.3
dma_mapfree(D3)	Free a VME DMA map.		5.3
drv_getparm(D3)	Retrieve kernel state information.	page 205	SV*, 5.3
drv_hztousec(D3)	Convert clock ticks to microseconds	page 216	SV, 5.3
drv_priv(D3)	Test for privileged user.	page 205	SV, 5.3
drv_setparm(D3)	Set kernel state information.	page 205	SV, 5.3
drv_usectohz(D3)	Convert microseconds to clock ticks.	page 216	SV, 5.3
drv_usecwait(D3)	Busy-wait for a specified interval.	page 216	SV, 5.3
dtimeout(D3)	Schedule a function execute on a specified processor after a specified length of time.	page 216	5.3
dupb(D3)	Duplicate a message block.		SV, 5.3
dupmsg(D3)	Duplicate a message.		SV, 5.3
eisa_dma_disable(D3)	Disable recognition of hardware requests on an EISA DMA channel.		5.3
eisa_dma_enable(D3)	Enable recognition of hardware requests on a EISA DMA channel.		5.3
eisa_dma_free_buf(D3)	Free a previously allocated EISA DMA buffer descriptor.		5.3

<b>Table A-4 (continued)</b>		Kernel Functions	
<b>Name</b>	<b>Summary</b>	<b>Discussed</b>	<b>Versions</b>
eisa_dma_free_cb(D3)	Free a previously allocated EISA DMA command block.		5.3
eisa_dma_get_buf(D3)	Allocate an EISA DMA buffer descriptor.		5.3
eisa_dma_get_cb(D3)	Allocate an EISA DMA command block.		5.3
eisa_dma_prog(D3)	Program an EISA DMA operation for a subsequent software request.		5.3
eisa_dma_stop(D3)	Stop software-initiated EISA DMA operation and release channel.		5.3
eisa_dma_swstart(D3)	Initiate an EISA DMA operation via software request.		5.3
eisa_dmachan_alloc()	Allocate a DMA channel for EISA slave DMA.		5.3
eisa_ivec_alloc()	Allocate an IRQ level for EISA.		5.3
eisa_ivec_set()	Associate a handler with an EISA IRQ.		5.3
enableloek(D3)	Allow a queue to be serviced.		SV, 5.3
enable_sysad_parity()	Reenable parity checking on SysAD bus.		
esballoc(D3)	Allocate a message block using an externally-supplied buffer.		SV, 5.3
esbcall(D3)	Call a function when an externally-supplied buffer can be allocated.		SV, 5.3
etoimajor(D3)	Convert external to internal major device number.	page 182	SV, 5.3
fast_itimeout(D3)	Same as itimeout() but takes an interval in "fast ticks."	page 216	6.2
fasthzto(D3)	Returns the value of a <i>struct timeval</i> as a count of "fast ticks."	page 216	6.2
flushband(D3)	Flush messages in a specified priority band.		SV, 5.3

**Table A-4 (continued)** Kernel Functions

<b>Name</b>	<b>Summary</b>	<b>Discussed</b>	<b>Versions</b>
flushbus(D3)	Make sure contents of the write buffer are flushed to the system bus	page 203	5.3
flushq(D3)	Flush messages on a queue.		SV, 5.3
freeb(D3)	Free a message block.		SV, 5.3
freemsg(D3)	Free a message.		SV, 5.3
freerbuf(D3)	Free a <i>buf_t</i> with no buffer.	page 192	SV, 5.3
freesema(D3)	Free the resources associated with a semaphore.	page 224	5.3*
freezestr(D3)	Freeze the state of a stream.		SV, 5.3
fubyte(D3)	Load a byte from user space.	page 195	5.3
fuword(D3)	Load a word from user space.	page 195	5.3
getblk(D3)	Get a <i>buf_t</i> with no buffer.	page 192	SV, 5.3
getemajor(D3)	Get external major device number.	page 182	SV, 5.3
geteminor(D3)	Get external minor device number.	page 182	SV, 5.3
geterror(D3)	retrieve error number from a buffer header	page 219	SV, 5.3
getmajor(D3)	Get internal major device number.	page 182	SV, 5.3
getminor(D3)	Get internal minor device number.	page 182	SV, 5.3
getnextpg(D3)	Return <i>pfdat</i> structure for next page.	page 202	5.3
getq(D3)	Get the next message from a queue.		SV, 5.3
getrbuf(D3)	Allocate a <i>buf_t</i> with no buffer.	page 192	SV, 5.3
hwcpin(D3)	Copy data from device registers to kernel memory.	page 195	5.3
hwcpout(D3)	Copy data from kernel memory to device registers.	page 195	5.3
initnsema(D3)	Initialize a semaphore to a specified count.	page 224	5.3
initnsema_mutex(D3)	Initialize a semaphore to a count of 1.	page 224	5.3

<b>Table A-4 (continued)</b>		Kernel Functions	
<b>Name</b>	<b>Summary</b>	<b>Discussed</b>	<b>Versions</b>
insq(D3)	Insert a message into a queue.		SV, 5.3
ip26_enable_ucmem(D3)	Change memory mode on IP26 processor.	page 29	6.2
ip26_return_ucmem(D3)	Change memory mode on IP26 processor.	page 29	SV, 5.3
is_sysad_parity_enabled( )	Test for parity checking on SysAD bus.		5.3
itimeout(D3)	Schedule a function to be executed after a specified number of clock ticks.	page 216	SV, 5.3
itoemajor(D3)	Convert internal to external major device number.	page 182	SV, 5.3
kern_calloc(D3)	Allocate and clear space from kernel memory.	page 190	5.3
kern_free(D3)	Free kernel memory space.	page 190	5.3
kern_malloc(D3)	Allocate kernel virtual memory.	page 190	5.3
kmem_alloc(D3)	Allocate space from kernel free memory.	page 190	SV, 5.3
kmem_free(D3)	Free previously allocated kernel memory.	page 190	SV, 5.3
kmem_zalloc(D3)	Allocate and clear space from kernel free memory.	page 190	SV, 5.3
kvtophys(D3)	Get physical address of kernel data.	page 202	5.3
linkb(D3)	Concatenate two message blocks.		SV*, 5.3*
LOCK(D3)	Acquire a basic lock, waiting if necessary.	page 208	SV*, 5.3*
LOCK_ALLOC(D3)	Allocate and initialize a basic lock.	page 208	SV*, 5.3*
LOCK_DEALLOC(D3)	Deallocate an instance of a basic lock.	page 208	SV*, 5.3*
LOCK_INIT(D3)	Initialize a basic lock that was allocated statically, or reinitialize an allocated lock.	page 208	6.2
LOCK_DESTROY(D3)	Uninitialize a basic lock that was allocated statically.	page 208	6.2

**Table A-4 (continued)** Kernel Functions

<b>Name</b>	<b>Summary</b>	<b>Discussed</b>	<b>Versions</b>
makedevice(D3)	Make device number from major and minor numbers.	page 182	SV, 5.3
max(D3)	Return the larger of two integers.		SV, 5.3
min(D3)	Return the lesser of two integers.		SV, 5.3
msgdsz(D3)	Return number of bytes of data in a message.		SV, 5.3
msgpullup(D3)	Concatenate bytes in a message.		SV, 5.3
MUTEX_ALLOC(D3)	Allocate and initialize a mutex lock.	page 210	6.2
MUTEX_DEALLOC(D3)	Deinitialize and free a dynamically allocated mutex lock.	page 210	6.2
MUTEX_DESTROY(D3)	Deinitialize a mutex lock.	page 210	6.2
MUTEX_INIT(D3)	Initialize an existing mutex lock.	page 210	6.2
MUTEX_ISLOCKED(D3)	Test if a mutex lock is owned.	page 210	6.2
MUTEX_LOCK(D3)	Claim a mutex lock.	page 210	6.2
MUTEX_MINE(D3)	Test if a mutex lock is owned by this process.	page 210	6.2
MUTEX_TRYLOCK(D3)	Conditionally claim a mutex lock.	page 210	6.2
MUTEX_UNLOCK(D3)	Release a mutex lock.	page 210	6.2
MUTEX_WAITQ(D3)	Get the number of processes blocked by mutex lock.	page 210	6.2
ngeteblk(D3)	Allocate a <i>buf_t</i> and a buffer of specified size.	page 192	SV, 5.3
noenable(D3)	Prevent a queue from being scheduled.		SV, 5.3
OTHERQ(D3)	Get a pointer to queue's partner queue.		SV, 5.3
pcmsg(D3)	Test whether a message is a priority control message.		SV, 5.3
phalloc(D3)	Allocate and initialize a pollhead structure.	page 192	SV, 5.3

<b>Table A-4 (continued)</b>		Kernel Functions	
<b>Name</b>	<b>Summary</b>	<b>Discussed</b>	<b>Versions</b>
phfree(D3)	Free a pollhead structure.	page 192	SV, 5.3
physiock(D3)	Validate and issue a raw I/O request	page 219	SV, 5.3
pio_andb_rmw(D3)	VME Byte read-and-write.		5.3
pio_andh_rmw(D3)	VME 16-bit read-and-write.		5.3
pio_andw_rmw(D3)	VME 32-bit read-and-write.		5.3
pio_badaddr(D3)	Check for VME bus error when reading an address.		5.3
pio_badaddr_val(D3)	Check for VME bus error when reading an address and return the value read.		5.3
pio_bcopyin(D3)	Copy data from a VME bus address to kernel's virtual space.		5.3
pio_bcopyout(D3)	Copy data from kernel's virtual space to a VME bus address.		5.3
pio_mapaddr(D3)	Convert a VME bus address to a virtual address.		5.3
pio_mapalloc(D3)	Allocate a VME PIO map.		5.3
pio_mapfree(D3)	Free a VME PIO map.		5.3
pio_orb_rmw(D3)	VME Byte read-or-write.		5.3
pio_orh_rmw(D3)	VME 16-bit read-or-write.		5.3
pio_orw_rmw(D3)	VME 32-bit read-or-write.		5.3
pio_wbadaddr(D3)	Check for VME bus error when writing to an address.		5.3
pio_wbadaddr_val(D3)	Check for VME bus error when writing a specified value to an address.		5.3
pollwakeup(D3)	Inform polling processes that an event has occurred.	page 159	SV, 5.3
pptophys(D3)	Convert page pointer to physical address.	page 202	SV, 5.3

**Table A-4 (continued)** Kernel Functions

<b>Name</b>	<b>Summary</b>	<b>Discussed</b>	<b>Versions</b>
proc_ref(D3)	Obtain a reference to a process for signaling.	page 205	SV, 5.3
proc_signal(D3)	Send a signal to a process.	page 205	SV, 5.3
proc_unref(D3)	Release a reference to a process.	page 205	SV, 5.3
psema(D3)	Perform a “P” or wait semaphore operation.	page 224	SV, 5.3
ptob(D3)	Convert size in pages to size in bytes.	page 200	SV, 5.3
pullupmsg(D3)	Concatenate bytes in a message.		SV, 5.3
putbq(D3)	Place a message at the head of a queue.		SV, 5.3
putctl(D3)	Send a control message to a queue.		SV, 5.3
putctl1(D3)	Send a control message with a one-byte parameter to a queue.		SV, 5.3
putnext(D3)	Send a message to the next queue.		SV, 5.3
putnextctl(D3)	Send a control message to a queue.		SV, 5.3
putnextctl1(D3)	Send a control message with a one-byte parameter to a queue.		SV, 5.3
putq(D3)	Put a message on a queue.		SV, 5.3
qenable(D3)	Schedule a queue’s service routine to be run.		SV, 5.3
qprocsoff(D3)	Enable put and service routines.		SV, 5.3
qprocson(D3)	Disable put and service routines		SV, 5.3
qreply(D3)	Send a message in the opposite direction in a stream.		SV, 5.3
qsize(D3)	Find the number of messages on a queue.		SV, 5.3
RD(D3)	Get a pointer to the read queue.		SV, 5.3
rmalloc(D3)	Allocate space from a private space management map.	page 193	SV, 5.3
rmallocmap(D3)	Allocate and initialize a private space management map.	page 193	SV, 5.3

<b>Table A-4 (continued)</b>		Kernel Functions	
<b>Name</b>	<b>Summary</b>	<b>Discussed</b>	<b>Versions</b>
rmalloc_wait(D3)	Allocate resources from a space management map.	page 193	SV, 5.3
rmfree(D3)	Release resources into a space management map.	page 193	SV, 5.3
rmfreemap(D3)	Free a private space management map.	page 193	SV, 5.3
rmvb(D3)	Remove a message block from a message.		SV, 5.3
rmvq(D3)	Remove a message from a queue.		SV, 5.3
RW_ALLOC(D3)	Allocate and initialize a reader/writer lock.	page 213	SV*, 5.3*
RW_DEALLOC(D3)	Deallocate a reader/writer lock.	page 213	SV*, 5.3*
RW_DESTROY(D3)	Deinitialize an existing reader/writer lock.	page 213	6.2
RW_INIT(D3)	Initialize an existing reader/writer lock.	page 213	6.2
RW_RDLOCK(D3)	Acquire a reader/writer lock as reader, waiting if necessary.	page 213	SV*, 5.3*
RW_TRYRDLOCK(D3)	Try to acquire a reader/writer lock as reader, returning a code if it is not free.	page 213	SV*, 5.3*
RW_TRYWRLOCK(D3)	Try to acquire a reader/writer lock as writer, returning a code if it is not free.	page 213	SV*, 5.3*
RW_UNLOCK(D3)	Release a reader/writer lock as reader or writer.	page 213	SV*, 5.3*
RW_WRLOCK(D3)	Acquire a reader/writer lock as writer, waiting if necessary.	page 213	SV*, 5.3*
SAMESTR(D3)	Test if next queue is of the same type.		SV, 5.3
scsi_abort()	Transmits a SCSI ABORT command.	page 305	5.3*
scsi_alloc(D3)	Open a connection between a driver and a target device.	page 305	5.3*
scsi_command(D3)	Transmit a SCSI command on the bus and return results.	page 305	5.3*
scsi_free(D3)	Release connection to target device.	page 305	5.3*

**Table A-4 (continued)** Kernel Functions

<b>Name</b>	<b>Summary</b>	<b>Discussed</b>	<b>Versions</b>
scsi_info(D3)	Issue the SCSI Inquiry command and return the results.	page 305	5.3*
scsi_reset()	Resets the SCSI adapter or bus.	page 305	5.3*
setgiovector()	Register a GIO interrupt handler.		5.3
setgioconfig()	Prepare a GIO slot for use.		5.3
sgset(D3)	Assign physical addresses to a vector of software scatter-gather registers.	page 202	5.3
sleep(D3)	Suspend process execution pending occurrence of an event.	page 221	SV, 5.3
SLEEP_ALLOC(D3)	Allocate and initialize a sleep lock.	page 212	SV*, 5.3*
SLEEP_DEALLOC(D3)	Deinitialize and deallocate a dynamically allocated sleep lock.	page 212	SV*, 5.3*
SLEEP_DESTROY	Deinitialize a sleep lock.	page 212	6.2
SLEEP_INIT(D3)	Initialize an existing sleep lock.	page 212	6.2
SLEEP_LOCK(D3)	Acquire a sleep lock, waiting if necessary until the lock is free.	page 212	SV*, 5.3*
SLEEP_LOCKAVAIL(D3)	Query whether a sleep lock is available.	page 212	SV*, 5.3*
SLEEP_LOCK_SIG(D3)	Acquire a sleep lock, waiting if necessary until the lock is free or a signal is received.	page 212	SV*, 5.3*
SLEEP_TRYLOCK(D3)	Try to acquire a sleep lock, returning a code if it is not free.	page 212	SV*, 5.3*
SLEEP_UNLOCK(D3)	Release a sleep lock.	page 212	SV*, 5.3*
splbase(D3)	Block no interrupts.	page 215	SV, 5.3
spltimeout(D3)	Block only timeout interrupts.	page 215	SV, 5.3
spldisk(D3)	Block disk interrupts.	page 215	SV, 5.3
splstr(D3)	Block STREAMS interrupts.	page 215	SV, 5.3
spltty(D3)	Block disk, VME, serial interrupts.	page 215	SV, 5.3

<b>Table A-4 (continued)</b>		Kernel Functions	
<b>Name</b>	<b>Summary</b>	<b>Discussed</b>	<b>Versions</b>
splhi(D3)	Block all I/O interrupts.	page 215	SV, 5.3
spl0(D3)	Same as <b>splbase()</b> .	page 215	SV, 5.3
splx(D3)	Restore previous interrupt level.	page 215	SV, 5.3
strcat(D3)	Append one string to another.		SV, 5.3
strcpy(D3)	Copy a string.		SV, 5.3
streams_interrupt(D3)	Synchronize interrupt-level function with STREAMS mechanism.		5.3
STREAMS_TIMEOUT(D3)	Synchronize timeout with STREAMS mechanism.		5.3
strlen(D3)	Return length of a string.		SV, 5.3
strlog(D3)	Submit messages to the log driver.		SV, 5.3
strncmp(D3)	Compare two strings for a specified length.		SV, 5.3
strncpy(D3)	Copy a string for a specified length.		SV, 5.3
strqget(D3)	Get information about a queue or band of the queue.		SV, 5.3
strqset(D3)	Change information about a queue or band of the queue.		SV, 5.3
subyte(D3)	Store a byte to user space.	page 195	5.3
suword(D3)	Store a word to user space.	page 195	5.3
SV_ALLOC(D3)	Allocate and initialize a synchronization variable.	page 222	SV*, 5.3*
SV_BROADCAST(D3)	Wake all processes sleeping on a synchronization variable.	page 222	SV*, 5.3*
SV_DEALLOC(D3)	Deinitialize and deallocate a synchronization variable.	page 222	SV*, 5.3*
SV_DESTROY	Deinitialize a synchronization variable.	page 222	6.2
SV_INIT	Initialize an existing synchronization variable.	page 222	6.2

---

**Table A-4 (continued)** Kernel Functions

<b>Name</b>	<b>Summary</b>	<b>Discussed</b>	<b>Versions</b>
SV_SIGNAL(D3)	Wake one process sleeping on a synchronization variable.	page 222	SV*, 5.3*
SV_WAIT(D3)	Sleep until a synchronization variable is signalled.	page 222	SV*, 5.3*
SV_WAIT_SIG(D3)	Sleep until a synchronization variable is signalled or a signal is received.	page 222	SV*, 5.3*
timeout(D3)	Schedule a function to be executed after a specified number of clock ticks.	page 216	SV, 5.3
TRYLOCK(D3)	Try to acquire a basic lock, returning a code if the lock is not currently free.	page 208	SV*, 5.3*
uiomove(D3)	Copy data using <i>uio_t</i> .	page 197	SV, 5.3
uiophysio(D3)	Validate a raw I/O request and pass to a strategy function.	page 219	5.3
unbufcall(D3)	Cancel a pending bufcall request.		SV, 5.3
undma(D3)	Unlock physical memory in user space.	page 219	5.3
unfreezestr(D3)	Unfreeze the state of a stream.		SV, 5.3
unlinkb(D3)	Remove a message block from the head of a message.		SV, 5.3
UNLOCK(D3)	Release a basic lock.	page 208	SV*, 5.3*
untimeout(D3)	Cancel a previous itimeout or fast_itimeout request.	page 216	SV*, 5.3*
ureadc(D3)	Copy a character to space described by <i>uio_t</i> .	page 197	SV, 5.3
userdma(D3)	Lock physical memory in user space. small number of	page 219	5.3
userabi()	Get data sizes for the ABI of the user process (32- or 64-bit).	page 173	6.2
uwritec(D3)	Return a character from space described by <i>uio_t</i> .	page 197	SV, 5.3

**Table A-4 (continued)**      Kernel Functions

Name	Summary	Discussed	Versions
v_getaddr(D3)	Get the user virtual address associated with a <i>vhandl_t</i> .	page 199	5.3
v_gethandle(D3)	Get a unique identifier associated with a <i>vhandl_t</i> .	page 199	5.3
v_getlen(D3)	Get the length of user address space associated with a <i>vhandl_t</i> .	page 199	5.3
v_mapphys(D3)	Map kernel address space into user address space.	page 199	5.3
valusema(D3)	Return the value associated with a semaphore.	page 224	5.3
vme_adapter(D3)	Determine VME adapter that corresponds to a given memory address.		5.3
vme_ivec_alloc(D3)	Allocate a VME bus interrupt vector.		5.3
vme_ivec_free(D3)	Free a VME bus interrupt vector.		5.3
vme_ivec_set(D3)	Register a VME bus interrupt vector.		5.3
vsema(D3)	Perform a "V" or signal semaphore operation.	page 224	5.3
wakeup(D3)	Waken a process waiting for an event.	page 221	SV, 5.3
wbadaddr(D3)	Test physical address for output.	page 198	SV, 5.3
wbadaddr_val(D3)	Test physical address for output of specific value.	page 198	SV, 5.3
WR(D3)	Get a pointer to the write queue.		SV, 5.3

The following SVR4 kernel functions are not implemented in IRIX: *bioreset*, *dma\_disable*, *dma\_enable*, *dma\_free\_buf*, *dma\_free\_cb*, *dma\_get\_best\_mode*, *dma\_get\_buf*, *dma\_get\_cb*, *dma\_pageio*, *dma\_prog*, *dma\_swstart*, *dma\_swsetup*, *drv\_gethardware*, *hat\_getkpfnum*, *hat\_getppfnum*, *inb*, *inl*, *inw*, *kvtoppid*, *mod\_drvattach*, *mod\_drvdetach*, *outb*, *outl*, *outw*, *physmap*, *physmap\_free*, *phystoppid*, *psignal*, *rdma\_filter*, *repinsb*, *repinsd*, *repinsw*, *repoutsb*, *repoutsd*, *repoutsw*, *rminit*, *rmsetwant*, *SLEEP\_LOCKOWNED*, *strncat*, *vtop*.

---

## New and Updated Reference Pages

This appendix contains the text of new and updated reference pages that have been created since IRIX 6.3 initially shipped.

- “Address/Length List Reference Pages” on page 539 displays two pages about alenlist concepts and functions.
- “PCI Infrastructure Reference Pages” on page 547 displays the pages about the PCI bus support functions.

### Address/Length List Reference Pages

- Example B-1 displays the text of an overview of address/length lists.
- Example B-2 displays the reference page listing operations on alenlists.

#### Example B-1 alenlist(d4x)

**NAME**

alenlist - overview of Address/Length Lists

**DESCRIPTION**

An Address/Length List or alenlist is an abstract object containing a sequential list of address/length pairs.

All addresses in the list belong to the same address space, for example kernel virtual address space, physical memory address space, PCI DMA space, VME DMA space, and so on. All lengths are byte counts. Each address/length pair added to a list describes one contiguous segment in the relevant address space. Together, all the pairs in the list describe a region of memory, typically an I/O buffer, that is logically, but not necessarily physically, contiguous.

When a driver receives a request for I/O, the request may be specified in one of several forms, for example a kernel virtual address, a user virtual address, or a vector of addresses from a device register array. The driver can convert any of these forms into a single alenlist for easy management.

Other kernel layers (notably, the PCI infrastructure, see *pciio*(D3)) accept alenlists. Translation from one address space to another, a constant challenge for device drivers, can be performed on an entire alenlist in a single operation, avoiding repeated bouncing up and down through layers of software for each segment.

An alenlist is created by the driver when needed. It expands automatically as address/length pairs are added to the list. Its memory can easily be released for recycling, or a list can be cleared and re-used. There are operations to append a new pair to a list, and to load a list based on a user or kernel virtual buffer address and size.

Address/length pairs can be read out from a list in sequence, and in length units different from the units that were initially loaded (for example, load the list based on a user-space buffer, then read out pairs in 512-byte segments).

A cursor marks a position within an alenlist, and is used to scan through the pairs in the list. A cursor can point to the beginning of its list, any number of bytes into the list, or at the end of the list. Every alenlist contains an implicit cursor for simple management, but there are operations to create, destroy, and initialize additional cursors, so multiple positions can be maintained at once.

#### Usage

In order to use any of the Address/Length types or interfaces, a driver should include these header files in this order:

```
#include <sys/types.h>
#include <sys/alenlist.h>
#include <sys/ddi.h>
```

#### Types

The following abstract types are used and returned by alenlist functions.

*alenaddr\_t* An abstract address in some address space. This is the type of any address added to a list. It is guaranteed to be big enough to hold an address in any supported address space (currently 64 bits).

*size\_t* The standard byte-count type. Every address/length pair in a list comprises one *alenaddr\_t* and one *size\_t*.

*alenlist\_t* A handle to an alenlist.

*alenlist\_cursor\_t*

A handle to a cursor, an object that designates a position within an alenlist. An *alenlist\_cursor\_t* value of NULL always designates the implicit cursor of a list. Non-null values designate cursor objects allocated separate from their lists.

**Operations**

The functions that operate on alenlists are detailed in *alenlist\_ops(D3X)*. Operations to allocate and destroy alenlists are as follows:

- o *alenlist\_create()* creates a new, empty list.
- o *alenlist\_destroy()* recycles the storage used by an alenlist.
- o *alenlist\_cursor\_create()* creates a new cursor object.
- o *alenlist\_cursor\_destroy()* recycles the storage used by a cursor.

Operations to clear and load pairs into an alenlist are as follows:

- o *alenlist\_clear()* empties a list.
- o *alenlist\_append()* appends one address/length pair to a list.
- o *kvaddr\_to\_alenlist()* converts a kernel virtual address and length into physical memory address/length pairs and appends the pairs to a list.
- o *uvaddr\_to\_alenlist()* converts a user process virtual address into physical memory address/length pairs and appends the pairs to a list.
- o *buf\_to\_alenlist()* examines a *buf* structure (see *buf(D4)*) and appends physical memory address/length pairs to describe the buffer to a list.

In addition, *pciio\_dma(D3)* documents *pciio\_dmamap\_list()*, a function that translates an alenlist of physical memory addresses into an alenlist of PCI bus addresses for DMA use.

Operations to read out address/length pairs from an alenlist are as follows:

- o *alenlist\_get()* returns an address/length pair from a list, based on

the implicit cursor or on a specified cursor, and optionally advances the cursor.

- o `alenlist_cursor_init()` initializes a specific cursor or an implicit cursor to a specified offset.
- o `alenlist_cursor_offset()` returns the current offset of a cursor.

**Flags**

The following flag values are used with `alenlist` functions:

`AL_NOSLEEP` Indicates that the caller does not wish to be put to sleep while waiting for resources such as memory. If an operation would entail sleeping, the function returns failure.

`AL_NOCOMPACT` Indicates that adjacent address/length pairs that happen to be contiguous in their address space should not be compacted and treated as a single unit (which is the default behavior). Rather, they should be treated as if there is a discontinuity. This flag is not needed by most drivers.

`AL_LEAVE_CURSOR` Indicates that the operation should not affect the position of the internal cursor.

**Error Codes**

The following error codes are used by all `alenlist` functions:

`ALENLIST_SUCCESS`  
Indicates a successful operation.

`ALENLIST_FAILURE`  
Indicates a failed operation.

**SEE ALSO**

`alenlist_ops(D3X)`, `pciio_dma(D3)`, IRIX Device Driver Programmer's Guide

**Example B-2** `alenlist_ops(d3x)`

**NAME**

`alenlist_ops`: `alenlist_append`, `alenlist_clear`, `alenlist_create`, `alenlist_cursor_create`, `alenlist_cursor_destroy`, `alenlist_cursor_init`, `alenlist_cursor_offset`, `alenlist_destroy`, `alenlist_get`, `kvaddr_to_alenlist`, `uvaddr_to_alenlist`, `buf_to_alenlist` - operations on address/length lists

**SYNOPSIS**

```
#include <sys/types.h>
#include <sys/alenlist.h>
#include <sys/ddi.h>

alenlist_t
alenlist_create(unsigned flags);

void
alenlist_clear(alenlist_t plist);

void
alenlist_destroy(alenlist_t plist);

alenlist_cursor_t
alenlist_cursor_create(alenlist_t plist, unsigned flags);

void
alenlist_cursor_destroy(alenlist_cursor_t ocursor);

int
alenlist_get(alenlist_t plist,
            alenlist_cursor_t icursor,
            size_t maxlength,
            alenaddr_t *paddr,
            size_t *plength);

int
alenlist_cursor_init(alenlist_t plist,
                    size_t offset,
                    alenlist_cursor_t icursor);

size_t
alenlist_cursor_offset(alenlist_t plist, alenlist_cursor_t icursor);

int
alenlist_append(alenlist_t plist,
               alenaddr_t address,
               size_t length);

alenlist_t
kvaddr_to_alenlist(caddr_t kvaddr, size_t length);

alenlist_t
```

```
uvaddr_to_alenlist( alenlist_t plist,  
                   uvaddr_t uvaddr,  
                   size_t length);
```

```
alenlist_t  
buf_to_alenlist(buf_t *buf);
```

**Arguments**

- address*      An address in some address space.
- buf*          Address of a *buf* struct.
- flags*        A Boolean combination of the flags declared in *sys/alenlist.h*.
- icursor*     A handle to an existing cursor, or NULL to indicate the implicit cursor in *plist*.
- kvaddr*      A valid address in kernel virtual memory.
- length*      The length related to an address.
- maxlength*   The maximum length allowed in the returned address/length pair, or 0 to indicate no maximum applies.
- ocursor*     A handle to an existing cursor, the target of the operation.
- offset*      An initial byte offset for *icursor*, usually 0.
- paddr*        A pointer to a variable to receive the address from an address/length pair.
- plength*     A pointer to a variable to receive the length from an address/length pair.
- plist*        A handle to an existing alenlist.
- uvaddr*      A valid address in user virtual memory for the in-context process.

**DESCRIPTION**

For an overview of address/length lists (alenlists) see *alenlist(D4X)*.

**Allocation and Release**

Create an empty list using *alenlist\_create()*. The *AL\_NOSLEEP* flag

ensures that the caller will not sleep waiting for memory. If memory cannot be allocated, NULL is returned. The returned list has no pairs in it, and its implicit cursor is initialized to point to the start of the (empty) list.

The functions *kvaddr\_to\_alenlist()*, *uvaddr\_to\_alenlist()*, and *buf\_to\_alenlist()*, when the *plist* argument is NULL, allocate a new alenlist and return it. However, these functions do not honor AL\_NOSLEEP, either when creating a list or when getting memory to extend an existing list. If it is important not to sleep, preallocate the list.

Empty a list by applying *alenlist\_clear()*. The implicit cursor is reset to the start of the (empty) list.

Release a list using *alenlist\_destroy()*. This lets the system know that the specified List is no longer in use.

Create a cursor by calling *alenlist\_cursor\_create()*. Pass AL\_NOSLEEP to avoid sleeping on memory allocation. When a cursor cannot be created, NULL is returned. The new cursor is associated with *plist* and is initialized to point to the head of that list. More than one cursor can point to a given list.

Use *alenlist\_cursor\_destroy()* to tell the system the cursor is no longer needed. The cursor must not be used after this call is made.

#### **Reading a List**

Use the *alenlist\_get()* function to retrieve pairs from the list. An address/length pair from *plist* is stored based on *paddr* and *plength*. The pair to retrieve is established by a cursor, either *icursor* or the implicit cursor in *plist* when *icursor* is NULL. The cursor used is updated by the length returned, provided the AL\_LEAVE\_CURSOR flag is not used.

It is not necessary to read out exactly the pairs that were added to the list. When *maxlength* is nonzero, it establishes the maximum length retrieved. When *maxlength* is also an integral power of 2, the returned length is further constrained so that the returned address and length do not cross a *maxlength* boundary. For example, when *maxlength* is 512, the address/length values returned are such that the next pair returned will begin on a 512-boundary.

Returns ALENLIST\_SUCCESS or ALENLIST\_FAILURE. The normal cause for failure is that the list is exhausted.

Call `alenlist_cursor_init()` to initialize a cursor to a specified *offset* (usually 0). If *cursorp* is NULL, the implicit internal cursor of *plist* is initialized.

To retrieve the current byte offset of a cursor, call `alenlist_cursor_offset()`. When *icursor* is NULL, the offset of the implicit cursor in *plist* is returned.

### Loading a List

The basic operation to add an address/length pair to a list is `alenlist_append()`. The list is expanded if necessary. Use the AL\_NOSLEEP flag to prevent sleeping on memory allocation that might be necessary to resize the list. Use AL\_NOCOMPACT to prevent the added pair from being merged into a logically-adjacent preceding pair. Returns ALENLIST\_SUCCESS or ALENLIST\_FAILURE.

To build an *alenlist* that describes a buffer in kernel virtual memory, call `kvaddr_to_alenlist()`. The specified kernel virtual address is converted into a list of physical address/length pairs and the pairs are appended to an *alenlist*. A handle to the *alenlist* is returned. When *plist* is NULL, a new list is created, and this list is returned.

To build an *alenlist* that describes a buffer in user virtual memory, call `uvaddr_to_alenlist()`. The specified user virtual address for the specified length is converted into a list of physical address/length pairs and the pairs are appended to an *alenlist*. A handle to the *alenlist* is returned. When *plist* is NULL, a new list is created, and this list is returned.

To build an *alenlist* that describes a buffer mapped by a *buf* structure, call `buf_to_alenlist()`. The memory described by the *buf* is converted into a list of physical address/length pairs and the pairs are appended to an *alenlist*. A handle to the *alenlist* is returned. When *plist* is NULL, a new list is created, and this list is returned.

### SEE ALSO

`alenlist(D4X)`, `buf(D4)`, IRIX Device Driver Programmer's Guide

### NOTES

In IRIX 6.3, the declaration of `buf_to_alenlist()` was omitted from the header file. Declare it manually as follows:

```
/* temporary */
extern alenlist_t buf_to_alenlist(buf_t *);
```

This prototype for `buf_to_alenlist()`, as well as the prototypes shown above for `uvaddr_to_alenlist()`, `kvaddr_to_alenlist()`, `alenlist_append()`, and `alenlist_get()` are all different in future releases of IRIX. There is no simple workaround. These functions are extremely useful, but you should place a comment near each one indicating that when porting to IRIX 6.4 or later, additional arguments are required.

## PCI Infrastructure Reference Pages

- Example B-3 displays the reference page for registering and unregistering a PCI driver.
- Example B-4 displays the reference page about configuration space access.
- Example B-5 displays the reference page about setting up and using DMA maps.
- Example B-6 displays the reference page documenting the PCI error handler.
- Example B-7 displays the reference page for PCI query functions.
- Example B-8 displays the reference page documenting the PCI interrupt handler.
- Example B-9 displays the reference page about setting up and using PIO maps.

### Example B-3 pciio(d3)

#### NAME

pciio: pciio\_add\_attach, pciio\_driver\_register, pciio\_driver\_unregister - control PCI driver infrastructure

#### SYNOPSIS

```
#include <sys/PCI/pciio.h>

int
pciio_add_attach(__int32_t (*attach)(vertex_hdl_t),
                __int32_t (*detach)(vertex_hdl_t),
                pciio_error_handler_f *error,
                char *driver_prefix,
                int major)

int
pciio_driver_register(
    pciio_vendor_id_t vendor_id,
    pciio_device_id_t device_id,
    char *driver_prefix,
    unsigned flags);
```

```
void
pciio_driver_unregister(char *driver_prefix);

void
pciio_reset (vertex_hdl_t pconn);
```

**Arguments**

*attach*        Address of the driver's *attach()* entry point.  
*detach*        Address of the driver's *detach()* entry point, or NULL.  
*error*         Address of the driver's *error()* entry point, or NULL.  
*driver\_prefix*  
                The prefix string for the driver's standard entry points as  
                configured in */var/sysgen/system*.  
*major*         Major device number configured for this driver.  
*vendor\_id*  
*device\_id*     Values that the PCI device will present in its configuration  
                space as its vendor and device ID codes.  
  
*flags*         Normally passed as zero.  
  
*pconn*         is an appropriate PCI connection point.

**DESCRIPTION**

The PCI infrastructure is a package of kernel services used by drivers for PCI devices to set up services for their devices. These services include:

- o     Manipulating the PCI configuration space for the device (see *pciio\_config(D3)*).
- o     Constructing physical addresses to use for PIO access to the device (see *pciio\_pio(D3)*).
- o     Constructing PCI addresses for the device to use for DMA access to memory (see *pciio\_dma(D3)*).
- o     Arranging for a function to be called when the device requests interrupt service (see *pciio\_intr(D3)*).
- o     Arranging for a function to be called when an error occurs during PIO to, or DMA from the device (see *pciio\_error(D3)*).
- o     Accessing useful fields in some otherwise opaque data structures (see *pciio\_get(D3)*).

### Driver Registration

The first function a PCI driver must call, usually in the *init()* entry point, is *pciio\_add\_attach()*, to introduce the driver to the PCI infrastructure.

NOTE: This function is only called in the IRIX 6.3 implementation on the O2 workstation. The call to this function can be enclosed in a conditional test of the variable `_EARLY_PCI`, as follows:

```
xxxx_init() {
#ifdef _EARLY_PCI /* 6.3, must do add_attach */
    pciio_add_attach(xxxx_attach,NULL,NULL,"xxxx",XXXX_MAJNO);
#endif
}
```

*pciio\_driver\_register()* is used by a PCI driver in IRIX 6.3 or any later release to inform the infrastructure that it handles all PCI devices designated by specified *device\_id* and *vendor\_id* values. The infrastructure associates the specified ID numbers with the specified device driver prefix. When a device with these IDs is discovered, the infrastructure calls the attach entry point for the driver with that driver prefix, passing the hardware graph connection point vertex as the only parameter. This connection point is then used in most calls to the infrastructure to identify the PCI device of interest.

A loadable device driver calls *pciio\_driver\_register()* from its *reg()* entry point. A driver prelinked into the kernel should also make the call from *reg()* for consistency, but may call from the *init()* entry point if necessary.

Device drivers may make multiple calls with different vendor and device ID numbers, representing several compatible PCI devices.

Wildcard values `PCIIO_VENDOR_ID_NONE` and `PCIIO_DEVICE_ID_NONE` may be used if cards from any vendor or cards with any device code are of supported. When both vendor and device are wildcarded, the *attach()* routine is called for every PCI device connected to the system.

When a loadable device driver calls *pciio\_driver\_register()*, one or more calls to the driver's *attach()* function can occur before the infrastructure returns control to the caller. On some large systems, the *attach()* calls can be executed by other threads and possibly on other processors, concurrently with continued execution of the *reg()* entry point.

`pciio_driver_unregister()` should be called by any unloadable device driver from within the driver's `unreg()` entry point. This triggers calls to the driver's `detach()` entry point, and removes the association between the driver and any vendor and device IDs.

#### Resetting a PCI card

`pciio_reset()` is used to attempt to activate the PCI Reset line connected to a specific card without affecting other devices on the PCI bus. When reset is possible, the device is reset and basic configuration information is reloaded.

#### EXAMPLES

Here is how a typical driver might make use of these functions:

```
static char    pcifoo_prefix[] = "pcifoo_";
static char    pcifoo_edge[] = "foo";
#ifdef _EARLY_PCI /* 6.3, must do add_attach, need init() */
pcifoo_init(void)
{
    pciio_add_attach(pcifoo_attach,NULL,NULL,"pcifoo",42);
}
#endif
pcifoo_reg(void)
{
    pciio_driver_register(
        PCIFOO_VENDOR_ID,
        PCIFOO_DEVICE_ID,
        "pcifoo", 0);
}
pcifoo_unreg(void)
{
    pciio_driver_unregister("pcifoo");
}
```

#### SEE ALSO

`pciio_config(D3)`, `pciio_dma(D3)`, `pciio_error(D3)`, `pciio_get(D3)`, `pciio_intr(D3)`, `pciio_pio(D3)`.

#### Example B-4 pciio\_config(d3)

#### NAME

`pciio_config`: `pciio_config_get`, `pciio_config_set` - access PCI Configuration register

**SYNOPSIS**

```
#include <sys/PCI/pciio.h>

u_int32_t
pciio_config_get(
    volatile unsigned char *bus_addr,
    int cfg_reg)

int
pciio_config_set(
    volatile unsigned char *bus_addr,
    int cfg_reg,
    __int32_t value)
```

**Arguments**

*bus\_addr*     PIO target address of PCI configuration space for a device, as returned by `pciio_piomap_addr()` or `pciio_piotrans_addr()`.

*cfg\_reg*     Byte offset of the register of interest in the PCI address space.

*value*       Value to be written to the specified register.

**DESCRIPTION**

Various SGI platforms introduce complexities and restrictions in how Configuration Space cycles are generated on the PCI bus. Some platforms may require all PCI Configuration accesses to be done using 32-bit wide accesses. Others require more than a simple load or store to trigger the actual cycle, so that configuration access cannot be performed using normal PIO loads and stores. (Both of these restrictions are true of the O2 workstation.)

The functions described here were introduced to allow the kernel to trigger a PCI bus Configuration Cycle based on a PIO address.

The *cfg\_reg* value specifies the offset of the target value in configuration space. Registers defined by the PCI standard are 1, 2, 3, 4, or 8 bytes, but these functions support only 1-4 bytes. (Eight-byte registers can be fetched in two calls.)

Because configuration space is accessed in 32-bit units on 32-bit boundaries, when *reg* specifies a standard PCI configuration register, `pciio_config_get()` shifts and masks appropriately to return just the value of the register. Similarly, `pciio_config_set()` executes a read-merge-write operation to place the value data in the correct portion of the word.

### Standard PCI Configuration Registers

To access vendor-specific registers, specify the base address in PCI configuration space, bearing in mind that PCI places the least significant data in the lowest offset.

The following constants are declared in the header file `sys/PCI/PCI_defs.h` for use as the `cfg_reg` value to specify a standard register in the Type 00 PCI configuration space:

```
PCI_CFG_VENDOR_ID
PCI_CFG_DEVICE_ID
PCI_CFG_COMMAND
PCI_CFG_STATUS

PCI_CFG_REV_ID
PCI_CFG_BASE_CLASS
PCI_CFG_SUB_CLASS
PCI_CFG_PROG_IF

PCI_CFG_CACHE_LINE
PCI_CFG_LATENCY_TIMER
PCI_CFG_HEADER_TYPE
PCI_CFG_BISTPCI_CFG_BIST

PCI_CFG_BASE_ADDR_0
PCI_CFG_BASE_ADDR_1
PCI_CFG_BASE_ADDR_2
PCI_CFG_BASE_ADDR_3
PCI_CFG_BASE_ADDR_4
PCI_CFG_BASE_ADDR_5
PCI_CFG_BASE_ADDR(n)

PCI_CFG_CARDBUS_CIS
PCI_CFG_SUBSYS_VEND_ID
PCI_CFG_SUBSYS_ID
PCI_EXPANSION_ROM

PCI_INTR_LINE
PCI_INTR_PIN
PCI_MIN_GNT
PCI_MAX_LAT
```

Use `PCI_CFG_VEND_SPECIFIC` to specify the first vendor-specific register word.

The configuration functions deduce the length of the register from its offset. For a nonstandard or vendor-specific register the functions assume a 32-bit word.

#### NOTES

Logical byte order is preserved; that is, the most significant byte of a (big-endian) program word is sent as the most significant byte of a (little-endian) PCI bus word.

Writes to registers of less than 32 bits are synthesized by reading the word containing the register, modifying the proper bits in the word, then rewriting the entire bus word. The read-modify-write code knows about the special handling of the STATUS register. However, if other registers in your card's configuration space are sensitive to being rewritten, you should access them using full four-byte-wide accesses, manipulating the word data appropriately.

#### Subsequent Releases

In all platforms supported IRIX 6.4, configuration access is possible with normal PIO. Accordingly, these functions were not defined in IRIX 6.4 and will result in an unresolved extern if porting is attempted. In releases after IRIX 6.4, these configuration access functions are returned, but with different arguments and with more capabilities:

- o They take the register size as an argument, from 1-8 bytes.
- o They do not require use of a mapped PIO address, but take only the vertex handle of the device.

It is possible to code configuration access macros so that they compile properly in all releases from 6.3 onward. The macro code would be similar to the following:

```
/* PCI Config Space Access Macros
** for source compatibility in drivers
** that need to use the same source
** for IRIX 6.3, IRIX 6.4, and IRIX 6.5
**
** PCI_CFG_BASE(conn)
** PCI_CFG_GET(conn,base,offset,type)
** PCI_CFG_SET(conn,base,offset,type,value)
**
** Use PCI_CFG_BASE once during attach to get the
** base value to be used for the specific device.
```

```

** Later, use PCI_CFG_GET to read and PCI_CFG_SET
** to write config registers.
**
** NOTE: Irix 6.3 determines the size of the register
** directly on its own, based on the layout of a Type 00
** PCI Configuration Space Header. If you specify a
** nonstandard size, you will get different results
** depending on the system revision number.
*/
#if IRIX6_3
#define PCI_CFG_BASE(c)          pciio_piотrans_addr(c,0,PCIIO_SPACE_CFG,0,256,0)
#define PCI_CFG_GET(c,b,o,t)    pciio_config_get(b,o)
#define PCI_CFG_SET(c,b,o,t,v)  pciio_config_set(b,o,v)
#elif IRIX6_4
#define PCI_CFG_BASE(c)          pciio_piотrans_addr(c,0,PCIIO_SPACE_CFG,0,256,0)
#define PCI_CFG_GET(c,b,o,t)    ((*t *)((char *) (b)+(o)))
#define PCI_CFG_SET(c,b,o,t,v)  ((*t *)((char *) (b)+(o))) = v
#else /* starting in IRIX 6.5 */
#define PCI_CFG_BASE(c)          NULL
#define PCI_CFG_GET(c,b,o,t)    pciio_config_get(c,o,sizeof(t))
#define PCI_CFG_SET(c,b,o,t,v)  pciio_config_set(c,o,sizeof(t),v)
#endif

```

The macros would be used approximately as follows:

```

pcifoo_attach(vertex_hdl_t conn)
{
    void * config_base = PCI_CFG_BASE(conn);
    ...
    /* retrieve current device revision */
    foo_soft->fs_revision =
        PCI_CFG_GET(conn, config_base, PCI_CFG_REV_ID, uchar);
    ...
    /* write 0x5555AAAA test pattern to first
    ** vendor specific register */
    PCI_CFG_SET(conn, config_base, PCI_CFG_VEND_SPECIFIC, uint32_t,
        0x5555AAAA);
}

```

**SEE ALSO**

pciio(D3), pciio\_config(D3), pciio\_dma(D3), pciio\_error(D3),  
pciio\_get(D3), pciio\_intr(D3). pciio\_pio(D3).

**Example B-5** pciio\_dma(d3)**NAME**

```
pciio_dma: pciio_dmatrans_addr, pciio_dmatrans_list, pciio_dmamap_alloc,  
pciio_dmamap_addr, pciio_dmamap_list, pciio_dmamap_done,  
pciio_dmamap_free, - manage DMA on PCI bus
```

**SYNOPSIS**

```
#include <sys/PCI/pciio.h>
```

```
iopaddr_t
```

```
pciio_dmatrans_addr(  
    vertex_hdl_t vhdl,  
    device_desc_t desc,  
    iopaddr_t addr,  
    size_t size,  
    unsigned flags)
```

```
alenlist_t
```

```
pciio_dmatrans_list(  
    vertex_hdl_t vhdl,  
    device_desc_t desc,  
    alenlist_t list,  
    unsigned flags)
```

```
pciio_dmamap_t
```

```
pciio_dmamap_alloc(  
    vertex_hdl_t vhdl,  
    device_desc_t desc,  
    size_t max,  
    unsigned flags)
```

```
iopaddr_t
```

```
pciio_dmamap_addr(  
    pciio_dmamap_t map,  
    iopaddr_t addr,  
    size_t size);
```

```
alenlist_t
```

```
pciio_dmamap_list(  
    pciio_dmamap_t map,  
    alenlist_t list);
```

```
void
```

```
pciio_dmamap_done(pciio_dmamap_t map)
```

```
void  
pciio_dmamap_free(pciio_dmamap_t map)
```

**Arguments**

*addr* The DMA buffer address in system physical address space.

*desc* A device descriptor, usually zero.

*flags*

Attributes of the mapping.

*list* An address/length list as prepared by one of the *alenlist* construction functions (see *alenlist(D4)*).

*map* A dma map as returned by *pciio\_dmamap\_alloc()*.

*max* The maximum range of addresses this map will cover at any one time.

*size* The size of the mapped buffer in bytes.

*vhdl* The device connection point as passed to the *attach()* entry point.

**DESCRIPTION**

When a device driver wishes to use Direct Memory Access (DMA) to communicate with a device, the system needs to have a chance to set up any appropriate mapping registers. The work to be done varies with the available hardware and with the version of IRIX.

The functions described here provide an abstract interface to the creation of DMA mapping objects that is consistent across most hardware. These functions always do the least possible work given the available hardware. (In IRIX 6.3 and on the O2 hardware, the amount of work is minimal. In later releases and on multiprocessor platforms, the work can be considerable.)

There are two different models for setting up a DMA map, one simple but fallible and the other more general. In both models, the final goal is to retrieve an address in PCI bus address space that can be used by a PCI device to write into, or read from, system physical memory.

**Simple Model**

The simple model provides permanent mappings through fixed mapping resources that may or may not exist in a given system at a given time.

`pciio_dmatrans_addr()` is the one-stop shopping place for using system fixed shareable mapping resources to construct a DMA address. Such resources are not always available, in which case, the function returns `NULL`.

`pciio_dmatrans_list()` is similar, but operates on an address/length list of blocks of memory and returns a list of blocks in PCI address space.

When they work, these functions allow the driver to set up DMA with the fewest complications. Typically the functions always succeed in some platforms (those having simple hardware mappings of PCI to memory), and always fail in other platforms (where multiple layers of hardware mappings must be configured dynamically). However, drivers that hope to be portable must be coded as if the functions could succeed or fail alternately in the same system.

#### **General Model**

It is not always possible to establish DMA mappings using common shared system resources, so the concept of a DMA channel that preallocates scarce mapping resources is provided.

Such a channel is allocated using `pciio_dmamap_alloc()`, which is given the maximum size to be mapped. `pciio_dmamap_addr()` or `pciio_dmamap_list()` is then applied to the map to actually establish the proper mappings for a DMA target. Given the base address and block size of the buffer for DMA (or a list of buffers), the functions hand back the base PCI address to use for accessing that buffer (or a list of PCI addresses).

When all DMA to a given buffer (or list) is complete, `pciio_dmamap_done()` should be called to idle any mapping hardware (and possibly flush out any pipes or buffers along the path that might do unexpected things when mapping registers are modified). Later, `pciio_dmamap_addr()` or `pciio_dmamap_list()` can again be called, specifying the same or a different buffer area.

When a driver is completely finished with a DMA channel -- because the channel is used only for initialization of the device, because the driver's `close()` entry point is called so it is known that the device will be idle for some time, or because the device or the driver is being shut down -- the DMA channel resources should be released using `pciio_dmamap_free()`.

#### **DMA Attribute Flags**

The following attributes can be specified in the flags argument:

- PCIIO\_DMAMAP\_CMD This map is primarily for transfer of device-to-driver command and status data. (Flag ignored in this release.)
- PCIIO\_DMAMAP\_DATA This map is primarily for transfer of user data. (Flag ignored in this release.)
- PCIIO\_DMAMAP\_A64 The device is capable of using PCI-64 addressing. (Flag ignored in this release.)
- PCIIO\_DMAMAP\_LITTLEEND  
demands that any byte-swapping hardware along this DMA path be organized so that an ordered stream of bytes from the device are deposited in order in system memory. This is the typical setting for data streams. If this endianness cannot be supplied, then the service call fails.
- PCIIO\_DMAMAP\_BIGEND  
demands that any byte-swapping hardware along this DMA path be initialized so that 32-bit quantities on PCI-bus 32-bit boundaries maintain their binary values. This is the typical setting for command-type transactions because command words exchanged with a little-endian PCI device retain their binary values. If this endianness cannot be supplied, then the service call fails.

When PCIIO\_DMAMAP\_LITTLEEND is used, the bytes of multibyte values embedded in input data are found at their original offsets. Multibyte values from little-endian devices may require programmed swapping before use.

When PCIIO\_DMAMAP\_BIGEND is used,

- o Single bytes in input data are found at the offset the device places them, exclusive-or with 3.
- o 16-bit quantities in input data are found at the offset used by the device, exclusive-or with 2, and do not need to be byteswapped.
- o 32-bit values are found at the expected offset, and do not need to be byteswapped.

- o 64-bit values are found at the expected offset, and their 32-bit halves need to be swapped before use.

#### Source Compatibility With IRIX 6.4

In IRIX 6.4 and subsequent releases, the flag names given above are renamed, and additional useful flags are defined and supported. For easy source compatibility with later releases, insert the following code:

```
#ifdef _EARLY_PCI
# define PCIIIO_DMA_CMD PCIIIO_DMAMAP_CMD
# define PCIIIO_DMA_DATA PCIIIO_DMAMAP_DATA
# define PCIIIO_DMA_A64 PCIIIO_DMAMAP_A64
# define PCIIIO_BYTE_STREAM PCIIIO_DMAMAP_LITTLEEND
# define PCIIIO_WORD_VALUES PCIIIO_DMAMAP_BIGEND
#endif
```

Then use the defined names such as `PCIIIO_BYTE_STREAM` when creating DMA maps, instead of the flag values mentioned above.

When porting to a later release of IRIX, be sure to read this reference page for that release to see what flags are supported.

In IRIX 6.4 and later releases, `pciio_dmamap_list()` takes a third argument, `flags`, with the same meaning as the other flags arguments. In order to simplify source compatibility with later releases, you could use the `_EARLY_PCI` identifier to code a macro in this form:

```
#ifdef _EARLY_PCI
# define PCIIIO_DMAMAP_LIST(a,b) pciio_dmamap_list(a,b)
#else
# define PCIIIO_DMAMAP_LIST(a,b) pciio_dmamap_list(a,b,0)
#endif
```

However, the flags supported in later releases are sufficiently useful you should probably recode the calls to use nonzero flags.

#### EXAMPLES

Here is one way that a driver might make use of `dmamap` and `dmatrans` calls.

```
#ifdef _EARLY_PCI
# define PCIIIO_DMA_CMD PCIIIO_DMAMAP_CMD
# define PCIIIO_DMA_DATA PCIIIO_DMAMAP_DATA
# define PCIIIO_DMA_A64 PCIIIO_DMAMAP_A64
# define PCIIIO_BYTE_STREAM PCIIIO_DMAMAP_LITTLEEND
```

```
# define PCIIO_WORD_VALUES PCIIO_DMAMAP_BIGEND
#endif
pcifoo_attach(vertex_hdl_t vhdl)
{
    pciio_dmamap_t command_map;
    iopaddr_t command_dma;
    struct pcifoo_regs *reg_pio;
    struct pcifoo_ring *command_ring;
    ...
    /*
     * This driver has decided to use a dmamap
     * to get to its command rings, which contain
     * things like DMA addresses and counts; we
     * set PCIIO_WORD_VALUES so we don't have to
     * byteswap the 32-bit values.
     *
     * We still have to swap the upper and lower
     * halves of the 64-bit values.
     */
    /* allocate the channel
     */
    command_map = pciio_dmamap_alloc(
        vhdl, 0,
        RINGBYTES,
        PCIIO_DMA_CMD | PCIIO_WORD_VALUES);
    command_dma = pciio_dmamap_addr(
        command_map,
        kvtophys(command_ring),
        RINGBYTES);
    /* tell the device where it can find
     * it's command rings.
     */
    reg_pio->command_dma = command_dma;
    ...
}
{
    caddr_t data_buffer;
    size_t data_size;
    ...
    data_dma = pciio_dmatrans_addr(
        vhdl, 0,
        kvtophys(data_buffer), data_size,
        PCIIO_DMA_DATA | PCIIO_DMA_A64 | PCIIO_BYTE_STREAM);
    command_ring->data_dma_lo = data_dma & 0xFFFFFFFF;
    command_ring->data_dma_hi = data_dma >> 32;
}
```

```

        command_ring->data_dma_size = data_size;
        command_ring->ready = 1;
    }

```

**SEE ALSO**

alenlist(D3), pciio(D3), pciio\_config(D3), pciio\_error(D3),  
pciio\_get(D3), pciio\_intr(D3), pciio\_pio(D3).

**DIAGNOSTICS**

*pciio\_dmatrans\_addr()* returns zero if shared (fixed) resources can not be used to construct a valid PCI address that maps to the desired range of physical addresses. (Fixed resources are always available in IRIX 6.3 for O2, but may not be in other systems.)

*pciio\_dmatrans\_list()* returns a null pointer if any of the requested physical address blocks can not be reached using shared fixed resources, or if unable to allocate a return list.

*pciio\_dmamap\_alloc()* returns a null pointer if resources can not be allocated to establish DMA mappings of the requested size, or if the parameters are inconsistent.

*pciio\_dmamap\_addr()* returns zero if the specified target address can not be mapped using the specified DMA channel. This would usually be due to specifying a target block that is outside the previously specified target area or is larger than the previously specified maximum mapping size. It may also return a null pointer if the DMA channel is currently in use and has not been marked idle by a call to *pciio\_dmamap\_done()*.

*pciio\_dmamap\_list()* can return a null pointer for all the reasons mentioned above, or if it is unable to allocate the return list.

**Example B-6** pciio\_error(d3)**NAME**

pciio\_error - IRIX 6.3 PCI error interface

**SYNOPSIS**

```

#include <sys/PCI/pciio.h>

int
typedef int
(pciio_error_handler_f)( vertex_hdl_t   vhdl,
                        int             error_code,
                        ioerror_mode_t mode,
                        ioerror_t      *ioerror);

```

**Arguments**

- vhdl*      The connection point of the PCI device with the error.
  
- error\_code*  
           Bit-set describing the error type.
  
- mode*      Code for the type of device access when the error was detected.
  
- ioerror*   Error mode structure with more information.

**DESCRIPTION**

A PCI device driver can pass the address of an error-handling function to the *pciio\_add\_attach()* function documented in *pciio(d3)*. The error-handling function must have the prototype shown; that is, it must agree with type *pciio\_error\_handler\_f*. When a NULL is passed to *pciio\_add\_attach()*, the PCI infrastructure handles all errors.

When an error occurs, the handler is called. The *error\_code* value contains a set of the bits defined in *sys/mace.h*, as follows:

```
#define PERR_MASTER_ABORT          0x80000000
#define PERR_TARGET_ABORT          0x40000000
#define PERR_DATA_PARITY_ERR       0x20000000
#define PERR_RETRY_ERR             0x10000000
#define PERR_ILLEGAL_CMD           0x08000000
#define PERR_SYSTEM_ERR            0x04000000
#define PERR_INTERRUPT_TEST        0x02000000
#define PERR_PARITY_ERR            0x01000000
#define PERR_OVERRUN               0x00800000
#define PERR_RSVD                  0x00400000
#define PERR_MEMORY_ADDR           0x00200000
#define PERR_CONFIG_ADDR           0x00100000
#define PERR_MASTER_ABORT_ADDR_VALID 0x00080000
#define PERR_TARGET_ABORT_ADDR_VALID 0x00040000
#define PERR_DATA_PARITY_ADDR_VALID 0x00020000
#define PERR_RETRY_ADDR_VALID      0x00010000
```

The *mode* value is one of the following, defined in *sys/PCI/pci\_compat.h* (which in turn is included by *sys/PCI/pciio.h*):

```
typedef enum {
    MODE_DEVPROBE,          /* Probing mode. Errors not fatal */
    MODE_DEVEERROR,        /* Error while system is running */
    MODE_DEVREENABLE       /* Reenable pass          */
}ioerror_mode_t;
```

The structure addressed by *ioerror* is also declared in `sys/PCI/pci_compat.h` and contains numerous fields that may give the handler additional information.

The handler returns nonzero when it has handled the error adequately. The handler returns zero when it wants the PCI infrastructure to handle the error.

**NOTES**

This error-handling interface is unique to IRIX 6.3 for O2. In IRIX 6.4 and onward, a different interface is used: the driver calls a function `pciio_error_register()`, to register an error handler that takes fewer, and different, arguments.

The error-handling function in IRIX 6.3 receives information that is unique to the hardware of the O2 workstation. For example, all the declarations in `sys/mace.h` are unique to the MACE chip that manages bus and memory access in the O2. None of this information is relevant in other platforms such as the OCTANE workstation. None of this information is available in later releases of IRIX.

There is no simple way to make PCI error-handling in IRIX 6.3 source-compatible with later releases.

**SEE ALSO**

`pciio(D3)`, `pciio_config(D3)`, `pciio_dma(D3)`, `pciio_get(D3)`, `pciio_intr(D3)`, `pciio_pio(D3)`.

**Example B-7** `pciio_get(d3)`**NAME**

`pciio_get`: `pciio_info_get`, `pciio_info_bus_get`, `pciio_intr_cpu_get`, `pciio_dma_dev_get`, `pciio_info_dev_get`, `pciio_intr_dev_get`, `pciio_pio_dev_get`, `pciio_dma_slot_get`, `pciio_info_slot_get`, `pciio_pio_slot_get`, `pciio_pio_mapsz_get`, `pciio_pio_pciaddr_get`, `pciio_pio_space_get`, `pciio_info_vendor_id_get`, `pciio_info_device_id_get`, `pciio_info_function_get` - interrogate PCI infrastructure

**SYNOPSIS**

```
#include <sys/PCI/pciio.h>

vertex_hdl_t
pciio_intr_cpu_get(pciio_intr_t intr)
```

```
vertex_hdl_t
pciio_intr_dev_get(pciio_intr_t intr)

vertex_hdl_t
pciio_pio_dev_get(pciio_piomap_t piomap)

pciio_slot_t
pciio_pio_slot_get(pciio_piomap_t piomap)

pciio_space_t
pciio_pio_space_get(pciio_piomap_t piomap)

iopaddr_t
pciio_pio_pciaddr_get(pciio_piomap_t piomap)

ulong
pciio_pio_mapsz_get(pciio_piomap_t piomap)

vertex_hdl_t
pciio_dma_dev_get(pciio_dmamap_t dmamap)

pciio_slot_t
pciio_dma_slot_get(pciio_dmamap_t dmamap)

pciio_info_t
pciio_info_get(vertex_hdl_t vhdl)

vertex_hdl_t
pciio_info_dev_get(pciio_info_t info)

pciio_slot_t
pciio_info_bus_get(pciio_info_t info)

pciio_function_t
pciio_info_function_get(pciio_info_t info)

pciio_slot_t
pciio_info_slot_get(pciio_info_t info)

pciio_vendor_id_t
pciio_info_vendor_id_get(pciio_info_t info)

pciio_device_id_t
pciio_info_device_id_get(pciio_info_t info)
```

**Arguments**

- intr* A PCI interrupt object handle returned by *pciio\_intr\_alloc()*.
- piomap* A PCI PIO map returned by *pciio\_piomap\_alloc()*.
- dmamap* is a *pciio\_dmamap\_t* that was created by *pciio\_dmamap\_alloc()*
- vhdl* A pci connection point in the hardware graph, obtained as the parameter to the attach call.
- info* A PCI info object returned by *pciio\_info\_get()*.

**DESCRIPTION**

These routines are used to pull specific useful bits of information out of the various opaque data structures used by the PCI infrastructure. Few drivers will need to make use of these routines, but having them available might save the driver from doing extra bookkeeping.

**Interrupt Queries**

Two functions fetch parameters from an interrupt object:

- o *pciio\_intr\_dev\_get()* returns the connection point of the interrupt device.
- o *pciio\_intr\_cpu\_get()* returns the CPU that is the target of interrupts for that PCI bus.

**PIO Map Queries**

Several functions return items based on a PIO map (see *pciio\_pio(D3)*):

- o *pciio\_pio\_dev\_get()* returns the connection point of the mapped device.
- o *pciio\_pio\_mapsz\_get()* returns the map maximum size.
- o *pciio\_pio\_pciaddr\_get()* returns the base address specified for the map.
- o *pciio\_pio\_space\_get()* returns the target address space that was specified.
- o *pciio\_pio\_slot\_get()* returns the slot number on the PCI bus for a device.

#### DMA Map Queries

Two functions return items based on a DMA as map (see *pciio\_dma(D3)*):

- o *pciio\_dma\_dev\_get()* returns the connection point of the mapped device.
- o *pciio\_dma\_slot\_get()* returns the slot number on the PCI bus for a device.

#### Info Structure Queries

The PCI infrastructure stores a version-dependent information structure in the connection point for a PCI device. Several functions are provided to retrieve and interrogate this structure. Those most likely to be useful to a device driver are:

- o *pciio\_info\_get()* returns a handle to the information structure. The driver can save this handle at attach time to avoid the small overhead of looking it up each time it is needed.
- o *pciio\_info\_dev\_get()* returns the vertex handle of the connection point (from which the information structure was originally retrieved).
- o *pciio\_info\_bus\_get()* returns the bus number, always 0 unless the system has more than one PCI bus. Bus numbers are arbitrary, not necessarily sequential.
- o *pciio\_info\_slot\_get()* returns the PCI card slot number of the device.
- o *pciio\_info\_function\_get()* returns the PCI function number, 0 unless the device is part of a multifunction card.
- o *pciio\_info\_vendor\_id\_get()* returns the vendor ID configuration value of the device.
- o *pciio\_info\_device\_id\_get()* returns the device ID configuration value of the device.

#### SEE ALSO

*pciio(D3)*, *pciio\_config(D3)*, *pciio\_dma(D3)*, *pciio\_error(D3)*,  
*pciio\_intr(D3)*, *pciio\_pio(D3)*.

#### DIAGNOSTICS

*pciio\_info\_get()* returns NULL if there is no *pciio* info structure attached to that vertex.

Do not pass info as NULL to any of these functions, that would cause a kernel panic.

**Example B-8**    pciio\_intr(d3)

**NAME**

pciio\_intr: pciio\_intr\_alloc, pciio\_intr\_connect, pciio\_intr\_disconnect, pciio\_intr\_free - manage PCI Interrupts

**SYNOPSIS**

```
#include <sys/PCI/pciio.h>

pciio_intr_t
pciio_intr_alloc(
    vertex_hdl_t vhdl,
    device_desc_t desc,
    pciio_intr_line_t lines,
    vertex_hdl_t owner)

int
pciio_intr_connect(
    pciio_intr_t intr,
    intr_func_t func,
    intr_arg_t arg,
    void *thread)

void
pciio_intr_disconnect(pciio_intr_t intr)

void
pciio_intr_free(pciio_intr_t intr)
```

**Arguments**

*arg*        A parameter to pass to func() when this particular interrupt occurs, commonly a pointer to a driver-private data structure.

*desc*       A device descriptor containing an interrupt priority level.

*func*       The function to perform interrupt service.

*intr*       The interrupt channel handle returned by pciio\_intr\_alloc().

*lines*      Specifies one or more of the PCI Interrupt pins used by the device.

*owner* An appropriate vertex handle to use when printing messages about this particular interrupt, and is usually a vertex created by the device driver.

*vhdl* The PCI device connection point as passed to the driver `attach()` entry point.

*thread* Reserved, should be NULL.

#### DESCRIPTION

When a device driver wishes to accept interrupt events from a device, the system needs to make sure that there is a path from the PCI interrupt pin to the appropriate CPU interrupt hardware. This is split into two phases -- establishing the channel and connecting a service function -- so that the service function can be changed or disconnected without losing the allocated hardware resources.

The driver is responsible for connecting an interrupt handler when the device needs one, and for disconnecting the handler when it does not.

The interrupt delivery mechanism depends on the address of the interrupt function. It is important to disconnect interrupts before a driver unloads, otherwise the PCI infrastructure might call a nonexistent function. (A driver cannot be auto-loaded when an interrupt occurs.)

The necessary sequence of calls is based on the use of the driver entry points, as follows.

At the `reg()` or `init()` entry point the driver registers to handle a class of PCI devices, triggering `attach()` calls.

At the `attach()` entry, the driver calls `pciio_intr_alloc()` to establish interrupt connectivity between the device and the processor. The designated interrupts are disabled at this point. If interrupts can occur and are needed at this time, a call to `pciio_intr_connect()` enables interrupts and directs them to the designated handler.

At the `unload()` entry, the driver text is going to be removed, so it is important for all interrupts to be disconnected by calling `pciio_intr_disconnect()` as appropriate. It is not necessary to call `pciio_intr_free()` at this time.

Some devices do not require interrupt service when they are not open. Leaving an interrupt allocated but not connected keeps the interrupt

disabled, possibly reducing impact on the system from handling interrupts from devices that do not actually need service.

If this is the situation, then the scenario above may be somewhat simplified:

```
attach()    Allocate the interrupt to establish a connection and disable
            the interrupt. Only connect the interrupt if interrupts are
            required as part of device initialization; then disconnect
            it.

open()      If the interrupt is not yet connected, connect it.

close()     No processes have the device open; disconnect the interrupt
            when all pending I/O is complete or purged.

unload()    The driver is not called to unload when one of its devices is
            open, so no interrupts should be connected.
```

#### Specifying PCI Interrupt Lines

The lines parameter is formed by or-ing together appropriate flags:

```
PCIIO_INTR_LINE_A
PCIIO_INTR_LINE_B
PCIIO_INTR_LINE_C
PCIIO_INTR_LINE_D
```

#### Specifying the Device Descriptor

The desc value must be the address of a `device_desc_t` structure in which the `intr_swlevel` field has been assigned a value. The data type of this field, `pl_t`, is declared in `sys/types.h`. The header `sys/ddi.h` includes `sys/types.h`, and also declares several external objects of type `pl_t`.

In IRIX 6.4 and later, it is not required to supply a device descriptor to this function; `NULL` may be passed instead. Also in IRIX 6.4 and later, a default device descriptor is readily available by a function call. The following example shows how to code the call to `pciio_intr_alloc()` in a source-compatible way:

```
foo_attach(vertex_hdl_t connpt)
{
#ifdef _EARLY_PCI
    device_desc_t work_desc = {0};
#endif
    device_desc_t *pdesc;
```

```
vertex_hdl_t foo_chardev; /* our char device */
pciio_intr_t intr;
...
/* allocate an interrupt object. This device
 * uses both INTA and INTB, and routes both
 * interrupts to the same function.
 */
#ifdef _EARLY_PCI
    pdesc = &work_desc;
    pdesc -> intr_swlevel = plhi; /* in ddi.h */
#else
    pdesc = NULL; /* or: pdesc = device_desc_default_get(connpt) */
#endif
intr = pciio_intr_alloc(connpt, pdesc,
    PCIIO_INTR_LINE_A | PCIIO_INTR_LINE_B,
    foo_chardev);
pciio_intr_connect(intr,
    foo_int_hdlr,
    (intr_arg_t)foo_dev_info,
    (void *)0);
...
```

**SEE ALSO**

`pciio(D3)`, `pciio_config(D3)`, `pciio_dma(D3)`, `pciio_error(D3)`,  
`pciio_get(D3)`, `pciio_pio(D3)`.

**DIAGNOSTICS**

`pciio_intr_alloc()` returns a null value if it can not allocate memory.

`pciio_intr_connect()` returns a zero for success or a negative value on failure. Since the channel is preallocated, the only interesting failure for this function is the attempt to use a null interrupt handle value.

**Example B-9** `pciio_pio(d3)`**NAME**

`pciio_pio`: `pciio_piotrans_addr`, `pciio_piomap_alloc`, `pciio_piomap_addr`,  
`pciio_piomap_done`, `pciio_piomap_free`, `pciio_piospace_alloc`,  
`pciio_piospace_free` - programmed I/O to PCI bus

**SYNOPSIS**

```
#include <sys/PCI/pciio.h>
```

```
caddr_t
pciio_piotrans_addr(
```

```
vertex_hdl_t vhdl,  
device_desc_t desc,  
pciio_space_t space,  
iopaddr_t addr,  
size_t size,  
unsigned flags)  
  
pciio_piomap_t  
pciio_piomap_alloc(  
    vertex_hdl_t vhdl,  
    device_desc_t desc,  
    pciio_space_t space,  
    iopaddr_t addr,  
    size_t size,  
    size_t max,  
    unsigned flags)  
  
caddr_t  
pciio_piomap_addr(  
    pciio_piomap_t map,  
    iopaddr_t addr,  
    size_t size);  
  
void  
pciio_piomap_done(pciio_piomap_t map)  
  
void  
pciio_piomap_free(pciio_piomap_t map)  
  
iopaddr_t  
pciio_piospace_alloc(  
    vertex_hdl_t vhdl,  
    device_desc_t desc,  
    pciio_space_t space,  
    size_t size,  
    size_t align)  
  
void  
pciio_piospace_free(  
    vertex_hdl_t vhdl,  
    pciio_space_t space,  
    iopaddr_t addr,  
    size_t size)
```

**Arguments**

- addr* The offset within the given space.
- align* A desired alignment in PCI address space.
- desc* A device descriptor with the one field *intr\_swlevel* set to *plhi*.
- flags* Flags describing the use of the PIO map.
- max* The maximum size within space to be mapped at any one time.
- map* The map address returned by *pciio\_piomap\_alloc()*.
- mapp* A pointer variable to receive the address of an allocated map.
- size* The size of the region to be mapped.
- space* Specifies the target PCI address space.
- vhdl* The PCI connection point as given to the *attach()* entry point.

**DESCRIPTION**

When a device driver wishes to use Programmed I/O (PIO) to communicate with a device, the system needs to have a chance to set up any appropriate mapping registers. The work to be done varies with the available hardware and with the version of IRIX. The functions described here provide an abstract interface that is consistent across most hardware. These functions always do the least possible work given the available hardware.

There are two models for setting up a PIO map, one simple but fallible, and one more general. In both models, the final goal is to retrieve a physical address that, when used as the operand of a store or fetch, will access a word in PCI bus address space rather than in CPU memory address space.

**Simple Model**

The simple model provides permanent mappings through fixed mapping resources that may or may not exist in a given system at a given time. *pciio\_piotrans\_addr()* attempts to use shared hardware resources to construct a physical address that, whenever used, routes the transaction to the proper target on the PCI bus. This is not always possible. When it is not, the function returns NULL.

When it works, *pciio\_piotrans\_addr()* allows the driver to do PIO with the fewest complications. Typically *pciio\_piotrans\_addr()* always succeeds in some platforms (those having a simple mapping of PCI bus to memory), and always fails in others (where multiple layers of hardware mappings must be configured dynamically). However, a driver that uses it should be coded as if it could succeed or fail alternately in the same system (which it could).

### General Model

It is not always possible to establish a PIO mapping using common shared system resources, so the concept of a PIO channel that preallocates scarce mapping resources is provided.

Such a channel is allocated using *pciio\_piomap\_alloc()*, which is given the limits of the region that will be mapped and the maximum size to be mapped at any time within that region. The model assumes that many channels may be created, but that not all channels will be actively in use at any time.

*pciio\_piomap\_addr()* is used to actually establish the proper mappings for a PIO target. Given the offset within the target address space and the size of the region for PIO, it returns the base address to be used for accessing that region.

After all PIO transactions to that region are executed, *pciio\_piomap\_done()* should be called to idle any mapping hardware and possibly to flush out any pipes or buffers along the path that might do unexpected things when mapping registers are modified.

Later, *pciio\_piomap\_addr()* can again be called, specifying the same or a new target area.

When a driver is completely finished with a PIO channel -- either because the channel is used only for initialization of the device, or because the device or the driver is being shut down -- the PIO channel resources should be released using *pciio\_piomap\_free()*.

### Utility Functions

*pciio\_piospace\_alloc()* can be used to find a block of PCI address space that nobody else is using, which can then be used for whatever the device and driver wish to use it for. The PCI infrastructure preallocates PCI address space regions based on the device configuration BASE registers at the time the bus is discovered. As a result this function is needed only to manage a device that does not completely declare its address space usage in its hardware configuration registers.

`pciio_piospace_free()` is used to release an allocation made previously by `pciio_piospace_alloc()`.

### Specifying PCI Address Spaces

The space parameter takes on of the following values:

`PCIIO_SPACE_WIN(n)`  
specifies one of the regions on the PCI bus decoded by the PCI card's BASE registers. The address specified is the offset within the decoded area, and the entire PIO region must fit within the decoded area.

`PCIIO_SPACE_CFG`  
requests a pointer handle that can be used to access the configuration space for the card, via the `pciio_config_get()` and `pciio_config_set()` functions documented in `pciio_config(D3)`.

Other space types are rarely needed but can be used:

`PCIIO_SPACE_IO`  
requests a mapping into somewhere in the PCI bus I/O address space.

`PCIIO_SPACE_MEM`  
requests a mapping into somewhere in the PCI bus Memory space. Since PCI bus address space is preallocated by the kernel, this is a dangerous function to use.

### PIO Attribute Flags

There are no useful values for the flags argument in this release. Specify the argument as 0. In IRIX 6.4 and onward, some usable flags are available.

### EXAMPLES

Here is a contrived example of how one might initialize a very strange PCI card. It is not clear that this would be the best way to do it, but it does give an example of the relationship between the various functions.

```
pcifoo_attach(vertex_hdl_t vhdl)
{
    unsigned *cfgspace;
    struct pcifoo_devregs *devregs;
```

```
pciio_piomap_t pmap;
pciio_piomap_t cmap;
struct pcifoo_chan_config *tune;
...
/* Get the configuration space base
 * pointer.
 */
cfgspace = pciio_piotrans_addr
    (vhd1, 0, PCIIIO_SPACE_CFG, 0, 256, 0);
if (cfgspace == NULL) {
    cmn_err(CE_ALERT,
        "pcifoo_attach: pciio_piotrans_addr failed");
    return -1;
}
/* Get a pointer we can use for PIO to our
 * device's control registers. This call
 * attempts to use fixed shared resources,
 * but will allocate unshared mapping resources
 * if required.
 */
devregs = pciio_pio_addr
    (vhd1, 0,
    PCIIIO_SPACE_WIN(0), 0,
    sizeof (struct pcifoo_devregs),
    &pmap, 0);
if (devregs == NULL) {
    cmn_err(CE_ALERT,
        "pcifoo_attach: pciio_pio_addr failed");
    return -1;
}
/* save cfgspace and devregs for use;
 * save pmap for pciio_dmamap_free
 * call if/when we are unregistered.
 */
...
/* pretend our "channel" space is too big
 * to successfully map with piotrans, so
 * we have to use piomap, and that it is
 * too big for us to get it in one call
 * to piomap_addr.
 */
cmap = pciio_piomap_alloc(vhd1, 0,
    PCIIIO_SPACE_WIN(2), 0, CHAN_SEP * CHANS,
    sizeof (struct pcifoo_chan_config), 0);
for (chan = 0; chan < chans; ++chan) {
```

```

        tune = (struct pcifoo_chan_config *)
                pciio_piomap_addr(cmap, CHAN_SEP * chan,
                                sizeof (struct pcifoo_chan_config));
        /* now fiddle with this particular channel */
        tune->chan = chan + 2;
        tune->volume = 5;
        tune->balance = 0;
        pciio_piomap_done(cmap);
    }
    pciio_piomap_free(cmap);
    ...
}

```

**NOTES**

It is not necessary to separately establish mappings for each individual PIO target register. It is customary and more efficient to use a single mapping to cover the entire register set of a device.

**SEE ALSO**

`pciio(D3)`, `pciio_config(D3)`, `pciio_dma(D3)`, `pciio_error(D3)`,  
`pciio_get(D3)`, `pciio_intr(D3)`.

**DIAGNOSTICS**

`pciio_piotrans_addr()` returns a null pointer when shared (fixed) resources can not be used to construct a valid physical address that maps to the desired range of PCI addresses.

`pciio_pio_addr()` returns a null pointer when the target PCI address can not be mapped either with shared (fixed) resources, or with unshared mapping resources. If this happens, and the object being mapped is large, it might be possible to set up mappings to smaller regions of the target space.

`pciio_piomap_alloc()` returns a null pointer when resources can not be allocated to establish PIO mappings to the described region, or if the function parameters are inconsistent.

`pciio_piomap_addr()` returns a null pointer when the specified target address can not be mapped using the specified PIO channel. This would usually be due to specifying a target block that is outside the previously specified target area or is larger than the previously specified maximum mapping size. It may also return a null pointer if the PIO channel is currently in use and has not been marked idle by a `pciio_piomap_done()` call.

---

## Glossary

### **ABI**

Application Binary Interface, a defined interface that includes an *API*, but adds the further promise that a compiled object file will be portable; no recompilation will be required to move to any supported platform.

### **API**

Application Programming Interface, a defined interface through which services can be obtained. A typical API is implemented as a set of callable functions and header files that define the data structures and specific values that the functions accept or return. The promise behind an API is that a program that compiles and works correctly will continue to compile and work correctly in any supported environment (however, recompilation may be required when porting or changing versions). See *ABI*.

### **big-endian**

The hardware design in which the most significant bits of a multi-byte integer are stored in the byte with the lowest address. Big-endian is the default storage order in MIPS processors. Opposed to *little-endian*.

### **block**

As a verb, to suspend execution of a process. See *sleep*.

### **block device**

A device such as magnetic tape or a disk drive, that naturally transfers data in blocks of fixed size. Opposed to *character device*.

### **block device driver**

Driver for a block device. A block device's driver is not allowed to support the *ioctl()*, *read()* or *write()* entry points, but does have a *strategy()* entry point. See *character device driver*.

**bus master**

An I/O device that is capable of generating a sequence of bus operations—usually a series of memory reads or writes—independently, once programmed by software. See *direct memory access*.

**bus-watching cache**

A *cache memory* that is aware of bus activity and, when the I/O system performs a DMA write into physical memory or another CPU in a multiprocessor system modifies *virtual memory*, automatically invalidates any copy of the same data then in the cache. This hardware function eliminates the need for explicit data cache write back or invalidation by software.

**cache coherency**

The problem of ensuring that all cached copies of data are true reflections of the data in memory. The usual solution is to ensure that, when one copy is changed, all other copies are automatically marked as invalid so that they will not be used.

**cache line**

The unit of data when data is loaded into a *cache memory*. Typically 128 bytes in current CPU models.

**cache memory**

High-speed memory closely attached to a CPU, containing a copy of the most recently used memory data. When the CPU's request for instructions or data can be satisfied from the cache, the CPU can run at full rated speed. In a multiprocessor or when DMA is allowed, a *bus-watching cache* is needed.

**character device**

A device such as a terminal or printer that transfers data as a stream of bytes, or a device that can be treated in this way under some circumstances. For example, a disk (normally a *block device*) can be treated as a character device for purposes of reading diagnostic information.

**character device driver**

The kernel-level device driver for a *character device* transfers data in bytes between the device and the user program. A *STREAMS driver* works with a character driver. Note that a *block device* such as magnetic tape or disk drives can also support character access through a character driver. Each disk device, for example, is represented as two different device special files, one managed by a *block device driver* and one by a character device driver.

**close**

Relinquish access to a resource. The user process invokes the **close()** system call when it is finished with a device, but the system does not necessarily execute your *drvclose()* entry point for that device.

**data structure**

Contiguous memory used to hold an ordered collection of fields of different types. Any *API* usually defines several data structures. The most common data structure in the *DDI/DKI* is the *buf\_t*.

**DDI/DKI**

Device Driver Interface/Device Kernel Interface; the formal *API* that defines the services provided to a device driver by the kernel, and the rules for using those services. *DDI/DKI* is the term used in the UNIX System V documentation. The IRIX version of the *DDI/DKI* is close to, but not perfectly compatible with, the System V interface.

**deadlock**

The condition in which two or more processes are blocked, each waiting for a lock held by the other. Deadlock is prevented by the rule that a driver upper-half entry point is not allowed to hold a lock while sleeping.

**devflag**

A public global flag word that characterizes the abilities of a device driver, including the flags *D\_MP*, *D\_WBACK* and *D\_OLD*.

**device driver**

A software module that manages access to a hardware device, taking the device in and out of service, setting hardware parameters, transmitting data between memory and the device, sometimes scheduling multiple uses of the device on behalf of multiple processes, and handling I/O errors.

**direct memory access**

When a device reads or writes in memory, asynchronously and without specific intervention by a CPU. In order to perform DMA, the device or its attachment must have some means of storing a memory address and incrementing it, usually through *mapping registers*. The device writes to physical memory and in so doing can invalidate *cache memory*; a *bus-watching cache* compensates.

**device number**

Each *device special file* is identified by a pair of numbers: the *major device number* identifies the device driver that manages the device, and the *minor device number* identifies the device to the driver.

**device special file**

A filename in the */dev* directory that represents a hardware device. A device special file does not specify data on disk, but rather identifies a particular hardware unit and the device driver that handles it. The *inode* of the file contains the *device number* as well as permissions and ownership data.

**downstream**

The direction of STREAMS messages flowing through a write queue from the user process to the driver.

**EISA bus**

Enhanced Industry Standard Architecture, a bus interface supported by certain Silicon Graphics systems.

**EISA Product Identifier (ID)**

The four-byte product identifier returned by an EISA expansion board.

**file handle**

An integer returned by the **open()** kernel function to represent the state of an open file. When the file handle is passed in subsequent kernel services, the kernel can retrieve information about the file, for example, when the file is a *device special file*, the file handle can be associated with the major and minor *device number*.

**gigabyte**

See kilobyte.

**GIO bus**

Graphics I/O bus, a bus interface used on Indigo, Indigo<sup>2</sup>, and Indy workstations.

**I/O operations**

Services that provide access to shared input/output devices and to the global data structures that describe their status. I/O operations open and close files and devices, read data from and write data to devices, set the state of devices, and read and write system data structures.

**inode**

The UNIX and IRIX disk object that represents the existence of a file. The inode records the owner and group IDs and permissions. For regular disk files, the inode distinguishes files from directories and has other data that can be set with `chmod`. For device special files, the inode contains the major and minor device numbers and distinguishes block from character files.

**inter-process communication**

System calls that allow a process to send information to another process. There are several ways of sending information to another process: signals, pipes, shared memory, message queues, semaphores, streams, or sockets.

**interrupt**

A hardware signal that causes a CPU to set aside normal processing and begin execution of an interrupt handler. An interrupt is parameterized by the type of bus and the *interrupt level*, and possibly with an *interrupt vector* number. The kernel uses this information to select the interrupt handler for that device.

**interrupt level**

A number that characterizes the source of an *interrupt*. The VME bus provides for seven interrupt levels. Other buses have different schemes.

**interrupt priority level**

The relative priority at which a bus or device requests that the CPU call an interrupt process. Interrupts at a higher level are taken first. The interrupt handler for an interrupt can only be preempted on its CPU by an interrupt handler for an interrupt of higher level.

**interrupt vector**

A number that characterizes the specific device that caused an *interrupt*. Most VME bus devices have a specific vector number set by hardware, but some can have their vector set by software.

**ioctl**

Control a character device. Character device drivers may include a "special function" entry point, `pfxioc()`.

**IRQ**

Interrupt Request Input, a hardware signal that initiates an interrupt.

**k0**

Virtual address range within the kernel address space that is cached but not mapped by translation look-aside buffers. Also referred to as kseg0.

**k1**

Virtual address range within the kernel address space that is neither cached nor mapped. Also called kseg1.

**k2**

Virtual address range within the kernel address space that can be both cached and mapped by translation look-aside buffers. Also called kseg2.

**kernel level**

The level of privilege at which code in the IRIX kernel runs. The kernel has a private address space, not acceptable to processes at *user-level*, and has sole access to physical memory.

**kilobyte (KB)**

1,024 bytes, a unit chosen because it is both an integer power of 2 ( $2^{10}$ ) and close to 1,000, the basic scale multiple of engineering quantities. Thus 1,024 KB,  $2^{20}$ , is 1 megabyte (MB) and close to  $1e6$ ; 1,024 MB,  $2^{30}$ , is 1 gigabyte (GB) and close to  $1e9$ ; 1,024 GB,  $2^{40}$ , is 1 terabyte (TB) and close to  $1e12$ . In the MIPS architecture using 32-bit addressing, the user segment spans 2 GB. Using 64-bit addressing, both the user segment and the range of physical addresses span 1 TB.

**kseg $n$**

See k0, k1, k2.

**little-endian**

The hardware design in which the least significant bits of a multi-byte integer are stored in the byte with the lowest address. Little-endian order is the normal order in Intel processors, and optional in MIPS processors. Opposed to *big-endian*. (These terms are from Swift's *Gulliver's Travels*, in which the citizens of Lilliput and Blefuscu are divided by the burning question of whether one's breakfast egg should be opened at the little or the big end.)

**lock**

A data object that represents the exclusive right to use a resource. A lock can be implemented as a *semaphore* (q.v.) with a count of 1, but because of the frequency of use of locks, they have been given distinct software support (see LOCK(D3)).

**major device number**

A number that specifies which device driver manages the device represented by a *device special file*. In IRIX 6.2, a major number has at most 9 bits of precision (0-511). Numbers 60-79 are used for OEM drivers. See also *minor device number*.

**map**

In general, to translate from one set of symbols to another. Particularly, translate one range of memory addresses to the addresses for the corresponding space in another system. The *virtual memory* hardware maps the process address space onto pages of physical memory. The *mapping registers* in a DMA device map bus addresses to physical memory corresponding to a buffer. The `mmap(2)` system call maps part of process address space onto the contents of a file.

**mapping registers**

Registers in a DMA device or its bus attachment that store the address translation data so that the device can access a buffer in physical memory.

**megabyte**

See kilobyte.

**minor device number**

A number that, encoded in a *device special file*, identifies a single hardware unit among the units managed by one device driver. Sometimes used to encode device management options as well. In IRIX 6.2, a minor number may have up to 18 bits of precision. See also major device number.

**mmapped device driver**

A driver that supports mapping hardware registers into process address space, permitting a user process to access device data as if it were in memory.

**module**

A STREAMS module consists of two related queue structures, one for upstream messages and one for downstream messages. One or more modules may be pushed onto a stream between the stream head and the driver, usually to implement and isolate a communication protocol or a line discipline.

**open**

Gain access to a device. The kernel calls the *pxopen()* entry when the user process issues an *open()* system call.

**page**

A block of virtual or physical memory, of a size set by the operating system and residing on a page-size address boundary. The page size is 4,096 ( $2^{12}$ ) bytes when in 32-bit mode; the page size in 64-bit mode can range from  $2^{12}$  to  $2^{20}$  at the operating system's choice (see the *getpagesize(2)* reference page).

**PIO**

Programmed I/O, meaning access to a VME device by mapping device registers into process address space, and transferring data by storing and loading single bytes or words.

**poll**

Poll entry point for a non-stream character driver. A character device driver may include a *drvpoll()* entry point so that users can use *select(2)* or *poll(2)* to poll the file descriptors opened on such devices.

**prefix**

Driver prefix. The name of the driver must be the first characters of its standard entry point names; the combined names are used to dynamically link the driver into the kernel. Specified in the *master.d* file for the driver. Throughout this manual, the prefix *px* represents the name of the device driver, as in *pxopen()*, *pxioctl()*.

**primary cache**

The *cache memory* most closely attached to the CPU execution unit, usually in the processor chip.

**primitives**

Fundamental operations from which more complex operations can be constructed.

**priority inheritance**

An implementation technique that prevents *priority inversion* when a process of lower priority holds a mutual exclusion *lock* and a process of higher priority is blocked waiting for the lock. The process holding the lock “inherits” or acquires the priority of the highest-priority waiting process in order to expedite its release of the lock. IRIX supports priority inheritance for mutual exclusion locks only.

**priority inversion**

The effect that occurs when a low-priority process holds a *lock* that a process of higher priority needs. The lower priority process runs and the higher priority process waits, inverting the intended priorities. *See* priority inheritance.

**process control**

System calls that allow a process to control its own execution. A process can allocate memory, lock itself in memory, set its scheduling priorities, wait for events, execute a new program, or create a new process.

**protocol stack**

A software subsystem that manages the flow of data on a communications channel according to the rules of a particular protocol, for example the TCP/IP protocol. Called a “stack” because it is typically designed as a hierarchy of layers, each supporting the one above and using the one below.

**pseudo-device**

Software that uses the facilities of the *DDI/DKI* to provide specialized access to data, without using any actual hardware device. Pseudo-devices can provide access to system data structures that are unavailable at the user-level. For example, the *fsctl* driver gives superuser access to filesystem data (see *fsctl(7)*) and the inode monitor pseudo-device allows access to file activity (see *imon(7)*).

**read**

Read data from a device. The kernel executes the *pfxread()* entry point whenever a user process calls the *read()* system call.

**scatter/gather**

An I/O operation in which what to the device is a contiguous range of data is distributed across multiple pages that may not be contiguous in physical memory. On input to memory, the device scatters the data into the different pages; on output, the device gathers data from the pages.

**SCSI**

Small Computer System Interface, the bus architecture commonly used to attach disk drives and other block devices.

**SCSI driver interface**

A collection of machine-independent input/output controls, functions, and data structures, that provides a standard interface for writing a SCSI driver.

**semaphore**

A data object that represents the right to use a limited resource, used for synchronization and communication between asynchronous processes. A semaphore contains a count that represents the quantity of available resource (typically 1). The **P** operation (mnemonic: dePlete) decrements the count and, if the count goes negative, causes the caller to wait (see `psema(D3X)`, `cpsema(D3X)`). The **V** operation (mnemonic: reVive) increments the count and releases any waiting process (see `vsema(D3X)`, `cvsema(D3X)`). *See also* lock.

**signals**

Software interrupts used to communicate between processes. Specific signal numbers can be handled or blocked. Device drivers sometimes use signals to report events to user processes. Device drivers that can wait have to be sensitive to the possibility that a signal could arrive.

**sleep**

Suspend process execution pending occurrence of an event. The term “block” is also used.

**socket**

A software structure that represents one endpoint in a two-way communications link. Created by `socket(2)`.

**spl**

Set priority level, a function that was formerly part of the *DDI/DKI*, and used to lock or allow interrupts on a processor. It is not possible to use `spl` effectively in a multiprocessor system, so it has been superceded by more sophisticated means of synchronization such as the *lock* and *semaphore*.

**strategy**

In general, the plan or policy for arbitrating between multiple, concurrent requests for the use of a device. Specifically in disk device drivers, the policy for scheduling multiple, concurrent disk block-read and block-write requests.

**STREAM**

A linked list of kernel data structures that provide a full-duplex data path between a user process and a device. Streams are supported by the STREAMS facilities in UNIX System V Release 3 and later.

**STREAM head**

The stream head, which is inserted by the STREAMS subsystem, processes STREAMS-related system calls and performs data transfers between user space and kernel space. It is the component of a stream closest to the user process. Every stream has a stream head.

**STREAMS**

A kernel subsystem used to build a stream, which is a modular, full-duplex data path between a device and a user process. In IRIX 5.x and later, the TCP/IP stack sits on top of the STREAMS stack. The Transport Layer Interface (TLI) is fully supported.

**STREAMS driver**

A software module that implements one stage of a STREAM. A STREAMS driver can be “pushed on” or “popped off” any *STREAM*.

**TCP/IP**

Transmission Control Protocol/Internet Protocol.

**terabyte**

See *kilobyte (KB)*.

**TFP**

The internal code name for the MIPS R8000 processor, used in some Silicon Graphics publications.

**TLI**

Transport Interface Layer.

**user-level**

The privilege level of the system at which user-initiated programs run. A user-level process can access the contents of one address space, and can access files and devices only by calling kernel functions. Contrast to *kernel level*.

**unmap**

Disconnect a memory-mapped device from user process space, breaking the association set by mapping it.

**VME bus**

VERSA Module Eurocard bus, a bus architecture supported by the Silicon Graphics Challenge and Onyx systems.

**VME-bus adapter**

A hardware conduit that translates host CPU operations to VME-bus operations and decodes some VME-bus operations to translate them to the host side.

**virtual memory**

Memory contents that appear to be in contiguous addresses, but are actually mapped to different physical memory locations by hardware action of the translation lookaside buffer (TLB) and page tables managed by the IRIX kernel. The kernel can exploit virtual memory to give each process its own address space, and to load many more processes than physical memory can support.

**virtual page number**

The most significant bits of a virtual address, which select a *page* of memory. The processor hardware looks for the VPN in the TLB; if the VPN is found, it is translated to a physical page address. If it is not found, the processor traps to an exception routine.

**volatile**

Subject to change. The volatile keyword informs the compiler that a variable could change value at any time (because it is mapped to a hardware register, or because it is shared with other, concurrent processes) and so should always be loaded before use.

**wakeup**

Resume suspended process execution.

**write**

Write data to a device. The kernel executes the *pfxread()* or *pfxwrite()* entry points whenever a user process calls the **read()** or **write()** system calls.



---

# Index

## Numbers

- 32-bit address space
  - See* address space, 32-bit
- 64-bit address space
  - See* address space, 64-bit
- 64-bit mode, 26

## A

- address exception, 9
- addressing, 3-29
- address space
  - 32-bit, 16-20
    - embedding in 64-bit, 22
    - kseg0, 19
    - kseg1, 20
    - kseg2, 19
    - kuseg, 18
    - segments of, 16
    - virtual mapping, 18
  - 64-bit, 20-25
    - cache-controlled, 24
    - segments of, 20-25
    - sign extension, 22
    - virtual mapping, 22
    - xkseg, 24
    - xksegs, 23
    - xkuseg, 23
- bus virtual, 12
- data transfer between, 194

- device address, 4
- kernel, 19, 24
  - map to user, 27
- locking in memory, 129
- memory address, 4
- physical, 4, 202
- supervisor, 23
- user process, 18, 23
- alternate console, 249
- ASSERT macro, 254
- audio not covered, 81
- Audio/Serial Option (ASO), 47
- authorized binary interface (ABI), 173

## B

- bdevswtable*, 141
- block device, 48
  - combined with character, 58, 151
  - driver must be MP-aware, 146
  - used when mounting filesystem, 279
  - versus character, 35
- buffer (*buf\_t*)
  - See* data types, *buf\_t*
- bus adapter
  - translates addresses, 12
- bus virtual address, 12

**C**

cache, 15-16  
  64-bit access, 24  
  alignment of buffers, 204  
  coherency, 15  
  control functions, 203  
  device access always uncached, 10  
  primary, 7  
  secondary, 7  
cache algorithm, 25  
*cdevsw* table, 141  
Challenge/Onyx  
  no uncached memory, 28  
character device, 48  
  combined with block, 58, 151  
  used with *mkfs*, 278  
  versus block, 36  
COFF file format not supported, 230  
command  
  *See* IRIX commands  
compiler options  
  32-bit, 232  
  64-bit, 233  
  for loadable driver, 240  
  for network driver, 345  
compiler variables, 231  
configuration files, 39-41  
  */dev/MAKEDEV*, 229, 243  
  */etc/inittab*, 249  
  */etc/rc2.d*, 38  
  */etc/rc2/S23autoconfig*, 241  
  */usr/cpu/sysgen/IPnnboot*, 235  
  */usr/lib/X11/input/config*, 41  
  */var/sysgen/boot*, 40, 234  
  */var/sysgen/Makefile.kernio*, 230  
  */var/sysgen/master.d*, 40, 228, 234, 235-238  
    dependencies, 236  
    example, 280  
    format, 235, 240

  stubs, 237  
  variables, 237  
  */var/sysgen/master.d/mem*, 28  
  */var/sysgen/mtune/\**, 41  
  */var/sysgen/system*, 41, 234  
  example, 280  
  */var/sysgen/system/irix.sm*, 64, 69  
  for debugging, 247  
  for SCSI, 304  
configuration flags, 236  
configuring a driver  
  loadable, 239-242  
  nonloadable, 234-239  
CPU, 5-16  
  device access, 10  
  IP26, 29  
  memory access by, 6  
  model number from inventory, 33  
  processors in, 5  
  type numbers, 5  
  watchpoint registers, 261

**D**

D\_MP flag, 145, 176  
D\_MT flag, 146  
D\_OLD flag, 146, 150  
D\_WBACK flag, 146  
Data Link Provider Interface (DLPI), 339  
data transfer, 194-197  
data types  
  summary table, 523  
  *buf\_t*, 158, 185-187  
  BP\_ISMAPPED, 187  
  displaying, 270  
  for synchronization, 178  
  functions, 202  
  interrupt handling, 169  
  management, 219

- data types (*continued*)
  - cred\_t*, 152, 205
  - dev\_t*, 37, 151, 182
  - struct dsconf*, 92
  - struct dsreq*, 85, 92
    - ds\_flags*, 87
    - ds\_msg*, 91
    - ds\_ret*, 89
    - ds\_status*, 91
  - edt\_t*, 149
  - iovec\_t*, 184
  - lock\_t*, 187, 208
  - major\_t*, 36, 182
  - minor\_t*, 37, 182
  - mrlock\_t*, 187, 214
  - mutex\_t*, 187, 211
  - struct pollhead*, 159
  - proc\_t* (*not available*), 206
  - struct scsi\_request*, 311, 316
  - struct scsi\_target\_info*, 308
  - sema\_t*, 187
  - sleep\_t*, 213
  - sv\_t*, 187, 222
  - uio\_t*, 155, 184, 197
  - \_\_userabi\_t*, 174
  - vhandl\_t*, 163, 199
- debugging kernel, 245-250
- device access, 10
- device address, 4
- device number
  - See* major device number, minor device number
- device special file, 34-39
  - as normal file, 35
  - defining, 229-230
  - /dev/dsk*, 38
  - /dev/ei*, 118, 128
  - /dev/kmem*, 27
  - /dev/mem*, 27
  - /dev/mmdev*, 28
  - /dev/scsi/\**, 82-85
  - /dev/tty\**, 47
  - /dev/vme/\**, 129
  - EISA mapping, 70
  - for user-level interrupt, 128
  - inode contents, 35
  - multiple names for, 39
  - name format, 38, 83
  - PCI mapping, 79
  - VME mapping, 65
- device special file */dev/kmem*, 28
- digital media not covered, 81
- Direct Memory Access (DMA), 11, 54-56
  - buffer alignment for, 204
  - cache control, 203
  - maximum size, 204
  - setting up, 201-204
  - user-level, 72-78
  - user-level SCSI, 89
- disk volume header, 246, 277
- driver
  - compiling, 231-233, 345
  - configuring, 234-242
  - debugging, 245-271
  - examples
    - network, 348-372
    - RAM drive, 273-296
    - SCSI bus, 318-322
  - flag constant, 145-147, 239, 500
  - initialization, 147-149
  - lower half, 56, 57
  - prefix, 140, 234
    - in *master.d*, 40
  - process context, 205
  - registration, 241
  - types of, *xxv*, 43-60
    - block, 48
    - character, 48
    - kernel-level, *xxv*, 28, 47-60
    - layered, 58
    - loadable, 59

## driver

types of (*continued*)

network, 337-372

process-level, xxv

pseudo-device, 54

SCSI bus, 317-318

STREAMS, xxv, 48

upper half, 56

in multiprocessor, 175, 176

user-level, 27, 43-47

*See also* entry points*See also* loadable driver

## driver debugging

alternate console, 249

breakpoints, 259

circular buffer output, 252

lock metering, 248

memory display, 262

multiprocessor, 255

*setsym* use, 249

stopping during bootstrap, 256

symbol lookup, 258

symbols, 247

symmon use, 254

system log output, 251

## driver operations, 48-56

DMA, 54

ioctl, 50

mmap, 53

open, 49

read, 51

write, 51

## dslib library, 94-107

function summary, 94

data transfer options, 89

**doscsireq()**, 97**ds\_ctostr()**, 99**ds\_vtostr()**, 99**dsfclose()**, 95**dsopen()**, 95**filldsreq()**, 97**fillg0cmd()**, 98**fillg1cmd()**, 98**inquiry12()**, 100**modeselect15()**, 100**modesense1a()**, 101**read08()**, 102**readcapacity25()**, 103**readextended28()**, 102**releaseunit17()**, 104**requestsense03()**, 104**reserveunit16()**, 104**senddiagnostic1d()**, 105**testunitready00()**, 106**write0a()**, 106**writeextended2a()**, 106

## dsreq driver, 82

data transfer options, 89

DS\_ABORT, 93

DS\_CONF, 92

DS\_RESET, 94

exclusive open, 85

flags, 87

return codes, 89

scatter/gather, 89

*struct dsconf*, 92*struct dsreq*, 85-92*ds\_flags*, 87*ds\_msg*, 91*ds\_ret*, 89*ds\_status*, 91

## E

## EISA bus

mapping into user process, 44

PIO bandwidth, 71

user-level PIO, 68-72

ELF object format, 230

entry points  
 summary table, 142, 522  
 close, 153-154, 163, 502  
 devflag, 145-147, 239  
 edtinit, 148, 241  
 halt, 171-172  
 info, 500  
 init, 148, 241, 501  
 interrupt, 167-170  
 ioctl, 154-155, 173  
 map, 163-165  
 mmap, 165  
 mversion, 239  
 open, 150-153, 501  
   mode flag, 152  
   type flag, 151  
 poll, 158-161  
   and interrupts, 169  
 print, 172  
 read, 155-157  
 size, 153, 172  
 start, 149, 241, 501  
 strategy, 157-158  
   and interrupts, 169  
   called from read or write, 156  
   design models, 219  
 unload, 163, 170-171, 242  
 unmap, 166  
 usage, 143  
 write, 155-157

example driver, 273, 318, 348

execution model, 173-174

external interrupt, 46, 118-123  
 generate, 118  
 input is level-triggered, 119  
 pulse widths, 120  
 set pulse widths, 121  
 user-level handler, 131

**F**

*fmodsw* table, 141

function  
 See IRIX functions, kernel functions

**H**

hardware inventory, 31-34  
 adding entries to, 34  
 contents, 32  
*hinv* displays, 32  
 network driver use, 344  
 software interface to, 33

header files  
 summary table, 188  
*dslib.h*, 94  
 for network drivers, 341  
*sgidefs.h*, 26  
*sys/cmnerr.h*, 251  
*sys/debug.h*, 254  
*sys/file.h*, 152  
*sys/immu.h*, 200  
*sys/major.h*, 36  
*sys/open.h*, 151  
*sys/param.h*, 185  
*sys/poll.h*, 159  
*sys/region.h*, 165  
*sys/scsi.h*, 303  
*sys/sem.h*, 187  
*sys/sysmacros.h*, 36, 37, 200  
*sys/types.h*, 26, 36, 37, 182  
*sys/uiio.h*, 184  
*sys/var.h*, 41

- I**
  - idbg debugger, 247-249, 264-271
    - command line use, 265
    - command syntax, 266-271
    - configuring in kernel, 248
    - display I/O status, 269
    - display process data, 267
    - interactive mode, 264
    - invoking, 264
    - loading, 264
    - lock meter display, 269
    - log file output, 265
    - memory display, 267
  - ide PROM monitor, 246
  - include file
    - See header files
  - INCLUDE statement, 148, 238, 241
  - initialization, 147-149
  - inode, 35, 49
  - interrupt, 57
    - and strategy entry point, 169
    - associating to a driver, 167
    - concurrent with processing, 175
    - enabled during initialization, 147
    - latency, 169
    - on multiprocessor, 168
    - See also user-level interrupt (ULI)
  - inventory
    - See hardware inventory
  - IP26 CPU, 29
  - IRIX commands
    - autoconfig*, 234, 238, 249
    - dvhtool*, 246-247
    - hinv*, 32
      - and MAKEDEV, 38, 229
      - for CPU type, 6
    - install*, 38, 229, 243
    - lboot*, 41
    - builds switch tables, 141
    - driver prefix with, 140
    - loads SCSI driver, 304
  - mkfs*, 278
  - mknod*, 38, 229, 243
  - ml*, 59, 243
  - mount*, 50, 150, 279
  - nvr*, 250
  - prtvtoc*, 277
  - setsym*, 249
  - system*, 41, 242, 252
    - max DMA size, 204
    - switch table size, 141
  - umount*, 153
  - uname*, 6
  - versions*, 246
- IRIX functions
    - close()**, 153
    - endinvent()**, 33
    - getinvent()**, 33
    - getpagesize()**, 23
    - ioctl()**, 45, 46, 51, 131
    - kmem\_alloc()**, 57
    - mmap()**, 27, 53, 162-163
      - EISA PIO, 70
      - PCI PIO, 79
      - VME PIO, 66
    - mpin()**, 129
    - munmap()**, 166
    - open()**, 49, 150
      - with dsreq driver, 85
    - plock()**, 129
    - poll()**, 159-160
    - read()**, 51, 54
    - setinvent()**, 33
    - syslog()**, 251
    - test\_and\_set()**, 133
    - ULI\_block\_intr()**, 132
    - ULI\_register\_ei()**, 131
    - ULI\_register\_vme()**, 131
    - ULI\_sleep()**, 128, 132

IRIX functions (*continued*)

**ULI\_wakeup()**, 132  
**write()**, 51, 54

## J

jag (SCSI-to-VME) adapter, 84  
jag (SCSI-to-VME adapter), 304

## K

kernel address space  
  driver runs in, 28  
  mapping to user space, 27  
kernel execution model, 173

kernel functions

  summary table, 525  
  **add\_to\_inventory()**, 34  
  **badaddr()**, 198  
  **bcopy()**, 196  
  **biodone()**, 169, 220  
  **bioerror()**, 169  
  **biowait()**, 220  
  **bp\_mapin()**, 203  
  **brelse()**, 193  
  **bzero()**, 196  
  **cmn\_err()**, 251-253  
    buffer output, 252  
    system log output, 251  
  **copyin()**, 154, 196  
  **copyout()**, 154, 196  
  **cvsema()**, 180  
  **dki\_dcache\_inval()**, 203  
  **dki\_dcache\_wb()**, 29, 203  
  **drv\_getparm()**, 205  
  **drv\_priv()**, 152, 205  
  **drvhztousec()**, 217  
  **drvusectohz()**, 217

**flushbus()**, 204  
**fubyte()**, 196  
**geteblk()**, 193  
**getemajor()**, 36, 183  
**getemisor()**, 183, 308  
**getinvent()**, 6  
**getrbuf()**, 193  
**initnsema()**, 180  
**initnsema\_mutex()** (not supported), 225  
**ip26\_enable\_ucmem()**, 29  
**ip26\_return\_ucmem()**, 29  
**itimeout()**, 160, 217  
**kern\_malloc()** (obsolete), 190  
**kmem\_alloc()**, 19, 190  
**kmem\_zalloc()**, 191  
**makedevice()**, 183  
**phalloc()**, 159, 192  
**phfree()**, 171, 192  
**physiock()**, 146, 156  
**pollwakeup()**, 159, 169  
**printf()**, 253  
**psema()**, 180, 225  
**ptob()**, 23  
**rmalloc()**, 194  
**rmallocmap()**, 194  
**rmfree()**, 194  
**sleep()**, 221  
**splhi()**  
  denigrated, 215  
  meaningless, 175  
**splnet()**  
  ineffective, 345  
**splvme()**  
  useless, 179  
**subyte()**, 196  
**timeout()**, 217  
**uiomove()**, 197  
**uiophysio()**, 156  
**untimeout()**, 217  
**userabi()**, 173  
**v\_getaddr()**, 199

kernel functions (*continued*)

- v\_gethandle()**, 200
  - v\_mapphys()**, 163, 199
  - vsema()**, 180, 225
  - vt\_gethandle()**, 165, 166
  - wakeup()**, 221
- kernel-level driver, xxv, 47-60, 139-180
- structure of, 140
- kernel mode of processor, 8
- kernel panic
- address exception, 9
  - moving data, 195
- kernel switch tables, 141

**L**

- layered driver, 58
- lboot*
- See* IRIX commands
- libc reentrant version, 127
- loadable driver, 59
- and switch table, 141
  - autoregister, 148
  - compiler options, 240
  - configuring, 239
  - initialization, 148
  - loading, 241
  - master.d, 240
  - mversion entry, 239
  - not in miniroot, 60
  - registration, 241
  - unloading, 242
- loading a driver, 241
- locking
- See* mutual exclusion
- locking memory, 129
- lock metering support, 248, 269
- lower half of driver, 57

**M**

- major device number, 36, 182
- available numbers, 36
  - block vs. character, 36
  - dynamic allocation, 243
  - external versus internal, 183
  - for STREAMS clone, 510, 511
  - in */dev/scsi*, 83
  - indexes switch table, 141
  - in inode, 35
  - in master.d, 40, 235
  - input to open, 49
  - in variables in master.d, 237
  - range of, 36
  - selecting, 228
- /dev/MAKEDEV*, 37-39, 82, 182, 229, 243
- adding to */dev/scsi*, 84
- master.d configuration files
- See* configuration files, */var/sysgen/master.d*
- memory, 3-29
- memory address
- cached, 19
  - physical, 4, 261
  - uncached, 20
- memory allocation, 189-194
- memory display, 262
- memory mapping, 27-28, 161-166
- miniroot
- no loadable drivers, 60
- minor device number, 37, 182
- encoding, 37
  - external versus internal, 183
  - for STREAMS clone driver, 510, 511
  - in */dev/scsi*, 83
  - in inode, 35
  - input to open, 49, 151
  - selecting, 229

multiprocessor
 

- block driver must support, 146
- converting to, 178
- driver design for, 174-180, 505
- driver flag `D_MP`, 145, 176
- drivers for, 59
- interrupt handling on, 168
- network drivers in, 344-347
- nonMP driver on CPU 0, 145, 176
- `splhi` useless in, 175
- synchronizing upper-half code, 176
- uniprocessor assumptions invalid, 174
- uniprocessor drivers use CPU 0, 59
- using `symmon` in, 255

 mutex locks, 210
   
 mutual exclusion, 207, 208-216
 

- basic locks, 208-209
- in multiprocessor drivers, 175
- in network driver, 346-347
- mutex locks, 210
- priority inheritance, 211
- reader/writer locks, 213
- semaphore, 225
- sleep locks, 212

  
**N**
  
 names of devices, 35, 38, 83
   
 network, 337
 

- based on 4.3BSD, 340
- driver interfaces, 340-344
- example driver, 348
- header files, 341
- multiprocessor considerations, 344
- overview, 338
- STREAMS protocol stack, 339

 network driver
 

- debugging, 271
- must be MP-aware, 146

 Network File System (NFS), 339

**P**

page size
 

- I/O, 200
- macros, 200
- memory, 23, 200

 PCI bus
 

- user-level PIO, 80

 pipe semantics, 507
   
 prefix, 40, 140, 234
   
 primary cache, 7
   
 priority inheritance, 211
   
 priority level functions, 215
   
 privilege checking, 205
   
 process, 205-206
 

- display data about, 267
- handle of, 206
- sending signal to, 206
- table of in kernel, 267

 process-level driver, xxv
   
 processor
 

- kernel mode, 8
- types, 5
- user mode, 8

 Programmed I/O (PIO), 10, 63
 

- EISA bus, 68-72
- PCI bus, 80
- VME bus, 64-68

 pseudo-device driver, 54
   
 putbuf circular buffer, 252, 267

**R**

RAM drive, 273
   
 raw device
 

- See* character device

 reader/writer locks, 213
   
 reentrant C library, 127
   
 registration of loadable driver, 241

**S**

*sash* standalone shell, 246

SCSI bus, 299-334

  adapter error codes, 326

  adapter number, 83, 307

  adapter type number, 306, 324

  command

    Inquiry, 100, 308

    Mode Select, 100

    Mode Sense, 101

    Read, 102

    Read Capacity, 103

    Request Sense, 104

    Reserve Unit, 104

    Send Diagnostic, 105

    Test Unit Ready, 106

    Write, 106

  display request structure, 269

  driver, 317-318

  error messages, 325-334

  example driver, 318-322

  hardware support overview, 300

  host adapter, 301

    functions of, 305

    initialization, 323

    number of, 303

    overview, 323

    purpose, 302

**scsi\_abort()**, 316

**scsi\_alloc()**, 309

**scsi\_command()**, 311

**scsi\_free()**, 310

**scsi\_info()**, 308

**scsi\_reset()**, 317

    vectors to, 306, 324

  kernel overview, 301

  LUN, 84, 304

  message string tables, 326

  sense codes, 327

  target ID, 83

  target number, 304

  user-level access, 45, 81-107

*See also* dsreq driver

secondary cache, 7

sector unit macros, 200

semaphore, 224-226

  for mutual exclusion, 225

  for waiting, 226

signal, 206

sign extension of 32-bit addresses, 22

SIGSEGV, 9

Silicon Graphics

  developer program, xxvii

  FTP server, xxvii

  WWW server, xxvii

64-bit address space

*See* address space, 64-bit

64-bit entries *see* Numbers

sleep locks, 212

socket interface, 339

STREAMS, 499-519

  function summary, 512

  clone driver, 510-511

  close entry point, 502

  debugging, 270

  display data structures, 270

  driver, xxv

  extended poll support, 507

*module\_info* structure, 500

  multiprocessor design, 505

  multithreaded monitor, 505

  open entry point, 501

  put functions, 502

  service scheduling, 508

  srv functions, 503

*streamtab* structure, 500

  supplied drivers, 508

STREAMS protocol stack, 339

structure of driver, 140

switch table, 141

symmon debugger, 246-247, 254-263  
 breakpoints, 259  
 command syntax, 257-258  
 how invoked, 255  
 in multiprocessor, 255  
 in uniprocessor, 255  
 invoking at bootstrap, 256  
 memory display, 262  
 prompt, 255  
 symbol lookup, 258  
 virtual memory commands, 261  
 watchpoint register use, 261

synchronization variable, 222

sysgen files  
*See* configuration files

system console  
 alternate, 249

system log display, 251

*systune*  
*See* IRIX commands

## T

terminal as console, 249

The, 271

32-bit entries *see* Numbers

tick, 217

time unit functions, 217

TLI interface, 339

Translate Lookaside Buffer (TLB), 7

Translation Lookaside Buffer (TLB), 8  
 maps kernel space, 24  
 maps kuseg, 19  
 number of entries in, 9

## U

udmalib, 72-78  
**dma\_allocbuf()**, 73  
**dma\_mkparms()**, 74  
**dma\_start()**, 75

uncached memory access  
 32-bit, 20  
 64-bit, 24  
 do not map, 164  
 IP26, 29  
 none in Challenge, 28

uniprocessor  
 converting driver, 178  
 using symmon, 255

unloading a driver, 242

upper half of driver, 56

upper half of of driver, 175, 176

user-level DMA, 72-78

user-level driver, 43-47

user-level interrupt (ULI), 46, 125-136  
 and debugging, 127  
 external interrupt with, 131  
 initializing, 128  
 interrupt handler function, 126-128  
 registration, 130  
 restrictions on handler, 126  
**ULI\_block\_intr()** function, 132  
**ULI\_register\_ei()** function, 131  
**ULI\_register\_vme()** function, 131  
**ULI\_sleep()** function, 128, 132  
**ULI\_wakeup()** function, 132  
 VME interrupt with, 131

user-level process, 43

user mode of processor, 8

USE statement, 148, 238

## V

- variables in master.d, 237
- VECTOR statement, 238, 241
  - edtinit entry point, 148
  - EISA PIO, 69
  - for SCSI host adapter, 304
  - use of `ctrl=`, 151
  - VME PIO, 64
- vfssw* table, 141
- virtual memory, 7, 8-10
  - 32-bit mapping, 18
  - 64-bit mapping, 22
  - debug display of, 261
  - page size, 23
- virtual page number (VPN)
  - 32-bit, 18

## VME bus

- adapter number, 84
- jag adapter, 84
- mapping into user process, 44
- PIO to
  - bandwidth, 68
  - user-level DMA, 45, 72-78
  - user-level DMA bandwidth, 75
  - user-level interrupt handler for, 131
  - user-level PIO, 64-68
- volatile keyword, 127
- volume header, 246, 277

## W

- waiting, 207, 216-224
  - for a general event, 221
  - for an interrupt, 219
  - for memory, 218
  - semaphore, 226
  - synchronization variables, 222
  - timed events, 217
  - time units, 217



---

## Tell Us About This Manual

As a user of Silicon Graphics products, you can help us to better understand your needs and to improve the quality of our documentation.

Any information that you provide will be useful. Here is a list of suggested topics:

- General impression of the document
- Omission of material that you expected to find
- Technical errors
- Relevance of the material to the job you had to do
- Quality of the printing and binding

Please send the title and part number of the document with your comments. The part number for this document is 007-3443-002.

Thank you!

## Three Ways to Reach Us

- To send your comments by **electronic mail**, use either of these addresses:
  - On the Internet: [techpubs@sgi.com](mailto:techpubs@sgi.com)
  - For UUCP mail (through any backbone site): *[your\_site]!sgi!techpubs*
- To **fax** your comments (or annotated copies of manual pages), use this fax number: 650-932-0801
- To send your comments by **traditional mail**, use this address:

Technical Publications  
Silicon Graphics, Inc.  
2011 North Shoreline Boulevard, M/S 535  
Mountain View, California 94043-1389