## **Digital UNIX**

## **Network Programmer's Guide**

Order Number: AA-PS2WD-TE

March 1996

Product Version: Digital UNIX Version 4.0 or higher

This manual describes the Digital UNIX network programming environment. It describes the sockets and STREAMS frameworks, including information about system calls, header files, libraries, and software bridges that allow sockets programs to use STREAMS drivers and STREAMS programs to use BSD-based drivers. Additionally, it describes how to write programs to the X/Open Transport Interface (XTI), as well as how to port sockets-based applications to XTI. It also describes the Digital UNIX eSNMP API.

Digital Equipment Corporation makes no representations that the use of its products in the manner described in this publication will not infringe on existing or future patent rights, nor do the descriptions contained in this publication imply the granting of licenses to make, use, or sell equipment or software in accordance with the description.

Possession, use, or copying of the software described in this publication is authorized only pursuant to a valid written license from Digital or an authorized sublicensor.

Portions of this document © Digital Equipment Corporation 1994, 1995, 1996. All rights reserved

Portions of this document are adapted from *A STREAMS-based Data Link Provider Interface* − *Version* 2 © 1991 UNIX International, Inc. Permission to use, copy, modify, and distribute this documentation for any purpose and without fee is hereby granted, provided that the above copyright notice appears in all copies and that both that copyright notice and this permission notice appear in supporting documentation, and that the name UNIX International not be used in advertising or publicity pertaining to distribution of the software without specific, written prior permission. UNIX International makes no representations about the suitability of this documentation for any purpose. It is provided "as is" without express or implied warranty.

The following are trademarks of Digital Equipment Corporation:

ALL–IN–1, Alpha AXP, AlphaGeneration, AlphaServer, AlphaStation, AXP, Bookreader, CDA, DDIS, DEC, DEC Ada, DEC Fortran, DEC FUSE, DECnet, DECstation, DECsystem, DECterm, DECUS, DECwindows, DTIF, LinkWorks, MASSBUS, MicroVAX, OpenVMS, POLYCENTER, Q–bus, StorageWorks, TruCluster, TURBOchannel, ULTRIX, ULTRIX Mail Connection, ULTRIX Worksystem Software, UNIBUS, VAX, VAXstation, VMS, XUI, and the DIGITAL logo.

AT&T is a registered trademark of American Telephone & Telegraph Co. BSD is a trademark of Uunet Technologies. IEEE is a registered trademark of the Institute of Electrical and Electronic Engineers, Inc. Intel is a trademark of Intel Corporation. Adobe, PostScript, and Display PostScript are registered trademarks of Adobe Systems, Inc. UNIX is a registered trademark in the United States and other countries licensed exclusively through X/Open Company Ltd. X/Open is a trademark of X/Open Company Ltd. Xerox is a registered trademark of Xerox Corporation.

All other trademarks and registered trademarks are the property of their respective holders.

## **Contents**

ADC	out This Manual			
Aud	ience	xvii		
New	ew and Changed Features			
Orga	anization	xviii		
Rela	ted Documents	xix		
Read	der's Comments	XX		
Con	ventions	xxi		
1	Introduction to the Network Programming Environment			
1.1	Data Link Interfaces	1–3		
1.2	Sockets and STREAMS Frameworks	1–3		
1.3	X/Open Transport Interface	1–5		
1.4	Extensible SNMP	1–7		
1.5	Sockets and STREAMS Interaction	1–7		
1.6	Putting It All Together	1–9		
2	Data Link Provider Interface			
2.1	Modes of Communication	2–3		
2.2	Types of Service	2–4		
	2.2.1 Local Managment Services	2-5		

	<ul> <li>2.2.2 Connection-Mode Services</li> <li>2.2.3 Connectionless-Mode Services</li> <li>2.2.4 Acknowledged Connectionless-Mode Data Transfer</li> </ul>	2–5 2–6 2–6
2.3	DLPI Addressing	2–7
2.4	DLPI Primitives	2–8
2.5	Identifying Available PPAs	2–11
3	X/Open Transport Interface	
3.1	Overview of XTI	3–2
3.2	XTI Features	3–4
	3.2.1 Modes of Service and Execution	3–4
	3.2.1.1 Connection-Oriented and Connectionless Service 3.2.1.2 Asynchronous and Synchronous Execution	3–4 3–5
	3.2.2 The XTI Library, TLI Library, and Header Files	3–6
	3.2.2.1 XTI and TLI Header Files	3–6 3–7
	3.2.3 Events and States	3–10
	3.2.3.1 XTI Events	3–10 3–13
	3.2.4 Tracking XTI Events	3–14
	3.2.4.1 Outgoing Events 3.2.4.2 Incoming Events	3–14 3–16
	3.2.5 A Map of XTI Functions, Events, and States	3–17 3–20
3.3	Using XTI	3–21
	<ul> <li>3.3.1 Guidelines for Sequencing Functions</li> <li>3.3.2 State Management by the Transport Provider</li> <li>3.3.3 Writing a Connection-Oriented Application</li> </ul>	3–21 3–23 3–24
	3.3.3.1 Initializing an Endpoint 3.3.3.4 Transferring Data 3.3.3.5 Releasing Connections	3–24 3–30 3–33

	3.3.3.6 Deinitializing Endpoints	3–36
	3.3.4 Writing a Connectionless Application	3–37
	3.3.4.1 Initializing an Endpoint 3.3.4.2 Transferring Data 3.3.4.3 Deinitializing Endpoints	3–37 3–37 3–39
3.4	Phase-Independent Functions	3–40
3.5	Porting to XTI	3–41
	3.5.1 Protocol Independence and Portability 3.5.2 XTI and TLI Compatibility 3.5.3 Rewriting a Socket Application to Use XTI	3–41 3–43 3–45
3.6	Differences Between XPG3 and XPG4	3–47
	3.6.1 Major Differences 3.6.2 Source Code Migration	3–48 3–48
	3.6.2.1 Use the Older Binaries of your Application	3–49 3–49 3–49
	3.6.3 Binary Compatibility 3.6.4 Packaging 3.6.5 Interoperability 3.6.6 Using XTI Options	3–49 3–50 3–50 3–50
	3.6.6.1 Using XTI Options in XPG4 General Information Format of Options Elements of Negotiation Multiple Options and Options Levels Illegal Options Initiating an Option Negotiation Responding to a Negotiation Proposal	3–51 3–51 3–52 3–53 3–53 3–54 3–56
	Retrieving Information About Options Privileged and Read-Only Options Option Management of a Transport Endpoint The Option Value T_UNSPEC The info Argument	3–57 3–59 3–60 3–62 3–62
	Portability Issues 3.6.6.2 Negotiating Protocol Options in XPG3	3–63 3–63

4. Sockets  4.1 Overview of the Sockets Framework  4.1.1 Communication Properties of Sockets  4.1.1.1 Socket Abstraction  4.1.1.2 Communication Domains  4.1.1.3 Socket Types  4.1.1.4 Socket Names  4.2 Application Interface to Sockets  4.2.1 Modes of Communication  4.2.1.1 Connection-Oriented Communication  4.2.1.2 Connectionless Communication  4.2.2 Client/Server Paradigm  4.2.3 System Calls, Library Calls, Header Files, and Data Structures  4.2.3.1 Socket System Calls  4.2.3.2 Socket Library Calls  4.2.3.3 Header Files  4.2.3.4 Socket Related Data Structures  4.3 Using Sockets  4.3.1 Creating Sockets	3–65
4.1 Overview of the Sockets Framework  4.1.1 Communication Properties of Sockets  4.1.1.1 Socket Abstraction 4.1.1.2 Communication Domains 4.1.1.3 Socket Types 4.1.1.4 Socket Names  4.2 Application Interface to Sockets  4.2.1 Modes of Communication  4.2.1.1 Connection-Oriented Communication 4.2.1.2 Connectionless Communication  4.2.2 Client/Server Paradigm  4.2.3 System Calls, Library Calls, Header Files, and Data Structures  4.2.3.1 Socket System Calls 4.2.3.2 Socket Library Calls 4.2.3.3 Header Files 4.2.3.4 Socket Related Data Structures  4.3 Using Sockets  4.3 Using Sockets  4.3.1 Creating Sockets	
4.1.1 Communication Properties of Sockets  4.1.1.1 Socket Abstraction 4.1.1.2 Communication Domains 4.1.1.3 Socket Types 4.1.1.4 Socket Names  4.2 Application Interface to Sockets  4.2.1 Modes of Communication  4.2.1.1 Connection-Oriented Communication 4.2.1.2 Connectionless Communication  4.2.2 Client/Server Paradigm 4.2.3 System Calls, Library Calls, Header Files, and Data Structures  4.2.3.1 Socket System Calls 4.2.3.2 Socket Library Calls 4.2.3.3 Header Files 4.2.3.4 Socket Related Data Structures  4.3 Using Sockets  4.3.1 Creating Sockets	
4.1.1.1 Socket Abstraction 4.1.1.2 Communication Domains 4.1.1.3 Socket Types 4.1.1.4 Socket Names  4.2 Application Interface to Sockets 4.2.1 Modes of Communication 4.2.1.1 Connection-Oriented Communication 4.2.1.2 Connectionless Communication 4.2.3 System Calls, Library Calls, Header Files, and Data Structures  4.2.3.1 Socket System Calls 4.2.3.2 Socket Library Calls 4.2.3.3 Header Files 4.2.3.4 Socket Related Data Structures  4.3 Using Sockets 4.3.1 Creating Sockets	4–2
4.1.1.2 Communication Domains 4.1.1.3 Socket Types 4.1.1.4 Socket Names  4.2 Application Interface to Sockets 4.2.1 Modes of Communication 4.2.1.1 Connection-Oriented Communication 4.2.1.2 Connectionless Communication 4.2.3 System Calls, Library Calls, Header Files, and Data Structures  4.2.3.1 Socket System Calls 4.2.3.2 Socket Library Calls 4.2.3.3 Header Files 4.2.3.4 Socket Related Data Structures  4.3.5 Using Sockets 4.3.6 Creating Sockets	4–3
4.2.1 Modes of Communication  4.2.1.1 Connection-Oriented Communication  4.2.1.2 Connectionless Communication  4.2.2 Client/Server Paradigm  4.2.3 System Calls, Library Calls, Header Files, and Data Structures  4.2.3.1 Socket System Calls  4.2.3.2 Socket Library Calls  4.2.3.3 Header Files  4.2.3.4 Socket Related Data Structures  4.3 Using Sockets  4.3.1 Creating Sockets	4–3 4–3 4–5 4–6
4.2.1.1 Connection-Oriented Communication 4.2.1.2 Connectionless Communication 4.2.2 Client/Server Paradigm 4.2.3 System Calls, Library Calls, Header Files, and Data Structures  4.2.3.1 Socket System Calls 4.2.3.2 Socket Library Calls 4.2.3.3 Header Files 4.2.3.4 Socket Related Data Structures  4.3 Using Sockets 4.3.1 Creating Sockets	4–6
4.2.1.2 Connectionless Communication  4.2.2 Client/Server Paradigm  4.2.3 System Calls, Library Calls, Header Files, and Data Structures  4.2.3.1 Socket System Calls  4.2.3.2 Socket Library Calls  4.2.3.3 Header Files  4.2.3.4 Socket Related Data Structures  4.3 Using Sockets  4.3.1 Creating Sockets	4–7
4.2.3 System Calls, Library Calls, Header Files, and Data Structures  4.2.3.1 Socket System Calls 4.2.3.2 Socket Library Calls 4.2.3.3 Header Files 4.2.3.4 Socket Related Data Structures  4.3 Using Sockets 4.3.1 Creating Sockets	4–7 4–8
4.2.3.2 Socket Library Calls 4.2.3.3 Header Files 4.2.3.4 Socket Related Data Structures  4.3 Using Sockets 4.3.1 Creating Sockets	4–8 4–9
4.3.1 Creating Sockets	4–9 4–10 4–15 4–16
	4–18
	4–19
4.3.1.1 Setting Modes of Execution	4–21
<ul> <li>4.3.3 Establishing Connections</li> <li>4.3.4 Accepting Connections</li> <li>4.3.5 Setting and Getting Socket Options</li> </ul>	4–22 4–23 4–25 4–27 4–29
4.3.6.2 Using the write System Call 4.3.6.3 Using the send, sendto, recv and recvfrom System	4–29 4–30 4–30

	4.3.6.4 Using the sendmsg and recvmsg System Calls	4–33
	4.3.7 Shutting Down Sockets 4.3.8 Closing Sockets	4–36 4–36
4.4	BSD Socket Interface	4–37
	<ul><li>4.4.1 Variable-Length Network Addresses</li><li>4.4.2 Receiving Protocol Data with User Data</li></ul>	4–38 4–38
4.5	Common Socket Errors	4–40
4.6	Advanced Topics	4-41
	<ul><li>4.6.1 Selecting Specific Protocols</li><li>4.6.2 Binding Names and Addresses</li></ul>	4–42 4–42
	4.6.2.1 Binding to the Wildcard Address 4.6.2.2 Binding in the UNIX Domain	4–42 4–44
	4.6.3 Out-of-Band Data 4.6.4 Internet Protocol Multicasting	4–44 4–46
	4.6.4.1 Sending IP Multicast Datagrams 4.6.4.2 Receiving IP Multicast Datagrams	4–47 4–49
	4.6.5 Broadcasting and Determining Network Configuration 4.6.6 The inetd Daemon 4.6.7 Input/Output Multiplexing 4.6.8 Interrupt Driven Socket I/O 4.6.9 Signals and Process Groups 4.6.10 Pseudoterminals	4–51 4–54 4–55 4–58 4–58 4–59
5	Digital UNIX STREAMS	
5.1	Overview of the STREAMS Framework	5–1
	5.1.1 A Review of STREAMS Components	5–2
5.2	Application Interface to STREAMS	5–5
	<ul><li>5.2.1 Header Files and Data Types</li><li>5.2.2 STREAMS Functions</li></ul>	5–5 5–6
	5.2.2.1 The open Function	5–6 5–7

	5.2.2.3 The read Function	5-8
	5.2.2.4 The write Function	5-8
	5.2.2.5 The ioctl Function	5–9
	5.2.2.6 The mkfifo Function	5–9
	5.2.2.7 The pipe Function	5-10
	5.2.2.8 The putmsg and putpmsg Functions	5–11
	5.2.2.9 The getmsg and getpmsg Functions	5–11
	5.2.2.10 The poll Function	5–12
	5.2.2.11 The isastream Function	5–13
	5.2.2.12 The fattach Function	5–14
	5.2.2.13 The fdetach Function	5–14
5.3	Kernel Level Functions	5–17
	5.3.1 Module Data Structures	5-17
	5.3.2 Message Data Structures	5–18
	5.3.3 STREAMS Processing Routines for Drivers and Modules	5-20
	5.3.3.1 open and close Processing	5-20
	5.3.3.2 Configuration Processing	5-21
	5.3.3.3 Read Side Put and Write Side Put Processing	5-22
	5.3.3.4 Read Side Service and Write Side Service Processing .	5–22
	5.3.4 Digital UNIX STREAMS Concepts	5–23
	5.3.4.1 Synchronization	5-23
	5.3.4.2 Timeout	5–25
5.4	Configuring a User-Written STREAMS-Based Module or Driver in	
	the Digital UNIX Kernel	5–25
5.5	Device Special Files	5-29
5.6	Error and Event Logging	5–31
6	Extensible SNMP Application Programming Interface	
6.1	Overview of eSNMP	6–2
	6.1.1 Components of eSNMP	6–2
	6.1.2 Architecture	6–3
	6.1.2 Architecture 6.1.3 SNMP Versions:	6–4
		0-4
6.2	Overview of the Extensible SNMP Application Programming	
	Interface	6–4

	6.2.1 6.2.2	Subtree Object	Tables	6–6 6–8
		.2.2.1	The subtree_tbl.h File The subtree_tbl.c File	6–8 6–10
	6.2.3 6.2.4		nenting a Subagent. ent Protocol Operations	6–12 6–15
	_	.2.4.1	Order of Operations	6–16 6–16
6.3	Extens	sible SN	MP Application Programming Interface	6–18
	6.3.1	Calling	g Interface	6–18
		.3.1.1	The esnmp_init Routine	6–19
		.3.1.2	The esnmp_register Routine	6–20
		.3.1.3	The esnmp_unregister Routine	6–22
		.3.1.4	The esnmp_poll Routine	6–23
		.3.1.5	The esnmp_are_you_there Routine	6–23
	6.	.3.1.6	The esnmp_trap Routine	6–24
	6.	.3.1.7	The esnmp_term Routine	6–25
	6.	.3.1.8	The esnmp_sysuptime Routine	6–25
	6.3.2	Metho	d Routine Calling Interface	6–26
	6.	.3.2.1	The *_get Routine	6–26
	6.	.3.2.2	The *_set Method Routine	6–28
	O	verall P	Processing of the *_set Routine	6-31
		.3.2.3	Method Routines	6–33
	Value I	Represe	ntation	6–35
	6.3.3	The lib	osnmp Support Routines	6–37
	6	.3.3.1	The o_integer Routine	6–38
	6	.3.3.2	The o_octet Routine	6–39
	6	.3.3.3	The o_oid Routine	6–40
	6	.3.3.4	The o_string Routine	6-41
	6	.3.3.5	The str2oid Routine	6-43
	6	.3.3.6	The sprintoid Routine	6–44
	6	.3.3.7	The instance2oid Routine	6–44
	6	.3.3.8	The oid2instance Routine	6–45
	6.	.3.3.9	The inst2ip Routine	6–46
		.3.3.10	The cmp_oid Routine	6–49
		3.3.11	The cmp_oid_prefix Routine	6-50

	6.3.3.12 The clone_old Routine 6.3.3.13 The free_old Routine 6.3.3.14 The clone_buf Routine 6.3.3.15 The mem2oct Routine	6–50 6–51 6–52 6–52
	6.3.3.16 The cmp_oct Routine	6–53 6–53
	6.3.3.18 The free_oct Routine	6–54
	6.3.3.19 The free_varbind_data Routine	6–55
	6.3.3.20 The set_debug_level Routine	6-55
	6.3.3.21 The is_debug_level Routine	6–56
	6.3.3.22 The ESNMP_LOG Routine	6–57
7	Digital UNIX STREAMS/Sockets Coexistence	
7.1	Bridging STREAMS Drivers to Sockets Protocol Stacks	7–2
	7.1.1 The STREAMS Driver	7–3
	7.1.1.1 Using the ifnet STREAMS Module	7–4 7–10
	7.1.1.2 Data Link Provider Interface Primitives	/-10
7.2	Bridging BSD Drivers to STREAMS Protocol Stacks	7–11
	<ul> <li>7.2.1 Supported DLPI Primitives and Media Types</li> <li>7.2.2 Using the STREAMS Pseudodriver</li> </ul>	7–12 7–12
Α	Sample STREAMS Module	
В	Socket and XTI Programming Examples	
B.1	Connection-Oriented Programs	B-2
	B.1.1 Socket Server Program	B-2
	B.1.2 Socket Client Program	B-6
	B.1.3 XTI Server Program	B-9
	B.1.4 XTI Client Program	B-14
B.2	Connectionless Programs	B-17
	B.2.1 Socket Server Program	B-17
	B.2.2 Socket Client Program	B-20

	B.2.3 XTI Server Program B.2.4 XTI Client Program	B-23 B-27
B.3	Common Code	B-30
	B.3.1 The common.h Header File	B-31
	B.3.2 The server.h Header File	B - 32
	B.3.3 The serverauth.c File	B-33
	B.3.4 The serverdb.c File	B-36
	B.3.5 The xtierror.c File	B-38
	B.3.6 The client.h Header File	B-39 B-39
	B.3.8 The clientdb.c File	B-39 B-41
	B.5.6 The chemids. The	D-41
С	TCP Specific Programming Information	
C.1	TCP Throughput and Window Size	C-1
C.2	Programming the TCP Socket Buffer Sizes	C-1
C.3	TCP Window Scale Option	C-2
	C.3.1 Increasing the Socket Buffer Size Limit	C-2
D	Information for Token Ring Driver Developers	
D.1	Enabling Source Routing	D-1
D.2	Using Canonical Addresses	D-2
D.3	Avoiding Unaligned Access	D-3
D.4	Setting Fields in the softc Structure of the Driver	D-4
E	The Data Link Interface	
E.1	Prerequisites for DLI Programming	E-1
E.2	DLI Overview	E-2
	E.2.1 DLI Services	E-3 E-3

	E.2.3 E.2.4		DLI to Access the Local Area Network ding Higher-Level Services	E-3 E-4
E.3	The DLI Socket Address Data Structure			
	E.3.1 E.3.2 E.3.3	How	ard Frame Formats the sockaddr_dl Structure Works Ethernet Substructure	E-4 E-6 E-8
	_	2.3.3.1	How Ethernet Frames Work Defining Ethernet Substructure Values	E-8 E-9
	E.3.4	The 8	02.2 Substructure	E-11
	E	2.3.4.1	Defining 802 Substructure Values	E-11
E.4	Writin	ng DLI	Programs	E-15
	E.4.1 E.4.2 E.4.3 E.4.4 E.4.5 E.4.6	Using Creat Settin Bindi	ying Data Link Services g Digital UNIX System Calls ing a Socket g Socket Options ng the Socket g in the sockaddr_dl Structure	E-16 E-16 E-17 E-18 E-19 E-19
	E	2.4.6.1 2.4.6.2 2.4.6.3	Specifying the Address Family Specifying the I/O Device ID Specifying the Substructure Type	E-20 E-20 E-20
	E.4.7 E.4.8 E.4.9	Trans	lating the Buffer Size ferring Data ivating the Socket	E-22 E-22 E-23
E.5	DLI P	rogram	ming Examples	E-23
	E.5.1 E.5.2 E.5.3 E.5.4 E.5.5	Samp Samp Samp	le DLI Client Program Using Ethernet Format Packets . le DLI Server Program Using Ethernet Format Packets . le DLI Client Program Using 802.3 Format Packets le DLI Server Program Using 802.3 Format Packets le DLI Program Using getsockopt and setsockopt	E-24 E-27 E-31 E-36 E-41

## Glossary

#### Index

## Examples

5-1: Sample Module	5–26
B-1: Connection-Oriented Socket Server Program	B-2
B-2: Connection-Oriented Socket Client Program	B-6
B-3: Connection-Oriented XTI Server Program	B-9
B-4: Connection-Oriented XTI Client Program	B-14
B-5: Connectionless Socket Server Program	B-17
B-6: Connectionless Socket Client Program	B-20
B-7: Connectionless XTI Server Program	B-23
B-8: Connectionless XTI Client Program	B-27
B-9: The common.h Header File	B-31
B-10: The server.h Header File	B-32
B-11: The serverauth.c File	B-33
B-12: The serverdb.c File	B-36
B-13: The xtierror.c File	B-38
B-14: The client.h File	B-39
B-15: The clientauth.c File	B-39
B-16: The clientdb.c File	B-41
E-1: Filling the sockaddr_dl structure for Ethernet	E-21
E-2: Filling the sockaddr_dl structure for 802.2	E-21

## **Figures**

1-1:	Sockets and STREAMS Frameworks	1–4
1-2:	XTI, STREAMS, and Sockets Interactions	1–6
1-3:	Bridging STREAMS Drivers to Sockets Protocol Stacks	1–8
1-4:	Bridging BSD Drivers to STREAMS Protocol Stacks	1–9
1-5:	The Network Programming Environment	1–10
2-1:	DLPI Interface	2–2
2-2:	DLPI Service Interface	2–3
2-3:	Identifying Components of a DLPI Address	2–7
3-1:	X/Open Transport Interface	3–2
3-2:	A Transport Endpoint	3–3
3-3:	State Transitions for Connection-Oriented Transport Services	3–22
3-4:	State Transitions for the Connectionless Transport Service	3–23
4-1:	The Sockets Framework	4–2
4-2:	4.3BSD and 4.4BSD sockaddr Structures	4–38
4-3:	4.3BSD, 4.4BSD, and XPG4 msghdr Structures	4–40
5-1:	The STREAMS Framework	5–2
5-2:	Example of a Stream	5–3
7-1:	The ifnet STREAMS module	7–3
7-2:	DLPI STREAMS Pseudodriver	7–11
D-1:	Typical Frame	D-4
E-1:	DLI and the Digital UNIX Network Programming Environment	E-2
E-2:	The Ethernet Frame Format	E-4
E-3:	The 802.3 Frame Format	E-5
E-4:	The FDDI Frame Format	E-5
E 5.	The 200 2 Structures	E 6

## **Tables**

1-1: Components of the Network Programming Environment	1-1
2-1: DLPI Primitives Supported in Digital UNIX	2–8
3-1: Header Files for XTI and TLI	3–7
3-2: XTI Library Calls	3–8
3-3: Asynchronous XTI Events	3–10
3-4: Asynchronous Events and Consuming Functions	3–12
3-5: XTI Functions that Return TLOOK	3–12
3-6: XTI States	3–13
3-7: Outgoing XTI Events	3–15
3-8: Incoming XTI Events	3–16
3-9: State Transitions for Initialization of Connection-Oriented or Connectionless Transport Services	3–17
3-10: State Transitions for Connectionless Transport Services	3–18
3-11: State Transitions for Connection-Oriented Transport Services: Part 1.	3–18
3-12: State Transitions for Connection-Oriented Transport Services: Part 2 .	3–19
3-13: Phase-Independent Functions	3–40
3-14: Comparison of XTI and Socket Functions	3–45
3-15: Comparison of Socket and XTI Messages	3–47
4-1: Characteristics of the UNIX and Internet Communication Domains	4–4
4-2: Socket System Calls	4–9
4-3: Socket Library Calls	4–14
4-4: Header Files for the Socket Interface	4–16
4-5: Common Errors and Diagnostics	4–40
5-1: STREAMS Reference Pages	5–15
E-1: Calling Sequence for DLI Programs	E-16
E-2: Data Transfer System Calls Used With DLI	E-22

## **About This Manual**

This manual explains how to write programs with the X/Open Transport Interface (XTI) calls, STREAMS I/O calls, and the Berkeley Software Distribution (BSD) socket calls. For XTI and sockets, it provides conceptual and programming information. Additionally, it explains how to port applications from Transport Layer Interface (TLI) to XTI and from sockets to XTI. For STREAMS, this manual explains any differences between the Digital UNIX® It also provides information on the Digital UNIX Extensible System Network Management Protocol (eSNMP) application programming interface.

After reading this manual, you should be able to:

- Understand the programming support provided in Digital UNIX for networking
- Write an XTI application by using either connection-oriented or connectionless service
- Understand the Digital UNIX implementation of STREAMS
- Write a socket application
- Understand the differences between TLI and XTI and between sockets and XTI
- Write an eSNMP application

#### Audience

This manual addresses experienced UNIX programmers. We assume you are familiar with the following:

- C language
- Programming interfaces for UNIX operating systems
- Basic networking concepts, including an understanding of the Open Systems Interconnection (OSI) 7-layer model
- Efforts required to write networking applications

## **New and Changed Features**

This revision of the *Network Programmer's Guide* contains the following changes:

- A new section has been added to Chapter 3 that provides information on the XPG4 version of XTI. This section includes the major differences between the XPG3 and XPG4 versions, migration, and maintaining existing XTI programs.
- Chapter 4 has been revised. The chapter now provides information on the XPG4 compliant version of the sockets programming interface. It also includes a section that compares the BSD, XPG3, and XPG4 interfaces.
- Chapter 6 has been added to provide information on the Extensible Simple Network Management Protocol (eSNMP) application programming interface.
- Chapter 7 has been added to provide information on the ifnet STREAMS module.

## Organization

This manual is organized as follows:

	8
Chapter 1	Provides an overview of XTI, STREAMS, sockets, and the programming tasks required for network applications.
Chapter 2	Describes the dlb pseudodriver, which implements a subset of the the Data Link Provider Interface (DLPI).
Chapter 3	Describes the fundamental concepts associated with XTI, how to write connection-oriented and connectionless applications, compatibility issues with TLI, and how to port applications to XTI. XTI errors are also covered in this chapter.
Chapter 4	Describes the concepts associated with the 4.3BSD socket interface, and how to write socket applications.
Chapter 5	Describes Digital UNIX's STREAMS implementation.
Chapter 6	Describes Digital UNIX's Extensible System Network Management Application Programming Interface.
Chapter 7	Describes the ifnet STREAMS module and dlb STREAMS pseudodriver communication bridges.
Appendix A	Provides a sample STREAMS module.
Appendix B	Provides XTI and sockets programming examples.
Appendix C	Provides Transport Protocol Control (TCP) specific programming information.
Appendix D	Provides information required by token ring driver developers.

Appendix E Describes the Data Link Interface (DLI) and provides programming examples.

This guide also contains a glossary of terms and an index.

#### **Related Documents**

For general information about programming with Digital UNIX, refer to the *Programmer's Guide*.

For additional information about XTI, refer to the following manuals:

- X/Open Portability Guide Volume 7: Networking Services, ISBN 0-13-685892-9
- Application Environment Specification (AES) Operating System Programming Interfaces Volume, ISBN 0-13-043522-8, published by Prentice-Hall, includes all of the mandatory XTI calls
- X/Open CAE Specification: Networking Services, Issue 4 ISBN 1-85912-049-0

For additional information about the STREAMS I/O framework, refer to the following manuals:

- *Programmer's Guide: STREAMS*. Englewood Cliffs:Prentice-Hall, Inc., 1990.
  - This manual explains how to write applications, modules, and device drivers with STREAMS.
- AT&T System V Release 4 Programmer's Reference Manual. Englewood Cliffs:Prentice-Hall, Inc., 1989.
  - This manual contains the reference pages for all programming interfaces, including those for STREAMS.
- *IAT&T System V Release 4 System Administrator's Reference Manual.* Englewood Cliffs:Prentice-Hall, Inc., 1989.
  - This manual contains the reference pages for STREAMS ioctl commands.
- *Transport Provider Interface (TPI) Specification*, yet to be published by AT&T.

For additional information about the 4.3BSD socket interface, refer to the following books:

- Internetworking with TCP/IP: Principles, Protocols, and Architecture. Englewood Cliffs:Prentice-Hall, Inc., 1988.
  - This book, by Douglas Comer, includes a chapter that describes the socket interface.
- Design and Implementation of the 4.3BSD UNIX Operating System. Reading: Addison-Wesley Publishing Company, 1989.
  - This book, by Leffler, McKusick, Karels, and Quarterman, includes information about the purpose and use of sockets.

For information about administering networking interfaces, refer to the *System Administration* guide and the *Network Administration* guide.

The printed version of the Digital UNIX documentation set is color coded to help specific audiences quickly find the books that meet their needs. (You can order the printed documentation from Digital.) This color coding is reinforced with the use of an icon on the spines of books. The following list describes this convention:

Audience	lcon	Color Code
General users	G	Blue
System and network administrators	S	Red
Programmers	P	Purple
Device driver writers	D	Orange
Reference page users	R	Green

Some books in the documentation set help meet the needs of several audiences. For example, the information in some system books is also used by programmers. Keep this in mind when searching for information on specific topics.

The *Documentation Overview, Glossary, and Master Index* provides information on all of the books in the Digital UNIX documentation set.

#### **Reader's Comments**

Digital welcomes any comments and suggestions you have on this and other Digital UNIX manuals.

You can send your comments in the following ways:

- Fax: 603-881-0120 Attn: UEG Publications, ZK03-3/Y32
- Internet electronic mail: readers\_comment@zk3.dec.com

A Reader's Comment form is located on line in the following location:

/usr/doc/readers\_comment.txt

• Mail:

Digital Equipment Corporation UEG Publications Manager ZK03-3/Y32 110 Spit Brook Road Nashua, NH 03062-9987

A Reader's Comment form is located in the back of each printed manual. The form is postage paid if you mail it in the United States.

Please include the following information along with your comments:

- The full title of the book and the order number. (The order number is printed on the title page of this book and on its back cover.)
- The section numbers and page numbers of the information on which you are commenting.
- The version of Digital UNIX that you are using.
- If known, the type of processor that is running the Digital UNIX software.

The Digital UNIX Publications group cannot respond to system problems or technical support inquiries. Please address technical questions to your local system vendor or to the appropriate Digital technical support office. Information provided with the software media explains how to send problem reports to Digital.

#### **Conventions**

This document uses the following typographic conventions:

% \$	A percent sign represents the C shell system prompt. A dollar sign represents the system prompt for the Bourne and Korn shells.
#	A number sign represents the superuser prompt.
% cat	Boldface type in interactive examples indicates typed user input.
file	Italic (slanted) type indicates variable values, placeholders, and function argument names.

[ ]	In syntax definitions, brackets indicate items that are optional and braces indicate items that are required. Vertical bars separating items inside brackets or braces indicate that you choose one item from among those listed.
	In syntax definitions, a horizontal ellipsis indicates that the preceding item can be repeated one or more times.
cat(1)	A cross-reference to a reference page includes the appropriate section number in parentheses. For example, cat(1) indicates that you can find information on the cat command in Section 1 of the reference pages.
Return	In an example, a key name enclosed in a box indicates that you press that key.
Ctrl/x	This symbol indicates that you hold down the first named key while pressing the key or mouse button that follows the slash. In examples, this key combination is enclosed in a box (for example, Ctrl/C).

# **Introduction to the Network Programming Environment**

1

The network programming environment includes the programming interfaces for application, kernel, and driver developers writing network applications and implementing network protocols. Additionally, it includes the kernel-level resources that an application requires to process and transmit data, some of which include libraries, data structures, header files, and transport protocols.

This chapter introduces Digital UNIX's network programming environment by focussing on how the data link and application programming interfaces work to get data from an application in user space, through the network layers in kernel space, out onto the network, and back again.

Information about the kernel resources that support the interfaces is included in later chapters in this book. Individual chapters describe the particular system and library calls, data structures, and other programming considerations for each interface.

The primary components of the network programming environment are summarized in Table 1-1.

Table 1-1: Components of the Network Programming Environment

Component	Interface	Description
Data Link Interfaces	Data Link Interface (DLI)	Allows programs written to DLI on the ULTRIX operating system to use DLI on the Digital UNIX operating system to access the data link layer. Digital UNIX provides DLI for backward compatibility with ULTRIX. See Appendix E.

Table 1-1: (continued)

Component	Interface	Description
	dlb interface	Kernel-level interface targeted for STREAMS protocol modules that either use or provide data link services. The dlb STREAMS pseudodriver implements a subset of the Data Link Provider Interface (DLPI). See Chapter 2 and the Data Link Provider Specification (dlpi.ps) located in the /usr/share/doclib/dlpi directory. Note that the OSFPGMK200 subset must be installed to access the DLPI specification online.
Application Programming Interfaces	Sockets	The de facto industry standard programming interface. Digital UNIX implements the 4.3BSD socket interface as its default. You can use a special option to access the 4.4BSD interface.
		The Internet Protocol Suite, which consists of TCP, UDP, IP, ARP, ICMP, and SLIP is implemented over sockets. See RFC 1200: <i>IAB Protocol Standards</i> and Chapter 4.
	STREAMS	A kernel mechanism that supports the implementation of device drivers and networking protocol stacks. The STREAMS framework defines interface standards for character input and output within the kernel as well as between the kernel and user levels. The Digital UNIX operating system provides an AT&T, System V Release 4.0 compatible version of STREAMS. See Chapter 5.
	XTI/TLI	A protocol independent, transport layer application interface that consists of a series of functions. XTI is based on the Transport Layer Interface (TLI) and the transport service definition for the Open Systems Interconnection (OSI) model. See Chapter 3 and the X/Open Portability Guide Networking Services, Part 7.
	eSNMP	A set of routines that enables you to extend the SNMP agent process by creating MIBs.

Table 1-1: (continued)

Component	Interface	Description
Communication Bridges Between STREAMS and Sockets	ifnet STREAMS module	Allows STREAMS-based network device drivers to access the sockets-based TCP/IP protocol stack provided on Digital UNIX. See Chapter 7.
	dlb pseudodriver	Allows applications that use STREAMS-based protocol stacks to access BSD-based drivers. The dlb pseudodriver implements a subset of the DLPI specification. See Chapter 7.

It is easiest to understand the network programming environment by examining each component. The following sections introduce the environment piece by piece, starting with the components closest to the network and working up.

#### 1.1 Data Link Interfaces

The Digital UNIX network programming environment supports both the Data Link Interface (DLI) and the Data Link Provider Interface (DLPI). DLI enables you to port programs that run on ULTRIX systems to Digital UNIX systems. See Appendix E for information about DLI.

DLPI is a kernel-level interface that maps to the data link layer of the OSI reference model. DLPI frees its users from specific knowledge of the characteristics of the data link provider, allowing them to be implemented independently of a specific communications medium. Chapter 2 describes in greater detail DLPI, Digital UNIX's dlb pseudodriver, and the supported primitives.

#### 1.2 Sockets and STREAMS Frameworks

The Digital UNIX operating system supports AT&T's System V Release 4 STREAMS and BSD sockets frameworks for writing networking applications and for doing kernel-level network input/output (I/O). A framework comprises a particular programming interface and the kernel-level resources that the system requires to transmit and receive data.

Sockets is the de facto industry standard interface for writing networking applications. The sockets framework is BSD-based, consisting of a series of system and library calls, header files, and data structures. Applications can access kernel-resident networking protocols, such as the Internet Protocol suite, through socket system calls. Applications can also use socket library

calls to manipulate network information; for example, mapping service names to service numbers or translating the byte order of incoming data to that appropriate for the local system's architecture.

The STREAMS framework provides an alternative to sockets. The STREAMS interface was developed by AT&T and consists of system calls, kernel routines, and kernel utilities that are used to implement everything from networking protocol suites to device drivers. Applications in user space access the kernel portions of the STREAMS framework using system calls such as open, close, putmsg, getmsg, and ioctl. Figure 1-1 illustrates the STREAMS and sockets frameworks.

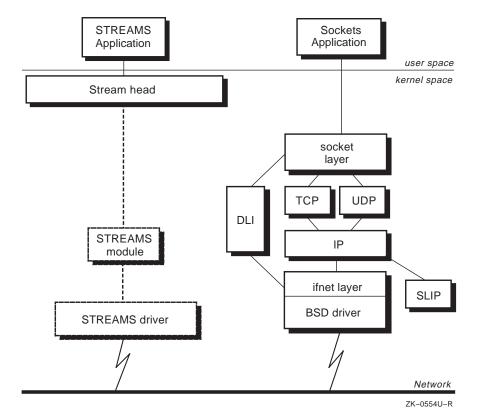


Figure 1-1: Sockets and STREAMS Frameworks

The dotted lines in Figure 1-1 indicate a supported configuration but not one that is provided by Digital UNIX.

With sockets, the application in user space passes data to the appropriate socket system calls, which then pass it to the network layer. Finally, the network layer passes it, via the ifnet layer, to the BSD driver, which puts it on to the network.

With STREAMS, the application in user space passes data to the Stream head, which passes it to any STREAMS modules that have been pushed on the Stream to process it. Each module passes the data to the next module until it finally reaches the STREAMS driver, which puts it out on to the network.

#### Note

Digital UNIX does not provide any STREAMS-based transport providers.

### 1.3 X/Open Transport Interface

The X/Open Transport Interface (XTI) defines a transport layer application interface that is independent of any transport provider. This means that programs written to XTI can be run over a variety of transport providers, such as the Transmission Control Protocol (TCP) or the User Datagram Protocol (UDP). The application specifies which transport provider to use.

Because XTI provides an interface that is independent of a transport provider, application developers are encouraged to write programs to XTI instead of STREAMS or sockets. Figure 1-2 illustrates the interaction between XTI and the STREAMS and sockets frameworks.

Sockets XTI/TLI Application user space kernel space Stream head timod socket xtiso layer TCP UDP DLI STREAMS ifnet layer module SLIP driver BSD driver STREAMS driver Network

Figure 1-2: XTI, STREAMS, and Sockets Interactions

Depending on the transport provider specified by the application, data can flow along one of two paths:

 If a STREAMS-based transport provider is specified, data follows the same route that it did for an application written to run over STREAMS. It passes first through the Stream head, then to any modules that the application pushed onto the Stream, and finally to the STREAMS driver, which puts it on to the network.

#### Note

Digital UNIX does not provide any STREAMS-based transport providers.

2. If a socket-based transport provider is specified (TCP or UDP), data is passed through timod and xtiso. The appropriate socket layer routines are called and the data is passed through the Internet protocols

and ifnet layer to the BSD-based driver, which puts it on to the network.

#### 1.4 Extensible SNMP

The Digital UNIX SNMP agent provides a framework for extensibility (called eSNMP). The SNMP daemon functions as an extensible masteragent, communicating with various subagents via the eSNMP protocol. The master agent implements the SNMP on behalf of the entire system, while subagents provide the actual MIB instrumentation. The eSNMP subagent development tools and API provide the mechanism for users to develop subagents that communicate with the master-agent and extend the MIB view on the Digital UNIX system.

#### 1.5 Sockets and STREAMS Interaction

Digital UNIX provides the ifnet STREAMS module to allow programs using Digital UNIX's BSD-based TCP/IP to access STREAMS-based drivers. It provides the dlb pseudodriver to allow programs using a STREAMS-based protocol stack to access BSD-based drivers provided on Digital UNIX.

Figure 1-3 illustrates an application using the BSD-based TCP/IP provided on Digital UNIX and accessing a STREAMS-based driver.

STREAMS Sockets Application Application user space kernel space Stream head socket layer TCP UDP DLI ΙP ifnet STREAMS ifnet layer **STREAMS** SLIP module driver module BSD driver **DLPI** Interface STREAMS driver Network ZK-0556U-R

Figure 1-3: Bridging STREAMS Drivers to Sockets Protocol Stacks

appropriate sockets system calls and is processed by the Internet protocols. Then the BSD ifnet layer of the networking subsystem, whose function is to map BSD ifnet messages to DLPI, passes the data to the ifnet STREAMS module. The ifnet STREAMS module processes it so that the STREAMS driver can put it on to the network. When information for the sockets-based application is returned, the STREAMS driver picks it up off of the network and passes it to the DLPI interface of the ifnet STREAMS module translates

In Figure 1-3, data travels from a sockets-based application through the

DLPI messages to BSD ifnet and passes it back to the BSD ifnet layer. The data is then processed by the Internet protocols and passed back to the application.

Figure 1-4 illustrates an application using a STREAMS-based protocol stack and accessing a BSD-based driver.

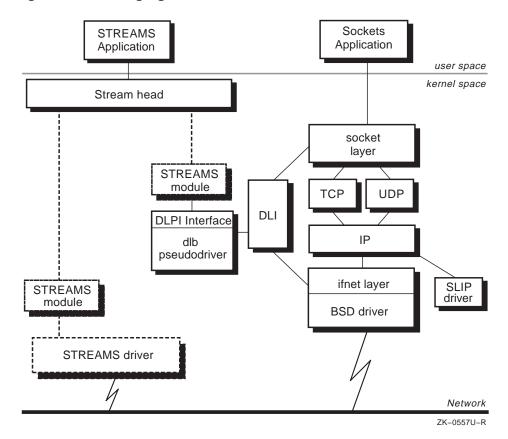


Figure 1-4: Bridging BSD Drivers to STREAMS Protocol Stacks

In Figure 1-4, data travels from a STREAMS-based application through the Stream head and is processed by whatever Streams modules have been pushed onto the stack. Instead of finally being passed to a STREAMS driver, the data is passed to the dlb STREAMS pseudodriver and is then forwarded to the ifnet layer of the sockets framework. From there it is further processed by a BSD driver and put on to the network.

## 1.6 Putting It All Together

Figure 1-5 represents the entire network programming environment. Variations of this figure appear in each chapter to give you perspective on the information being presented.

STREAMS Sockets Application XTI/TLI Application user space kernel space Stream head timod socket layer xtiso i STREAMS module TCP UDP DLI **DLPI** Interface dlb ΙP pseudodriver ifnet STREAMS STREAMS SLIP driver ifnet layer module module BSD driver **DLPI** Interface STREAMS driver Network ZK-0547U-R

Figure 1-5: The Network Programming Environment

## Data Link Provider Interface 2

Digital UNIX provides the dlb STREAMS pseudodriver, which is a partial implementation of the Data Link Provider Interface (DLPI).

This chapter describes the dlb STREAMS pseudodriver and the basics of DLPI. A PostScript copy of the DLPI specification (dlpi.ps) is located in the /usr/share/doclib/dlpi directory.

#### Note

You must have the OSFPGMK200 subset installed to access the DLPI specification online.

Figure 2-1 highlights the data link interfaces and shows their relationship to the rest of the network programming environment.

**STREAMS** Sockets XTI/TLI Application Application user space kernel space Stream head timod socket layer xtiso **STREAMS** module TCP UDP DLPI Interface DLI ΙP pseudodriver ifnet ifnet laver STREAMS SLIP **STREAMS** module module driver BSD driver DLPI Interface STREAMS driver ZK-0677U-R

Figure 2-1: DLPI Interface

#### Note

The dlb STREAMS pseudodriver supports a subset of DLPI primitives. See Section 2.4 for a list of the supported primitives.

The data link interface is the boundary between the network and data link layers of the OSI reference model. A network application, or data link service user (DLS user), uses the services of the data link interface. A driver, pseudodriver, or data link service provider (DLS provider), provides the services to the data link layer.

DLPI specifies a STREAMS kernel-level service interface that maps to the OSI reference model. It defines an interface to the services of the data link layer and provides a definition of service primitives that make up the interface.

Figure 2-2 shows the components of DLPI. The DLS user communicates with the DLS provider using request/response primitives; the DLS provider communicates with the DLS user with indication/confirmation primitives.

Request/Response Primitives

Data Link Service User

DLPI Interface

--Network Interface

Provider

Data Link Service

Primitives

Figure 2-2: DLPI Service Interface

The primitives that Digital UNIX supports are listed in Section 2.4.

#### 2.1 Modes of Communication

DLPI supports three modes of communication:

#### Connection

Enables a DLS user to establish a data link connection, transfer data over that connection, reset the link, and release the connection when the conversation has terminated.

The connection service establishes a data link connection between a local DLS user and a remote DLS user for the purpose of sending data. Only one data link connection is allowed on each Stream.

#### Connectionless

Enables a DLS user to transfer units of data to peer DLS users without incurring the overhead of establishing and releasing a connection. The connectionless service does not, however, guarantee reliable delivery of data units between peer DLS users (for instance, lack of flow control may cause buffer resource shortages that result in data being discarded).

Once a Stream has been initialized using the local management services, it may be used to send and receive connectionless data units.

ZK-0731U-R

#### Note

Digital UNIX supports only the connectionless mode of communication.

Acknowledged connectionless

Designed for general use for the reliable transfer of information between peer DLS users. These services are intended for applications that require acknowledgement of cross-LAN data unit transfer, but wish to avoid the complexity associated with the connection-mode services. Although the exchange service is connectionless, in-sequence delivery is guaranteed for data sent by the initiating station.

## 2.2 Types of Service

This section describes the types of service, or phases of communication, supported by DLPI. Note that the types of service available depend on the mode of communication (connection, connectionless, acknowledged connectionless) between the DLS provider and the DLS user.

DLPI supports the following types of service:

- Local management services
  - Information reporting service
  - Attach service
  - Bind service
- Connection-mode services
  - Connection establishment
  - Data transfer
  - Connection release
  - Reset service
- Connectionless-mode services
  - Connectionless data transfer
  - Quality of Service (QOS) management
  - Error reporting
- Acknowledged connectionless-mode services
  - Acknowledged connectionless-mode data transfer
  - Quality of service (QOS) management
  - Error reporting

## 2.2.1 Local Managment Services

The local management services apply to all three modes of communication supported by DLPI. They enable a DLS user to initialize a Stream that is connected to a DLS provider and to establish an identity with that provider. The local management services support the following:

Information reporting service
 Provides information about the DLPI Stream to the DLS user.

Attach service

Assigns a physical point of attachment (PPA) to a Stream. See Section 2.3 for more information.

Bind service

Associates a data link service access point (DLSAP) with a Stream.

## 2.2.2 Connection-Mode Services

The connection-mode services allow two DLS users to establish a data link connection between them to exchange data, and to reset the link and release the connection when the conversation is through. The connection-mode services support the following:

• Connection establishment service

Establishes a data link connection between a local DLS user and a remote DLS user for the purposes of sending data.

• Data transfer service

Provides for the exchange of user data in either direction or both directions simultaneously. Data is sent in logical groups called data link service data units (DLSDUs) and is guaranteed to be delivered in the order in which it was sent.

• Connection release service

Enables either the DLS user or DLS provider to break an established connection.

Reset service

Allows a DLS user to resynchronize the use of a data link connection, or a DLS provider to report detected loss of data unrecoverable within the data link service.

#### 2.2.3 Connectionless-Mode Services

The connectionless-mode services allow DLS users to exchange data without incurring the overhead of establishing and releasing a connection. The connectionless-mode services support the following:

- Connectionless data transfer service
  - Provides for the exchange of user data (DLSDU) in either direction or in both directions simultaneously.
- Quality of service (QOS) management service
   Enables a DLS user to specify the quality of service it can expect for each invocation of the connectionless data transfer service.
- Error reporting service

Provides a means to notify a DLS user that a previously sent data unit either produced an error or could not be delivered. However, the error reporting service does not guarantee that an error indication will be issued for every undeliverable data unit.

## 2.2.4 Acknowledged Connectionless-Mode Data Transfer

The acknowledged connectionless-mode data transfer services are designed for general use for the reliable transfer of data between peer DLS users. These services are intended for applications that require acknowledgment of data transfer between local area networks, but wish to avoid using the connection mode services. In-sequence delivery is guaranteed for data sent by the initiating station. The following services are supported:

- Acknowledged connectionless-mode data transfer service
   Provides for the exchange of DLSDUs which are acknowledged at the LLC sublayer.
- Quality of service (QOS) management service
   Enables a DLS user to specify the quality of service it can expect for each invocation of the connectionless data transfer service.
- Error reporting service

Provides a means to notify a DLS user that a previously sent data unit either produced an error or could not be delivered. However, the error reporting service does not guarantee that an error indication will be issued for every undeliverable data unit.

# 2.3 DLPI Addressing

Each DLPI user must establish an identity to communicate with other data link users. This identity consists of the following pieces of information:

- Physical attachment identification
  - This identifies the physical medium over which the DLS user communicates. The importance of identifying the physical medium is particularly evident on systems that are attached to multiple physical media. See Section 2.5 for information about identifying the available physical points of attachment (PPAs) on your system.
- Data link user identification
   The DLS user must register with the DLS provider so that the provider

can deliver protocol data units destined for that user.

The format of the DLSAP address is an unsigned character array containing the Medium Access Control (MAC) addresses followed by the bound Service Access Point (SAP). The SAP is usually two bytes in the case of Ethernet, or one byte in the case of ISO 8802-2 (IEEE 802.2). The one exception is when a HIERACHICAL DL\_SUBS\_BIND\_REQ is processed. In that case, the DLSAP address consists of the MAC address, the SNAP SAP (0xAA), and a five-byte SNAP.

Figure 2-3 illustrates the components of this identification approach.

DLS Users

DLSAP

DLS Provider

Physical Media

Figure 2-3: Identifying Components of a DLPI Address

ZK-0678U-R

The PPA is the point at which a system attaches itself to a physical communications medium. All communication on that physical medium funnels through the PPA. On systems where a DLS provider supports more

than one physical medium, the DLS user must identify the medium through which it will communicate. A PPA is identified by a unique PPA identifier.

DLPI defines the following two styles of DLS provider, which are distinguished by the way they enable a DLS user to choose a particular PPA:

- The style 1 provider assigns a PPA based on the major/minor device the DLS user opened. A style 1 driver can be implemented so that it reserves a major device for each PPA the data link driver would support.
  - This implementation of a style 1 driver allows the STREAMS clone open feature to be used for each PPA configured. Style 1 providers are appropriate when few PPAs are supported.
- The style 2 provider requires a DLS user to identify a PPA explicitly, using a special attach service primitive. For a style 2 driver, the open system call creates a Stream between the DLS user and DLS provider. Then, the attach primitive associates a particular PPA with that Stream. The format of the PPA identifier is specific to the DLS provider.
  - Digital UNIX supports only the style 2 provider because it is more suitable for supporting large numbers of PPAs.

## 2.4 DLPI Primitives

Table 2-1 lists and describes the DLPI primitives that Digital UNIX supports in the dlb STREAMS pseudodriver. For a complete list of DLPI primitives see the DLPI specification in the /usr/share/doclib/dlpi/dlpi.ps file.

Table 2-1: DLPI Primitives Supported in Digital UNIX

Primitive	Description
DL_ATTACH_REQ	Requests that the DLS provider associate a physical point of attachment (PPA) with a Stream. Used on style 2 providers only.
DL_BIND_REQ	Requests that the DLS provider bind a DLSAP to the Stream. The DLS user must identify the address of the DLSAP to be bound to the Stream.
DL_BIND_ACK	Reports the successful bind of a DLSAP to a Stream, and returns the bound DLSAP address to the DLS user. Generated in response to a DL_BIND_REQ.

Table 2-1: (continued)

Primitive	Description
DL_UNBIND_REQ	Requests that the DLS provider unbind the DLSAP that was bound by a previous DL_BIND_REQ from this Stream.
DL_DETTACH_REQ	Requests the DLS provider disassociate a physical point of attachment (PPA) with a stream.
DL_DISABMULTI_REQ	Request the DLS provider disable the multicast address.
DL_ENABMULTI_REQ	Request the DLS provider enable a specific multicast address. (The current implementation of the DLB driver requires the state to be DL_IDLE.)
DL_ERROR_ACK	Informs DLS user of a previously issued request which was invalid.
DL_INFO_ACK	Response to DL_INFO_REQ primitive; conveys information about the DLPI stream.
DL_INFO_REQ	Requests the DLS provider return information about the DLPI stream.
DL_OK_ACK	Acknowledges to the DLS user that a previously issued request primitive was successfully received.
DL_PHYS_ADDR_REQ	Requests that the DLS provider return either the default (factory) or current value of the physical address associated with the Stream, depending upon the value of the address type selected in the request.
DL_PHYS_ADDR_ACK	Returns the value for the physical address to the link user in response to a DL_PHYS_ADDR_REQ.
DL_SUBS_BIND_ACK	Is the positive response to a DL_SUBS_BIND_REQ from the DLS provider.

Table 2-1: (continued)

Primitive	Description
DL_SUBS_BIND_REQ	Requests the DLS provider bind a subsequent DLSAP to stream. There are two classes of subsequent bind requests: HIERACHICAL and PEER.
	HIERACHICAL requests are only valid for SNAPs (see the IEEE 802.1 specification) and you must have bound to the SNAP sap (0xAA) with a DL_BINDS_REQ before issuing the DL_SUBS_BIND_REQ for the SNAP.
	The PEER request binds to additional saps but does not change the DLSAP address of the stream.
DL_SUBS_UNBIND_REQ	Requests the DLS provider to unbind a sap which was previously bound by a DL_SUBS_BIND_REQ.
DL_TEST_CON	Conveys that a DLSDU TEST response was received in response to a DL_TEST_REQ.
DL_TEST_IND	Conveys to the DLS user that a TEST cmd DLSDU was received.
DL_TEST_REQ	Requests the DLS provider to transmit a TEST cmd DLSDU on behalf of the DLS user.
DL_TEST_RES	Requests the DLS provider to send a TEST response command on behalf of the DLS user.
DL_UDERROR_IND	Informs DLS user that a previously sent DL_UNITDATA_REQ failed.
DL_UNITDATA_REQ	Conveys one DLSDU from the DLS user to the DLS provider for transmission to a peer DLS user.
DL_UNITDATA_IND	Conveys one DLSDU from the DLS provider to the DLS user.
DL_XID_CON	Conveys that a XID DLSDU was received in response to a DL_XID_REQ.
DL_XID_IND	Conveys to the DLS user that a XID DLSDU was received.
DL_XID_REQ	Requests the DLS provider to transmit a XID DLSDU on behalf of the DLS user.

## Table 2-1: (continued)

# Primitive Description

DL\_XID\_RES

Requests the DLS provider to send a XID
DLSDU on behalf of the DLS user. This is in repsonse to a DL\_XID\_RES.

# 2.5 Identifying Available PPAs

When compiled and run as root on a Digital UNIX system, the following program opens the STREAMS device /dev/streams/dlb and prints to the screen the PPAs available on the system. The PPA number should be passed in using the dl\_ppa field of the DL\_ATTACH\_REQ DLPI primitive.

```
#include <sys/ioctl.h>
#include <stropts.h>
#include <errno.h>
#include <fcntl.h>
\#define ND\_GET ('N' << 8 + 0)
#define BUFSIZE 256
main()
{
        int i;
        int fd;
        char buf [BUFSIZE];
        struct strioctl stri;
        fd = open("/dev/streams/dlb", O_RDWR, 0);
        if (fd < 0) {
                perror("open");
                 exit(1);
        sprintf(buf, "dl_ifnames");
        stri.ic_cmd = ND_GET;
        stri.ic_timout = -1;
        stri.ic_len = BUFSIZE;
        stri.ic_dp = buf;
        if (ioctl(fd, I_STR, &stri) < 0) {</pre>
                 perror("ioctl");
                 exit(1);
        }
```

# X/Open Transport Interface 3

The X/Open Transport Interface (XTI) is a transport layer application interface that consists of a series of functions designed to be independent of the specific transport provider used. In the Digital UNIX operating system XTI is implemented according to the XPG3 and XPG4 specifications. XPG4 is the default. (XPG3 is provided for backward compatibility and is available by using a compiler switch.) For more information about XPG3 and XPG4, see the *X/Open Portability Guide Volume 7: Networking Services*. The Digital UNIX implementation of XTI is also thread safe.

Although similar in concept to the Berkeley socket interface, XTI is based on the AT&T Transport Layer Interface (TLI). TLI, in turn, is based on the **transport service** definition for the Open Systems Interconnection (OSI) model.

#### Note

Digital UNIX includes the Transport Control Protocol (TCP) and User Datagram Protocol (UDP) transport providers. Although the information provided in this chapter applies to all transport providers that Digital UNIX XTI supports, such as DECnet/OSI, the examples are specific to TCP or UDP. For more specific information using XTI over TCP and UDP, see the <code>xti\_internet(7)</code> reference page. For examples and information specific to other transport providers, see the documentation that accompanies their software.

This chapter contains the following information:

- Overview of XTI
- Description of XTI features
- Instructions on how to use XTI
- Instructions on how to port applications to XTI
- Information on the differences between XPG3 and XPG4.
- Explanation of XTI errors and error messages
- Information on configuring transport providers.

Figure 3-1 highlights XTI and its relationship to the Digital UNIX implementation of the Internet Protocol suite. It also shows how XTI and the Internet Protocol suite fit into the rest of the network programming environment.

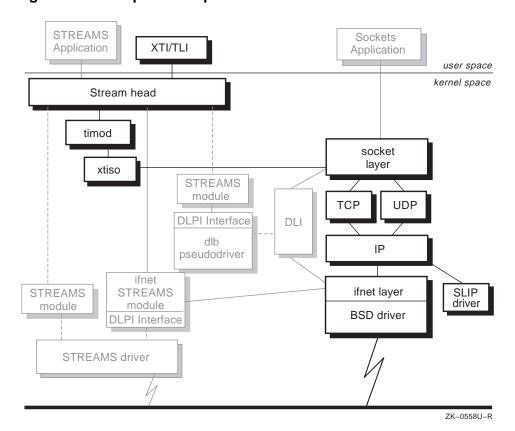


Figure 3-1: X/Open Transport Interface

## 3.1 Overview of XTI

XTI involves the interaction of the following entities:

Transport providers

A **transport provider** is a transport protocol, such as TCP or UDP, that offers transport layer services.

Transport users

A **transport user** is an application program that requires the services of a transport provider to send data to or receive data from another program.

A transport user communicates with a transport provider over a communications path identified by a transport endpoint.

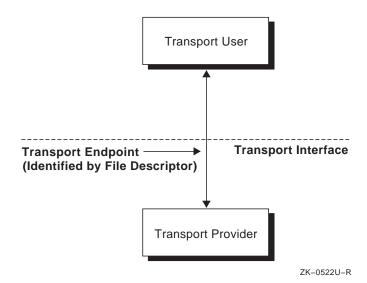
#### Transport endpoints

A **transport endpoint** is created when an application issues a t\_open library call. All of the transport user's requests to the transport provider pass through the endpoint associated with that provider.

The transport user activates a transport endpoint by binding a transport address to it. Once an endpoint is active, a transport user can send data over it. The transport provider routes the data to the appropriate peer user or other destination.

When using a connection-oriented transport service, such as TCP, the transport user must establish a connection between itself and a peer transport user with a t\_connect function, specifying an active endpoint, before sending data. In a transport connection, the transport user initiating the connection is the **active user**, or client, and the peer transport user responding to the connection request is the **passive user**, or server. Figure 3-2 summarizes the relationship between transport providers, transport users, and transport endpoints.

Figure 3-2: A Transport Endpoint



#### 3.2 XTI Features

XTI consists of library calls, header files, and the rules and restrictions elaborating how XTI processes work and interact. This section describes the library calls and header files, as well as the regulations that govern the interaction between communicating processes.

#### 3.2.1 Modes of Service and Execution

Transport users use different service modes and execution modes to determine how data is exchanged with transport providers. The following sections introduce the service modes and execution modes available in XTI.

#### 3.2.1.1 Connection-Oriented and Connectionless Service

In XTI, an endpoint can support one of the following modes of service:

• Connection-oriented transport service

A circuit-oriented service that transfers data over an established connection in a reliable, sequenced manner.

Connection-oriented transport is useful for applications that require long, order dependent and reliable, stream-oriented interactions. With connection-oriented transport, transport users and providers can negotiate the parameters and options that govern data transfer. In addition, because a connection provides identification of both parties, the transport user avoids the overhead of transmitting and resolving addresses during data transfer. A connection also provides a context that logically relates successive units of data.

Connectionless transport service

A message-oriented service that transfers data in self-contained units or datagrams, which have no logical sequence with respect to one another.

Connectionless transport is best suited for applications that have the following qualities:

- Short-term request and response interactions
- Dynamic reconfiguration of connections to multiple endpoints
- No need for the guaranteed, sequential delivery of data

Each data unit is self-contained and has no relationship to previous or successive data units, so the transport provider can route it independently.

## 3.2.1.2 Asynchronous and Synchronous Execution

Execution modes provide a means for transport users to handle completion of functions and receipt of events. An event is an occurrence or happening that is significant to a transport user. XTI supports two execution modes:

#### • Synchronous mode

Waits for transport primitives to complete before returning control to the transport user. Also known as **blocking mode**.

Synchronous mode is suited for applications that want to wait for functions to complete or maintain only a single transport connection. In synchronous mode, the transport user cannot perform other tasks while waiting for a function to complete. For example, if the transport user issues a t\_rcv function in synchronous mode, t\_rcv waits until data is received before returning control to the transport user.

Even while using synchronous mode, it is possible to get some event notification, which the transport user does not ordinarily expect. Such asynchronous events are returned to the user through a special error, TLOOK.

If an asynchronous event occurs while a function is executing, the function returns the TLOOK error; the transport user can then issue the t\_look function to retrieve the event.

#### Asynchronous mode

Returns control to the transport user before transport primitives complete. Also known as **nonblocking mode**.

Asynchronous mode is useful for applications that have long delays between completion of functions and other tasks to perform in the meantime. This mode is also useful for applications that handle multiple connections simultaneously. Many applications handle networking functions in asynchronous mode because they can perform useful work while waiting for particular networking functions to complete. For example, if a transport user issues a t\_rcv function call in asynchronous mode, the function returns control to the user immediately if no data is available. The user periodically polls for data until the data arrives.

By default, all functions that process incoming events operate in synchronous mode, blocking until the task completes. To select asynchronous mode, the transport user specifies the O\_NONBLOCK flag with the t\_open function when the endpoint is created or before executing a function or group of functions with the fcntl operating system call.

For a full discussion of the specific events supported by XTI, see Section 3.2.3.

## 3.2.2 The XTI Library, TLI Library, and Header Files

XTI functions are implemented as part of the XTI library, libxti.a. TLI functions are implemented in a separate TLI library, libtli.a. There are also shared versions of these libraries, libxti.so and libtli.so.

Digital UNIX provides shared library support by default when you link an XTI or TLI application with the XTI or TLI library.

For XTI or TLI applications built in the Digital UNIX Version 1.2 environment to use shared library support, you must relink the required object files with the appropriate library. You do not need to recompile source files.

The first of the following examples illustrates how to relink an XTI application's object files with the XTI shared library; the second illustrates how to relink a TLI application's object files with the TLI shared library:

```
% cc -o XTIapp XTIappmain.o XTIapputil.o -lxti
```

```
% cc -o TLIapp TLIappmain.o TLIapputil.o -ltli
```

To link programs statically with the XTI or TLI libraries (as was the default in Digital UNIX Version 1.2), use the non\_shared option to the cc command.

The following example illustrates how to link an XTI application's object files to the XTI library statically:

```
% cc -non_shared -o XTIapp XTIappmain.o XTIapputil.o -lxti
```

See the cc(1) reference page for more information.

To make a program thread safe, build the program with DECthreads pthreads routines. For more information, see *Guide to DECthreads*.

The few differences between XTI and TLI are described in Section 3.5.2, which also describes how to link your programs with the correct library at compile time.

#### 3.2.2.1 XTI and TLI Header Files

XTI and TLI header files contain data definitions, structures, constants, macros, and options used by the XTI and TLI library calls. An application program must include the appropriate header file to make use of structures or other information a particular XTI or TLI library call requires. Table 3-1 lists the XTI and TLI header files.

Table 3-1: Header Files for XTI and TLI

File Name	Description
<tiuser.h></tiuser.h>	Contains data definitions and structures for TLI applications. You must include this file for all TLI applications.
<xti.h></xti.h>	Contains data definitions and structures for XTI applications. You must include this file for all XTI applications.
<fcntl.h></fcntl.h>	Defines flags for modes of execution for the t_open function. You must include this file for all XTI and TLI applications.

#### Note

Typically, header file names are enclosed in angle brackets (< >). To obtain the absolute path to the header file, prepend /usr/include/ to the information enclosed in the angle brackets. For example, the absolute path for the tiuser.h file is /usr/include/tiuser.h.

### 3.2.2.2 XTI Library Calls

Some of the calls apply to connection-oriented transport (COTS), some to connectionless transport (CLTS), some to connection-oriented transport when used with the orderly release feature (COTS\_ORD), and some to all service modes. A small group of the calls are utility functions and do not apply to a particular service mode. Table 3-2 lists the name, purpose, and service mode of each XTI library call. Each call has an associated reference page by the same name.

Digital UNIX provides XTI reference pages only; it does not provide TLI reference pages. For information about TLI and for the TLI reference pages see the *UNIX System V Programmer's Guide: Networking Interfaces*, which is issued by UNIX System Laboratories, Inc. Digital UNIX provides reference pages for each of the functions. For more information, see the *X/Open CAE Specification: Networking Services*.

Table 3-2: XTI Library Calls

Name of Call	Purpose	Service Mode
t_accept	Accepts a connection request	COTS, COTS_ORD
t_alloc	Allocates memory for a library structure	All
t_bind	Binds an address to a transport endpoint	All
t_close	Closes a transport endpoint	All
t_connect	Establishes a connection with another transport user	COTS, COTS_ORD
t_error	Produces an error message	All
t_free	Frees memory previously allocated for a library structure	All
t_getinfo	Returns protocol-specific information	All
t_getprotaddr <sup>a</sup>	Returns the protocol address	All
t_getstate	Returns the current state for the transport endpoint	All
t_listen	Listens for a connection request	COTS, COTS_ORD
t_look	Returns the current event on the transport endpoint	All
t_open	Establishes a transport endpoint	All
t_optmgmt	Retrieves, verifies, or negotiates protocol options	All
t_rcv	Receives data or expedited data over a connection	COTS, COTS_ORD
t_rcvconnect	Receives the confirmation from a connection request	COTS, COTS_ORD
t_rcvdis	Identifies the cause of a disconnect, and retrieves information sent with a disconnect	COTS, COTS_ORD
t_rcvrel <sup>b</sup>	Acknowledges receipt of an orderly release indication	COTS_ORD
t_rcvudata	Receives a data unit	CLTS
t_rcvuderr	Receives information about an error associated with a data unit	CLTS

Table 3-2: (continued)

Name of Call	Purpose	Service Mode
t_snd	Sends data or expedited data over a connection	COTS, COTS_ORD
t_snddis	Initiates a release on an established connection, or rejects a connection request	COTS, COTS_ORD
$t\_{ t sndrel}^b$	Initiates an orderly release	COTS_ORD
t_sndudata	Sends a data unit	CLTS
t_strerror <sup>a</sup>	Produces and error message string	All
t_sync	Synchronizes the data structures in the transport library	All
t_unbind	Disables a transport endpoint	All

#### Table notes:

- a. This function is supported in XPG4 only.
- Digital UNIX as supplied by Digital does not provide a transport provider that supports the use of COTS\_ORD; therefore, this function returns an error.

XTI supports an orderly release mechanism, t\_sndrel and t\_rcvrel functions. (See Table 3-2 for more information.) However, if your applications need to be portable to the ISO transport layer, we recommend that you do not use this mechanism.

Finally, the XTI header file defines the following constants to identify service modes:

- T\_COTS Connection-oriented transport service (for example, OSI transport)
- T\_CLTS Connectionless transport service (for example, UDP)
- T\_COTS\_ORD Connection-oriented transport service with the orderly release mechanism implemented (for example, TCP)

These service modes are returned by the transport provider in the *servtype* field of the info structure when you create an endpoint with the t\_open function.

#### 3.2.3 Events and States

Each transport provider has a particular state associated with it, as viewed by the transport user. The state of a transport provider and its transition to the next allowable state is governed by outgoing and incoming events, which correspond to the successful return of specified user-level transport functions. Outgoing events correspond to functions that send a request or response to the transport provider, whereas incoming events correspond to functions that retrieve data or event information from the transport provider. This section describes the possible states of the transport provider, the outgoing and incoming events that can occur, and the allowable sequence of function calls.

#### 3.2.3.1 XTI Events

XTI applications must manage asynchronous events. An asynchronous event is identified by a mnemonic which is defined as a constant in the XTI header file. Table 3-3 lists the name, purpose, and service mode for each type of asynchronous event in XTI.

Table 3-3: Asynchronous XTI Events

<b>Event Name</b>	Purpose	Service Mode
T_CONNECT	The transport provider received a connection response. This event usually occurs after the transport user issues the t_connect function.	COTS, COTS_ORD
T_DATA	The transport provider received <b>normal data</b> , which is all or part of a <b>Transport Service Data Unit</b> (TSDU).	COTS, CLTS, COTS_ORD
T_DISCONNECT	The transport provider received a disconnect request. This event usually occurs after the transport user issues data transfer functions, the t_accept function, or the t_snddis function.	COTS, COTS_ORD
T_EXDATA	The transport provider received expedited data.	COTS, COTS_ORD
T_GODATA	The flow control restrictions on the flow of normal data are lifted. The transport user can send normal data again.	COTS, CLTS, COTS_ORD

Table 3-3: (continued)

<b>Event Name</b>	Purpose	Service Mode
T_GOEXDATA	The flow control restrictions on the flow of expedited data are lifted.  The transport user can send expedited data again.	COTS, COTS_ORD
T_LISTEN	The transport provider received a connection request from a remote user. This event occurs only when the file descriptor is bound to a valid address and no transport connection is established.	COTS, COTS_ORD
T_ORDREL	The transport provider received a request for an orderly release.	COTS_ORD
T_UDERR	An error was found on a datagram that was previously sent. This event usually occurs after the transport user issues the t_rcvudata or t_unbind functions.	CLTS

XTI stores all events that occur at a transport endpoint.

If using a synchronous mode of execution, the transport user returns from the function it was executing with a value of -1 and then checks for a value of TLOOK in t\_errno and retrieves the event with the t\_look function. In asynchronous mode, the transport user continues doing productive work and periodically checks for new events.

Every event at a transport endpoint is consumed by a specific XTI function, or it remains outstanding. Exceptions are the T\_GODATA and T\_GOEXDATA events, which are cleared by retrieving them with t\_look. Thus, once the transport user receives a TLOOK error from a function, subsequent calls to that function or a different function continue to return the TLOOK error until the transport user consumes the event. Table 3-4 summarizes the consuming functions for each asynchronous event.

**Table 3-4: Asynchronous Events and Consuming Functions** 

Event	Cleared by t_look	Consuming Function(s)
T_CONNECT	No	t_connect, t_rcvconnect
T_DATA	No	t_rcv, t_rcvudata
T_DISCONNECT	No	t_rcvdis
T_EXDATA	No	t_rcv
T_GODATA	Yes	t_snd, t_sndudata
T_GOEXDATA	Yes	t_snd
T_LISTEN	No	t_listen
T_ORDREL	No	t_rcvrel
T_UDERR	No	t_rcvuderr

Table 3-5 lists the events that cause a specific XTI function to return the TLOOK error. This information may be useful when you structure the event checking mechanisms in your XTI applications.

Table 3-5: XTI Functions that Return TLOOK

Function	<b>Events Causing TLOOK</b>		
t_accept	T_DISCONNECT, T_LISTEN		
t_connect	T_DISCONNECT, T_LISTEN a		
t_listen	T_DISCONNECT b		
t_rcv	T_DISCONNECT, T_ORDREL <sup>c</sup>		
t_rcvconnect	T_DISCONNECT		
t_rcvrel	T_DISCONNECT		
t_rcvudata	T_UDERR		
t_snd	T_DISCONNECT, T_ORDREL		
t_snddis	T_DISCONNECT		
t_sndrel	T_DISCONNECT		
t_sndudata	T_UDERR		
t unbind	T LISTEN, T DATA d		

#### Table notes:

- a. This event occurs only when  $t\_connect$  is issued for an endpoint that was bound with a qlen > 0, and has a pending connection indication.
- b. This event indicates a disconnect on an outstanding connection indication.
- c. This occurs only when all pending data has been read.
- d. T\_DATA may only occur for the connetionless mode.

Each XTI function manages one transport endpoint at a time. It is not possible to wait for several events from different sources, particularly from several transport connections at a time. The Digital UNIX implementation of XTI allows the transport user to monitor input and output on a set of file descriptors with the poll function. See poll(2) for more information.

#### 3.2.3.2 XTI States

XTI uses eight states to manage communication over a transport endpoint. Both the active and passive user have a unique state that reflects the function in process.

Table 3-6 describes the purpose of each XTI state. A service mode of COTS indicates the state occurs regardless of whether or not orderly service is implemented. A service mode of COTS\_ORD indicates the state occurs only when orderly service is implemented.

Table 3-6: XTI States

State	Description	Service Mode
T_UNINIT	Uninitialized. Initial and final state of the interface. To establish a transport endpoint, the user must issue a t_open.	COTS, CLTS, COTS_ORD
T_UNBIND	Unbound. The user can bind an address to a transport endpoint or close a transport endpoint.	COTS, CLTS, COTS_ORD
T_IDLE	Idle. The active user can establish a connection with a passive user (COTS), disable a transport endpoint (COTS, CLTS), or send and receive data units (CLTS). The passive user can listen for a connection request (COTS).	COTS, CLTS, COTS_ORD
T_OUTCON	Outgoing connection pending. The active user can receive confirmations for connection requests.	COTS, COTS_ORD

Table 3-6: (continued)

State	Description	Service Mode
T_INCON	Incoming connection pending. The passive user can accept connection requests.	COTS, COTS_ORD
T_DATAXFER	Data transfer. The active user can send data to and receive data from the passive user. The passive user can send data to and receive data from the active user.	COTS, COTS_ORD
T_OUTREL	Outgoing orderly release. The user can respond to an orderly release indication.	COTS_ORD
T_INREL	Incoming orderly release. The user can send an orderly release indication.	COTS_ORD

If you are writing a connection-oriented application, note that your program can release a connection at any time during the connection-establishment state or data-transfer state.

# 3.2.4 Tracking XTI Events

The XTI library keeps track of outgoing and incoming events to manage the legal states of transport endpoints. The following sections describe these outgoing and incoming events.

#### 3.2.4.1 Outgoing Events

Outgoing events are caused by XTI functions that send a request or response to the transport provider. An outgoing event occurs when a function returns successfully. Some functions produce different events, depending on the following values:

ocnt	A count of outstanding connection indications (those passed to the transport user but not yet accepted or rejected). This count is only meaningful for the current transport endpoint $(fd)$ .
fd	The file descriptor of the current transport endpoint.
resfd	The file descriptor of the endpoint where a connection will be accepted.

Table 3-7 describes the outgoing events available in XTI. A service mode of COTS indicates the event occurs for a connection-oriented service regardless of whether or not orderly service is implemented. A service mode of

COTS\_ORD indicates the event occurs only when orderly service is implemented.

**Table 3-7: Outgoing XTI Events** 

Event	Description	Service Mode
opened	Successful return of t_open function.	COTS, CLTS, COTS_ORD
bind	Successful return of t_bind function.	COTS, CLTS, COTS_ORD
optmgmt	Successful return of t_optmgmt function.	COTS, CLTS, COTS_ORD
unbind	Successful return of t_unbind function.	COTS, CLTS, COTS_ORD
closed	Successful return of t_close function.	COTS, CLTS, COTS_ORD
connect1	Successful return of t_connect function in synchronous execution mode.	COTS, COTS_ORD
connect2	The t_connect function returned the TNODATA error in asynchronous mode, or returned the TLOOK error because a disconnect indication arrived on the transport endpoint.	COTS, COTS_ORD
accept1	Successful return of $t_{accept}$ function, where $ocnt == 1$ and $fd == resfd$ .	COTS, COTS_ORD
accept2	Successful return of t_accept function, where $ocnt == 1$ and $fd != resfd$ .	COTS, COTS_ORD
accept3	Successful return of $t_{accept}$ function, where $ocnt > 1$ .	COTS
snd	Successful return of t_snd function.	COTS
snddis1	Successful return of t_snddis function, where ocnt <= 1.	COTS, COTS_ORD
snddis2	Successful return of t_snddis function, where $ocnt > 1$ .	COTS, COTS_ORD
sndrel	Successful return of t_sndrel function.	COTS_ORD
sndudata	Successful return of t_sndudata function.	CLTS

## 3.2.4.2 Incoming Events

Incoming events are caused by XTI functions that retrieve data or events from the transport provider. An incoming event occurs when a function returns successfully. Some functions produce different events, depending on the value of the ocnt variable. This variable is a count of outstanding connection indications (those passed to the transport user but not yet accepted or rejected). This count is only meaningful for the current transport endpoint (fd).

The pass\_conn incoming event is not associated directly with the successful return of a function on a given endpoint. The pass\_conn event occurs on the endpoint that is being passed a connection from the current endpoint. No function occurs on the endpoint where the pass\_conn event occurs.

Table 3-8 describes the incoming events available in XTI. A service mode of COTS indicates the event occurs regardless of whether or not orderly service is implemented. A service mode of COTS\_ORD indicates the event occurs only when orderly service is implemented.

Table 3-8: Incoming XTI Events

Event	Description	Service Mode
listen	Successful return of the t_listen function	COTS, COTS_ORD
rcvconnect	Successful return of the t_rcvconnect function	COTS, COTS_ORD
rcv	Successful return of the t_rcv function	COTS, COTS_ORD
rcvdis1	Successful return of the t_rcvdis function, where ocnt == 0	COTS, COTS_ORD
rcvdis2	Successful return of the t_rcvdis function, where ocnt == 1	COTS, COTS_ORD
rcvdis3	Successful return of the t_rcvdis function, where $ocnt > 1$	COTS, COTS_ORD
rcvrel	Successful return of the t_rcvrel function	COTS_ORD
rcvudata	Successful return of the t_rcvudata function	CLTS
rcvuderr	Successful return of the t_rcvuderr function	CLTS
pass_conn	Successfully received a connection that was passed from another transport endpoint	COTS, COTS_ORD

## 3.2.5 A Map of XTI Functions, Events, and States

This section describes the relationship among XTI functions, outgoing and incoming events, and states. Since XTI has well-defined rules about state transitions, it is possible to know the next allowable state given the current state and most recently received event. This section provides detailed tables that map the current event and state to the next allowable state.

This section excludes the t\_getstate, t\_getinfo, t\_alloc, t\_free, t\_look, t\_sync, and t\_error functions from discussions of state transitions. These utility functions do not affect the state of the transport interface, so they can be issued from any state except the uninitialized (T\_UNINIT) state.

To use Table 3-9, Table 3-10, and Table 3-11, find the row that matches the current incoming or outgoing event and the column that matches the current state. Go to the intersection of the row and column to find the next allowable state. A dash (—) at the intersection indicates an invalid combination of event and state. Some state transitions are marked by a number in parentheses that indicates an action that the transport user must take. The numbers and their meanings are listed at the end of the appropriate table.

Table 3-9 shows the state transitions for initialization and deinitialization functions, functions that are common to both the connection-oriented and connectionless modes of service. For example, if the current event and state are bind and T\_UNBND, the next allowable state is T\_IDLE. In addition, the transport user must set the count of outstanding connection indications to zero, as indicated by the numeral 1.

Table 3-9: State Transitions for Initialization of Connection-Oriented or Connectionless Transport Services

Event	T_UNINIT State	T_UNBND State	T_IDLE State
opened	T_UNBND	_	_
bind	_	T_IDLE <sup>a</sup>	_
unbind	_	_	T_UNBND
closed	_	T_UNINIT	T_UNINIT

#### Table notes:

a. Set the count of outstanding connection indications, ocnt, to 0.

Table 3-10 shows the state transitions for data transfer functions in connectionless transport services.

Table 3-10: State Transitions for Connectionless Transport Services

Event	T_IDLE State
sndudata	T_IDLE
rcvudata	T_IDLE
rcvuderr	T_IDLE

Table 3-11 and Table 3-12 show the transitions for connection, release, and data transfer functions in connection-oriented transport services for incoming and outgoing events. For example, if the current event and state are accept2 and T\_INCON, the next allowable state is T\_IDLE, providing the transport user decrements the count of outstanding connection indications and passes a connection to another transport endpoint.

Table 3-11: State Transitions for Connection-Oriented Transport Services: Part 1

Event	T_IDLE State	T_OUTCON State	T_INCON State	T_DATAXFER State
connect1	T_DATAXFER	_	_	_
connect2	T_OUTCON	_	_	_
rcvconnect	_	T_DATAXFER	_	_
listen	T_INCON (a)	_	T_INCON (a)	_
accept1	_	_	T_DATAXFER (a)	_
accept2	_	_	T_IDLE (b, c)	_
accept3	_	_	T_INCON (b, c)	_
snd	_	_	_	T_DATAXFER
rcv	_	_	_	T_DATAXFER
snddis1	_	T_IDLE	T_IDLE (b)	T_IDLE
snddis2	_	_	T_INCON (b)	_
rcvdis1	_	T_IDLE	_	T_IDLE
rcvdis2	_		T_IDLE (b)	_

Table 3-11: (continued)

Event	T_IDLE State	T_OUTCON State	T_INCON State	T_DATAXFER State
rcvdis3	_	_	T_INCON (b)	_
sndrel	_	_	_	T_OUTREL
rcvrel	_	_	_	T_INREL
pass_conn	T_DATAXFER	_	_	_
optmgmt	T_IDLE	T_OUTCON	T_INCON	T_DATAXFER
closed	T_UNINIT	T_UNINIT	T_UNINIT	T_UNINIT

## Table notes:

- a. Increment the count of outstanding connection indications.
- b. Decrement the count of outstanding connection indications.
- c. Pass a connection to another transport endpoint, as indicated in the  $t\_accept$  function.

Table 3-12: State Transitions for Connection-Oriented Transport Services: Part 2

Event	T_OUTREL State	T_INREL State	T_UNBND State
connect1	_	_	_
connect2	_	_	_
rcvconnect	_	_	_
listen	_	_	_
accept1	_	_	_
accept2	_	_	_
accept3	_	_	_
snd	_	T_INREL	_
rcv	T_OUTREL	_	_
snddis1	T_IDLE	T_IDLE	_
snddis2	_	_	_
rcvdis1	T_IDLE	T_IDLE	_

Table 3-12: (continued)

Event	T_OUTREL State	T_INREL State	T_UNBND State
rcvdis2	_	_	_
rcvdis3	_	_	_
sndrel	_	T_IDLE	_
rcvrel	T_IDLE	_	_
pass_conn	_	_	T_DATAXFER
optmgmt	T_OUTREL	T_INREL	T_UNBND
closed	T_UNINIT	T_UNINIT	_

# 3.2.6 Synchronization of Multiple Processes and Endpoints

In general, if you use multiple processes, you need to synchronize them carefully to avoid violating the state of the interface.

Although transport providers treat all transport users of a transport endpoint as a single user, the following situations are possible:

- One process can create several transport endpoints simultaneously.
- Multiple processes can share a single endpoint simultaneously.

For a single process to manage several endpoints in synchronous execution mode, the process must manage the actions on each endpoint serially instead of in parallel. Optionally, you can write a server to manage several endpoints at once. For example, the process can listen for an incoming connection indication on one endpoint and accept the connection on a different endpoint, so as not to block incoming connections. Then, the application can fork a child process to service the requests from the new connection.

Multiple processes that share a single endpoint must coordinate actions to avoid violating the state of the interface. To do this, each process calls the t\_sync function, which retrieves the current state of the transport provider, before issuing other functions. If all processes do not cooperate in this manner, another process or an incoming event can change the state of the interface.

Similarly, while several endpoints can share the same protocol address, only one can listen for incoming connections. Other endpoints sharing the protocol address can be in data transfer state or in the process of establishing a connection without causing a conflict. This means that an address can have only one server, but multiple endpoints can call the address at the same time.

# 3.3 Using XTI

This section presents guidelines to help you sequence functions, manage states, and use XTI options. It then describes the steps required to write both connection-oriented and connectionless programs to XTI.

# 3.3.1 Guidelines for Sequencing Functions

Figure 3-3 shows the typical sequence of functions and state transitions for an active user and passive user communicating with a connection-oriented transport service in nonblocking mode. The solid lines in the figure show the state transitions for the active user, while the dashed lines show the transitions for the passive user. Each line represents the call of a function, while each ellipse represents the resulting state. This example does not include the orderly release feature.

Figure 3-3: State Transitions for Connection-Oriented Transport Services

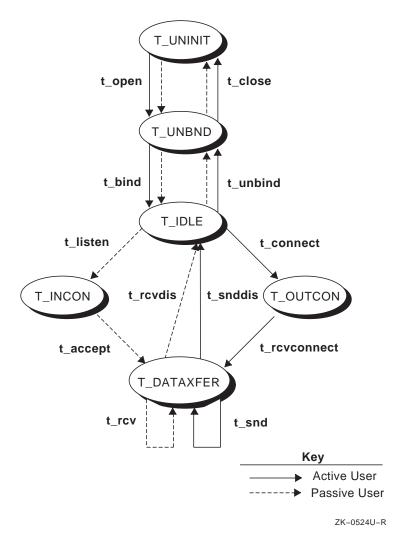


Figure 3-4 shows the typical sequence of functions and transitions in state for two users communicating with the connectionless transport service. Each line in the figure represents the call of a function, while each ellipse represents the resulting state. Both users are represented by solid lines.

t\_open t\_close

T\_UNINIT

t\_open t\_close

T\_UNBND

t\_bind t\_unbind

t\_rcvudata t\_sndudata

Figure 3-4: State Transitions for the Connectionless Transport Service

ZK-0525U-R

## 3.3.2 State Management by the Transport Provider

All transport providers take the following actions with respect to states:

- Keep a record of the state of the interface as seen by the transport user.
- Reject any requests or responses that would place the interface out of state and return an error. In this case, the state does not change. For example, if the user passes data with a function and the interface is not in T\_DATAXFER state, the transport provider does not accept or forward the data.

The uninitialized state (T\_UNINIT) serves two purposes:

- The initial state of a transport endpoint. The transport user must initialize
  and bind the transport endpoint before the transport provider views it as
  active.
- The final state of a transport endpoint. The transport provider must view the endpoint as unused. When the transport user issues the t\_close function, the transport provider is closed, and the resources associated with the transport library are freed for use by another endpoint.

## 3.3.3 Writing a Connection-Oriented Application

Follow these steps to write a connection-mode application:

- 1. Initialize an endpoint
- 2. Establish a connection
- 3. Transfer data
- 4. Release a connection
- 5. Deinitialize an endpoint

## 3.3.3.1 Initializing an Endpoint

To initialize an endpoint, complete the following steps:

- 1. Open the endpoint
- 2. Bind an address to the endpoint
- 3. Negotiate protocol options

Note that the steps described here for initializing an endpoint for connectionoriented service are identical for connectionless service.

### **Opening a Transport Endpoint**

Both connection-oriented and connectionless applications must open a transport endpoint using the t\_open function. The syntax of the t\_open function is as follows:

fd = t\_open (name, oflag, &info);

In the preceding statement:

fd

Identifies the file descriptor for the endpoint. You use the file descriptor in subsequent calls to identify this transport endpoint.

The t\_open function returns a file descriptor upon successful completion. Otherwise, t\_open returns a value of -1, and  $t_{errno}$  is set to one of the values described in Section 3.7. (For multithreaded applications,  $t_{errno}$  is thread specific.)

name

Identifies the transport provider to be accessed. Currently, the XTI implemented in Digital UNIX uses pathnames to device special files to identify transport providers, which is the same method as in the AT&T TLI. The device special files on a Digital UNIX system corresponding to TCP or UDP transport providers reside in the

/dev/streams/xtiso directory. If you use a different transport

provider, see its documentation for the correct device name.

#### Note

Using the special device with any mechanism other than XTI/TLI, for example, direct open, read, or write calls, is illegal and will generate undefined results.

#### oflag

Specifies whether the endpoint will block on functions to wait for completion. Specify O\_RDWR to indicate that the endpoint supports reading and writing by functions and blocks on them, or specify the bitwise inclusive OR of O\_RDWR and O\_NONBLOCK to indicate the endpoint supports reading and writing by functions but does not block on them. You must use O\_RDWR optionally with OR with O\_NONBLOCK for the mode flag passed to t\_open. In other words, the XTI specification forbids the use of O\_RDONLY or O\_WRONLY to make the endpoint either read-only or write-only as expected.

#### info

Returns the pointer to a structure containing the default characteristics of the transport provider. You use these characteristics to determine subsequent calls. The info parameter points to the  $t\_info$  structure. See  $t\_open(3)$  for more information.

If you are designing a protocol-independent program, you can determine data buffer sizes by accessing the information that the  $t\_open$  function returns about the  $t\_info$  structure. If the transport user exceeds the allowed data size, you receive an error. Alternatively, you can use the  $t\_alloc$  function to allocate data buffers.

See t\_open(3) for more information.

The following is an example of the t\_open function for the TCP transport provider:

```
if ( (newfd = t_open( "/dev/streams/xtiso/tcp" , O_RDWR , NULL) ) == -1 )
    {
       (void) t_error("could not open tcp transport");
       exit (1);
    }
```

#### Binding an Address to the Endpoint

Once you open an endpoint, you need to bind a protocol address to the endpoint. By binding the address, you activate the endpoint. In connection mode, you also direct the transport provider to begin accepting connection indications or servicing connection requests on the transport endpoint. To determine if the transport provider has accepted a connection indication, you

can issue the t\_listen function. In connectionless mode, once you bind the address, you can send or receive data units through the transport endpoint.

To bind an address to an endpoint, issue the t\_bind function with the following syntax:

t bind (fd,req,ret);

In the preceding statement:

fd

Identifies the file descriptor for the endpoint, which is returned by the t\_open function.

req

Specifies a pointer to the structure containing the address you wish to bind to the endpoint.

ret

Returns a pointer to the structure containing the address that XTI bound to the endpoint.

See t\_bind(3) for more information.

If the transport provider supports the automatic generation of addresses, you have the following choices in binding addresses:

- Set req to a null pointer if you do not wish to specify an address. The transport provider will assign an address to the transport endpoint.
- Set ret to a null pointer if you do not need to determine the actual address that was bound to the endpoint.
- Set req and ret to null pointers if you want the transport provider to both assign the address and not notify you of what it was.
- If the address that you requested in req is not available, the transport provider will assign an appropriate address.

To determine if the transport provider generates addresses, do not specify one in the t\_bind function (set req to a null pointer). If the transport provider supplies addresses, the function returns an assigned address in the ret field. If the transport provider does not supply addresses, the function returns an error of TNOADDR.

If you accept a connection on an endpoint that is used for listening for connection indications, the bound address is busy for the duration of the connection. You cannot bind any other endpoint for listening on that same address while the initial listening endpoint is actively transferring data or in T IDLE state.

You can use the gethostbyname routine, described in Section 4.2.3.2, to obtain host information when either TCP or UDP is the underlying transport provider.

If you use a method to retrieve host information other than the gethostbyname routine, consider the following:

- Your applications must pass XTI functions a socket address in the format that the transport provider expects. For XTI over TCP/IP, the expected address format is a sockaddr in structure.
- Your applications also need to pass a transport provider identifier to XTI functions. In Digital UNIX, this identifier must already be in the format of a pathname to the device special file for the transport provider.

The t\_bind function returns a value of 0 upon successful completion. Otherwise, it returns a value of -1, and  $t_{errno}$  is set to one of the values described in Section 3.7. (For multithreaded applications,  $t_{errno}$  is thread specific.)

#### 3.3.3.2 Using XTI Options

XPG3 and XPG4 implement option management differently.

In XPG3, option management is handled exclusively by the t\_optmgmt function. In XPG4, several functions contain an opt argument which is used to convey options between a transport user and the transport provider.

For more information, see Section 3.6.6.

## 3.3.3.3 Establishing a Connection

The connection establishment phase typically consists of the following actions:

- 1. A passive user, or server, listens for a connection request.
- 2. An active user, or client, initiates a connection.
- 3. A passive user, or server, accepts a connection request and a connection indication is received.

These steps are described in the following sections.

#### **Listening for Connection Indications**

The passive user issues the t\_listen function to look for enqueued connection indications. If the t\_listen function finds a connection indication at the head of the queue, it returns detailed information about the connection indication and a local sequence number that identifies the indication. The number of outstanding connection indications that can be

queued is limited by the value of the *qlen* parameter that was accepted by the transport provider when the t\_bind function was issued.

By default, the t\_listen function executes synchronously by waiting for a connection indication to arrive before returning control to the user. If you set the O\_NONBLOCK flag of the t\_open function or the fcntl function for asynchronous execution, the t\_listen function checks for an existing connection indication and returns an error of TNODATA if none is available.

To listen for connection requests, issue the t\_listen function with the following syntax:

#### t\_listen (fd,call);

In the preceding statement:

fd

Specifies the file descriptor of the endpoint where connection indications arrive.

call

Returns a pointer to information describing the connection indication.

See t\_listen(3) for more information.

The t\_listen function returns a value of 0 upon successful completion. Otherwise, it returns a value of -1, and  $t_{errno}$  is set to one of the values described in Section 3.7. (For multithreaded applications,  $t_{errno}$  is thread specific.)

## **Initiating Connections**

A connection is initiated in either synchronous or asynchronous mode. In synchronous mode, the active user issues the t\_connect function, which waits for the passive user's response before returning control to the active user. In asynchronous mode, t\_connect initiates a connection but returns control to the active user before a response to the connection arrives. Then, the active user can determine the status of the connection request by issuing the t\_rcvconnect function. If the passive user accepted the request, the t\_rcvconnect function returns successfully and the connection establishment phase is complete. If a response has not been received yet, the t\_rcvconnect function returns an error of TNODATA. The active user should issue the t\_rcvconnect function again later.

To initiate a connection, issue the t\_connect function with the following syntax:

t\_connect (fd,sndcall,rcvcall);

In the preceding statement:

fd

Specifies the file descriptor of the endpoint where the connection will be established.

sndcall

Points to a structure containing information that the transport provider needs to establish the connection.

rcvcall

Points to a structure containing information that the transport provider associates with the connection that was just established.

See t connect(3) for more information.

The t\_connect function returns a value of 0 upon successful completion. Otherwise, it returns a value of -1, and  $t_{errno}$  is set to one of the values described in Section 3.7. (For multithreaded applications,  $t_{errno}$  is thread specific.)

## **Accepting Connections**

When the passive user accepts a connection indication, it can issue the t\_accept function on the same endpoint (the endpoint where it has been listening with t\_listen) or a different endpoint.

If the passive user accepts on the same endpoint, the endpoint can no longer receive and enqueue incoming connection indications. The protocol address that is bound to the endpoint remains busy for the duration it is active. No other transport endpoints can be bound to the same protocol address as the listening endpoint. That is, no other endpoints can be bound until the passive user issues the t\_unbind function. Further, before the connection can be accepted on the same endpoint, the passive user must respond (with either the t\_accept or t\_snddis functions) to all previous connection indications that it has received. Otherwise, t accept returns an error of TBADF.

If the passive user accepts the connection on a different endpoint, the listening endpoint can still receive and enqueue incoming connection requests. The different endpoint must already be bound to a protocol address and be in the T\_IDLE state. If the protocol address is the same as for the endpoint where the indication was received, the *qlen* parameter must be set to zero (0).

For both types of endpoints, t\_accept will fail and return an error of TLOOK if there are connect or disconnect indications waiting to be received.

To accept a connection, issue the t\_accept function with the following syntax:

t\_accept (fd,resfd,call);

In the preceding statement:

fd

Specifies the file descriptor of the endpoint where the connection indication arrived.

resfd

Specifies the file descriptor of the endpoint where the connection will be established.

call

Points to information needed by the transport provider to establish the connection.

See t accept(3) for more information.

The t\_accept function returns a value of 0 upon successful completion. Otherwise, it returns a value of -1, and  $t_{errno}$  is set to one of the values described in Section 3.7. (For multithreaded applications,  $t_{errno}$  is thread specific.)

## 3.3.3.4 Transferring Data

Once a connection is established between two endpoints, the active and passive users can transfer data in full-duplex fashion over the connection. This phase of connection-oriented service is known as the data transfer phase. The following sections describe how to send and receive data during the data transfer phase.

## **Sending Data**

Transport users can send either normal or expedited data over a connection with the t\_snd function. Normally, t\_snd sends successfully and returns the number of bytes accepted if the transport provider can immediately accept all the data. If the data cannot be accepted immediately, the result of t\_snd depends on whether it is executing synchronously or asynchronously.

By default, the t\_snd function executes synchronously and waits if flow control conditions prevent the transport provider from accepting the data. The function blocks until one of the following conditions becomes true:

- The flow control conditions clear, and the transport provider can accept a new data unit. The t snd function returns successfully.
- A disconnect indication is received. The t\_snd function returns with an error of TLOOK. If you call the t\_look function, it returns the T\_DISCONNECT event. Any data in transit is lost.
- An internal problem occurs. The t\_snd function returns with an error of TSYSERR. Any data in transit is lost.

If the O\_NONBLOCK flag was set when the endpoint was created, t\_snd executes asynchronously and fails immediately if flow control restrictions exist. In some cases, only part of the data was accepted by the transport provider, so t\_snd returns a value that is less than the number of bytes that you requested to be sent. At this point, you can do one of the following:

- Issue t\_snd again with the remaining data.
- Check with the t\_look function to see if the flow control restrictions are lifted, then resend the data. The t\_look function is described at the end of this chapter.

To send data or expedited data over a connection, issue the t\_snd function with the following syntax:

**t\_snd** (fd,buf,nbytes,flags);

In the preceding statement:

fd

Specifies the file descriptor of the endpoint over which data should be sent.

buf

Points to the data.

nbytes

Specifies the number of bytes of data to be sent.

flags

Specifies any optional flags, such as the following:

• T EXPEDITED

Send the data as expedited data. The expedited data is subject to the interpretations of the transport provider. Some transport providers don't support expedited data.

T MORE

Indicates that another  $t\_snd$  function will follow with more data for the current TSDU or ETSDU. The end of the TSDU or ETSDU is indicated by a  $t\_snd$  function without  $T\_MORE$  set.

Some transport providers do not support the concept of a TSDU or ETSDU, so the T\_MORE flag is not meaningful. To find out if the transport provider supports TSDUs and ETSDUs, check the <code>info</code> argument of the t\_open or t\_getinfo function. If the <code>tsdu</code> field of <code>info</code> is greater than zero (0), the transport provider supports a record-oriented mode, and the return value indicates the maximum size of a TSDU. If the <code>tsdu</code> field is zero, the transport provider supports a stream-oriented mode of sending data. The T\_MORE

flag has no bearing on how the data is packaged for transfer at layers below the transport interface.

See t snd(3) for more information.

The t\_snd function returns a value of 0 upon successful completion. Otherwise, it returns a value of -1, and  $t_{errno}$  is set to one of the values described in Section 3.7. (For multithreaded applications,  $t_{errno}$  is thread specific.)

## **Receiving Data**

Transport users can receive either normal or expedited data over a connection with the t\_rcv function. Typically, if data is available, t\_rcv returns the data. If the connection has been disconnected, t\_rcv returns immediately with an error. If data is not available, but the connection still exists, t\_rcv behaves differently depending on the mode of execution:

- By default, t\_rcv executes synchronously and waits for one of the following to arrive:
  - Data
  - A disconnect indication
  - A signal

Instead of issuing t\_rcv and waiting, you can issue the t\_look function and check for the T\_DATA or T\_EXDATA events.

• If you set the O\_NONBLOCK flag, t\_rcv executes asynchronously and fails with an error of TNODATA if no data is available. You should continue to poll for data by issuing the t\_rcv or t\_look functions.

To receive data, issue the t\_rcv function with the following syntax:

t\_rcv (fd,buf,nbytes,flags);

In the preceding statement:

fd

Specifies the file descriptor of the endpoint through which data arrives.

buf

Points to a buffer where the data that is received will be placed.

nbytes

Specifies the size of the buffer.

flags

Returns the following optional flags that apply to the received data:

- T\_EXPEDITED Indicates the data received is expedited data.
- T\_MORE Indicates that there is more data for the TSDU or ETSDU that must be received by using additional t\_rcv functions. The end of the TSDU or ETSDU is indicated by a t\_rcv function with the T\_MORE flag not set. Some transport providers do not support the concept of a TSDU or ETSDU, so the T\_MORE flag is not meaningful. To find out if the transport provider supports TSDUs and ETSDUs, check the <code>info</code> argument of the t\_open or t\_getinfo function.

If you retrieve part of a TSDU and expedited data arrives, the receipt of the remainder of the TSDU is suspended until you process the ETSDU. For example, if you received data with T\_MORE set and then received data with T\_EXPEDITED and T\_MORE set, this indicates a situation where expedited data arrived in the middle of your receipt of a TSDU. After you retrieve the full ETSDU, you can retrieve the remainder of the TSDU. It is the responsibility of the application programmer to remember that the receipt of normal data has been interrupted.

See  $t_rcv(3)$  for more information.

The  $t\_rcv$  function returns a value of 0 upon successful completion. Otherwise, it returns a value of -1, and  $t\_errno$  is set to one of the values described in Section 3.7. (For multithreaded applications,  $t\_errno$  is thread specific.)

## 3.3.3.5 Releasing Connections

XTI supports two ways to release connections: abortive release and orderly release. All transport providers support abortive release. Orderly release is not provided by all transport providers. For example, the OSI transport supports only abortive release, while TCP supports abortive release and optionally, orderly release.

### **Abortive Release**

An abortive release, which can be requested by the transport user or the transport provider, aborts a connection immediately. Abortive releases cannot be negotiated, and once the abortive release is requested, there is no guarantee that user data will be delivered.

Transport users can request an abortive release in either the connection establishment or data transfer phases. During connection establishment, a transport user can use the abortive release to reject a connection request. In data transfer phase, either user can release the connection at any time. If a transport provider requests an abortive release, both users are informed that

the connection no longer exists.

To request an abortive release or to reject a connection indication, issue the t\_snddis function with the following syntax:

t snddis (fd,call);

In the preceding statement:

fd

Specifies the file descriptor of the endpoint.

call

Points to the information associated with the abortive release. This field is only meaningful if the transport user wants to send user data with the disconnect request, or if the transport user is rejecting a connection indication.

See t\_snddis(3) for more information.

Transport users are notified about abortive releases through the T\_DISCONNECT event. If your program receives a T\_DISCONNECT event, it must issue the t\_rcvdis function to retrieve information about the disconnect and to consume the T\_DISCONNECT event. The following is the syntax of the t\_rcvdis function:

t\_rcvdis (fd, discon);

In the preceding statement:

fd

Specifies the file descriptor of the endpoint where the connection existed.

discon

Points to information about the disconnect.

See t rcvdis(3) for more information.

Both t\_snddis and t\_rcvdis return a value of 0 upon successful completion. Otherwise, they return a value of -1, and  $t_{errno}$  is set to one of the values described in Section 3.7. (For multithreaded applications,  $t_{errno}$  is thread specific.)

### **Orderly Release**

An orderly release allows for release of a connection without loss of data. Orderly release is not provided by all transport providers. If the transport provider returned a service type of T\_COTS\_ORD with the t\_open or t\_getinfo functions, orderly release is supported. Transport users can

request an orderly release during the data transfer phase. The typical sequence of orderly release is as follows:

- 1. The active user issues the t\_sndrel function to request an orderly release of the connection.
- 2. The passive user receives the T\_ORDREL event indicating the active user's request for the orderly release and issues the t\_rcvrel function to indicate the request was received and consume the T\_ORDREL event.
- 3. When ready to disconnect, the passive user issues the t\_sndrel function.
- 4. The active user responds by issuing the t\_rcvrel function.

To initiate an orderly release, use the t\_sndrel function which has the following syntax:

### t sndrel (fd);

In the preceding statement:

fd

Specifies the field descriptor of the endpoint.

The transport user cannot send more data over the connection after it issues the t\_sndrel function. The transport user can, however, continue to receive data until it receives an orderly release indication (the T\_ORDREL event).

See t sndrel(3) for more information.

To acknowledge the receipt of an orderly release indication, issue the t\_rcvrel function with the following syntax:

### t\_rcvrel (fd);

In the preceding statement:

fd

Specifies the file descriptor of the endpoint.

After a transport user receives an orderly release indication (T\_ORDREL), it cannot receive more data. (If the user attempts to do so, the function blocks indefinitely.) The transport user can, however, continue to send data until it issues the t\_sndrel function.

See t\_rcvrel(3) for more information.

Both t\_sndrel and t\_rcvrel return a value of 0 upon successful completion. Otherwise, they return a value of -1, and t\_errno is set to one of the values described in Section 3.7. (For multithreaded applications, t\_errno is thread specific.)

## 3.3.3.6 Deinitializing Endpoints

When you are finished using an endpoint, you deinitialize it by unbinding and closing the endpoint with the t\_unbind and t\_close functions. Note that the steps described here for deinitializing an endpoint with connection-oriented service are identical to those for connectionless service.

When you unbind the endpoint, you disable the endpoint so that the transport provider no longer accepts requests for it. The syntax for t\_unbind is as follows:

#### t unbind (fd);

In the preceding statement:

fd

Specifies the file descriptor of the endpoint.

See t unbind(3) for more information.

By closing the endpoint, you inform the transport provider that you are finished with it and you free any library resources associated with the endpoint.

You should call t\_close when the endpoint is in the T\_UNBND state. However, this function does not check state information, so it may be called to close a transport endpoint from any state.

If you close an endpoint that is not in the T\_UNBND state, the library resources associated with the endpoint are freed automatically, and the file associated with the endpoint is closed. If there are no other descriptors in this or any other process that references the endpoint, the transport connection is broken.

To close the endpoint, issue the t\_close function. The syntax for t\_close is as follows:

### t\_close (fd);

In the preceding statement:

fd

Specifies the file descriptor of the endpoint.

See t\_close(3) for more information.

Both t\_unbind and t\_close return a value of 0 upon successful completion. Otherwise, they return a value of -1, and t\_errno is set to one of the values described in Section 3.7. (For multithreaded applications, t errno is thread specific.)

## 3.3.4 Writing a Connectionless Application

This section describes the steps required to write a connectionless mode application:

- 1. Initializing an endpoint
- 2. Transferring data
- 3. Deinitializing an endpoint

## 3.3.4.1 Initializing an Endpoint

Initializing an endpoint for connection-oriented and connectionless applications is the same. See Section 3.3.3.1 for information on how to initialize an endpoint for a CLTS application.

## 3.3.4.2 Transferring Data

The data transfer phase of connectionless service consists of the following:

- Sending data to other users
- Receiving data from other users
- Retrieving error information about previously sent data

Note that connectionless service:

- Does not support expedited data
- Reports only the T\_UDERR, T\_DATA, and T\_GODATA events

## **Sending Data**

The t\_sndudata function can execute synchronously or asynchronously. When executing synchronously, t\_sndudata returns control to the user when the transport provider can accept another datagram. In some cases, the function blocks for some time until this occurs. In asynchronous mode, the transport provider refuses to send a new datagram if flow control restrictions exist. The t\_sndudata function returns an error of TFLOW, and you must either try again later or issue the t\_look function to see when the flow control restriction is lifted, which is indicated by the T\_GODATA or T GOEXDATA events.

If you attempt to send a data unit before you activate the endpoint with the t\_bind function, the transport provider discards the data.

To send a data unit, issue the t\_sndudata function with the following syntax:

t sndudata (fd, unitdata);

In the preceding statement:

fd

Specifies the file descriptor of the endpoint through which data is sent.

unitdata

Points to the t unitdata structure.

See t sndudata(3) for more information.

The t\_sndudata function returns a value of 0 upon successful completion. Otherwise, it returns a value of -1, and  $t_{errno}$  is set to one of the values described in Section 3.7. (For multithreaded applications,  $t_{errno}$  is thread specific.)

## **Receiving Data**

When you call the t\_rcvudata function and data is available, t\_rcvudata returns immediately indicating the number of octets received. If data is not available, t\_rcvudata behaves differently depending on the mode of execution, as follows:

• Synchronous mode

The t\_rcvudata function blocks until either a datagram, error, or signal is received. As an alternative to waiting for t\_rcvudata to return, you can issue the t\_look function periodically for the T\_GODATA event, and then issue t\_rcvudata to receive the data.

• Asynchronous mode

The t\_rcvudata function returns immediately with an error. You then must either retry the function periodically or poll for incoming data with the t look function.

To receive data, issue the t revudata function with the following syntax:

t\_rcvudata (fd,unitdata,flags);

In the preceding statement:

fd

Specifies the file descriptor of the endpoint through which data is received.

unitdata

Points to the data to be sent, which consists of the following fields:

- addr Returns the address of the sender.
- opt Returns any protocol-specific options that apply to the data.

• udata Returns the data received.

flags

Indicates whether a complete data unit was received (no flag) or a portion of a data unit was received (T\_MORE flag).

See t\_rcvudata(3) for more information.

The t\_rcvudata function returns a value of 0 upon successful completion. Otherwise, it returns a value of -1, and  $t_{errno}$  is set to one of the values described in Section 3.7. (For multithreaded applications,  $t_{errno}$  is thread specific.)

## **Retrieving Error Information**

If you issue the t\_look function and receive the T\_UDERR event, previously sent data has generated an error. To clear the error and consume the T\_UDERR event, you should issue the t\_rcvuderr function. This function also returns information about the data that caused the error and the nature of the error, if you want.

To receive an error indication with information about data, issue the t\_rcvuderr function with the following syntax:

t\_rcvuderr (fd, uderr);

In the preceding statement:

fd

Specifies the file descriptor of the endpoint through which the error report is received.

uderr

Points to the t\_uderr structure, which identifies the error.

See t\_rcvuderr(3) for more information.

The t\_rcvuderr function returns a value of 0 upon successful completion. Otherwise, it returns a value of -1, and  $t_{errno}$  is set to one of the values described in Section 3.7. (For multithreaded applications,  $t_{errno}$  is thread specific.)

## 3.3.4.3 Deinitializing Endpoints

Deinitializing an endpoint for connection-oriented and connectionless applications is the same. See Section 3.3.3.6 for information on how to deinitialize an endpoint for a connectionless application.

# 3.4 Phase-Independent Functions

XTI provides a number of functions that can be issued during any phase of connection-oriented or connectionless service (except the uninitialized state) and do not affect the state of the interface. Table 3-13 lists and briefly describes these functions.

Table 3-13: Phase-Independent Functions

Function	Description
t_getinfo	Returns information about the characteristics of the transport provider associated with the endpoint.
t_getprotaddr <sup>a</sup>	Returns the protocol address.
t_getstate	Returns the current state of the endpoint.
t_strerror <sup>a</sup>	Produces and error message string.
t_sync	Synchronizes the data structures managed by the transport library with information from the transport provider.
t_alloc	Allocates storage for a specified data structure.
t_free	Frees storage for a data structure that was previously allocated by t_alloc.
t_error	Prints a message describing the last error returned by an XTI function. (Optional)
t_look	Returns the current event associated with the endpoint.

### Table notes:

a. This function is supported in XPG4 only.

The t\_getinfo and t\_getstate functions can be useful for retrieving important information. The t\_getinfo function returns the same information about the transport provider as t\_open. It offers the advantage that you can call it during any phase of communication, whereas you can call t\_open only during the initialization phase. If a function returns the TOUTSTATE error to indicate that the endpoint is not in the proper state, you can issue t\_getstate to retrieve the current state and take action appropriate for the state.

The t\_sync function can do the following:

• Synchronize data structures managed by the transport library with information from the underlying transport provider.

• Permit two cooperating processes to synchronize their interaction with a transport provider.

The t\_alloc and t\_free functions are convenient for allocating and freeing memory because you specify the names of the XTI structures rather than information about their size. If you use t\_alloc and t\_free to manage the memory for XTI structures, and the structures change in future releases, you will not need to change your program.

With t\_error you can print a user-supplied message (explanation) plus the contents of t errno to standard output.

Finally, t\_look is an important function for retrieving the current outstanding event associated with the endpoint. Typically, if an XTI function returns TLOOK as an error to indicate a significant asynchronous event has occurred, the transport user follows by issuing the t\_look function to retrieve the event. For more information about events, see Section 3.2.3.

# 3.5 Porting to XTI

This section provides the following:

- Guidelines for writing programs to XTI
- Information about XTI and TLI compatibility
- Information about rewriting sockets applications to use XTI

## 3.5.1 Protocol Independence and Portability

XTI was designed to provide an interface that is independent of the specific transport protocol used. You can write applications that can modify their behavior according to any subset of the XTI functions and facilities supported by each of the underlying transport providers.

Providers do not have to provide all the features of all the XTI functions. Therefore, Application programmers should follow these guidelines when writing XTI applications:

• Use only the functions that are commonly supported features of XTI.

If your application uses features that are not provided by all transport providers, it may not be able to use them with some transport providers or some XTI implementations.

For example, the orderly release facility (the t\_sndrel and t\_rcvrel functions) is not supported by all connection-based transport protocols; in particular it is not supported by ISO protocols. If your application runs in an environment with multiple protocols, make sure it does not use the orderly release facility.

As an alternative to using only the commonly supported features, write your application so that it modifies its behavior according to the subset of XTI functions supported by each transport provider.

Do not assume that logical data boundaries are preserved across a connection.

Some transport providers, such as TCP, do not support the concept of a TSDU, so they ignore the T\_MORE flag when used with the t\_snd, t sndudata, t rcv, and t rcvudata functions.

• Do not exceed the protocol-specific service limits returned on the t\_open and t\_getinfo functions.

Make sure your application retrieves these limits before transferring data and adheres to the limits throughout the communication process.

- Do not rely on options that are protocol-specific.
  - Although the t\_optmgmt function allows an application to access the default protocol options from the transport provider and pass them as an argument to the connection-establishment function, make sure your application avoids examining the options or relying on the existence of certain ones.
- Do not interpret the reason codes associated with the t\_rcvdis function or the error code associated with the t\_rcvuderr function.

  These codes depend on the underlying protocol so, to achieve protocol
  - These codes depend on the underlying protocol so, to achieve protocol independence, make sure your application does not attempt to interpret the codes.
- Perform only XTI operations on the file descriptor returned by the topen function.

If you perform other operations, the results can vary from system to system.

The following sections explain how to port applications from different transport-level programming interfaces to XTI. Specifically, they discuss how to port from the two most common transport-level programming interfaces: Transport Layer Interface (TLI), which many UNIX System V applications use, and the 4.3BSD socket interface, which many Berkeley UNIX applications use.

The information presented in the following sections presumes that you are experienced at programming with TLI or sockets and that you understand fundamental XTI concepts and syntax.

## 3.5.2 XTI and TLI Compatibility

This section discusses issues to consider before you recompile your TLI programs and explains how to recompile them. As a long-term solution, Digital recommends that you use the XTI interface instead of the TLI interface. As more applications and transport providers use XTI, you might find it advantageous to do so as well.

XTI and TLI support the same functions, states, and modes of service. Note that Digital UNIX provides shared library support by default when you link an XTI or TLI application with the XTI or TLI library. For more information on shared library support, see Section 3.2.2.

Before you recompile your TLI program, you should consider your program's current implementation of the following event management: The System V UNIX operating system provides the poll function as a tool for managing events. The Digital UNIX implementation of XTI supports the poll function, so if your application uses it, you can recompile. If your program uses a unique mechanism for managing events, you should port that mechanism to Digital UNIX or change to the polling mechanism provided with Digital UNIX.

Because the Digital UNIX implementation of TLI is compatible at the source level with AT&T TLI, you can recompile your TLI program with the Digital UNIX TLI library using the following steps:

 Make sure the TLI header file is included in your source code: #include <tli/tiuser.h>

2. Recompile your application using the following command syntax:

cc -o name name.c -Itli

If you decide to change your TLI application to an XTI application, be aware of the following minor differences between TLI and XTI.

- In XTI, t\_error is a function of type int that returns an integer value (0 for success and -1 for failure), while in TLI, it is a procedure of type void.
- In XTI, t\_look does not support the T\_ERROR event (as in TLI); it returns -1 and the t\_errno instead.
- For the *oflag* parameter of the t\_open function, the O\_NDELAY value in TLI is known as the O\_NONBLOCK value in XTI.
- XTI opens an endpoint with read-write access because most of its
  functions require read-write access to transport providers. TLI opens with
  read-only, write-only, or read-write access. Specifically, in the t\_open
  function, XTI uses the bitwise inclusive OR of O\_RDWR and
  O\_NONBLOCK as the value of the oflag parameter; TLI uses the

bitwise inclusive OR of O\_NDELAY and either O\_RDONLY, O\_WRONLY, or O\_RDWR. The O\_RDONLY and O\_WRONLY values are not available in XTI; O\_RDWR is the only valid value for access to an endpoint.

- TLI assumes the transport provider has an automatic address generator;
   XTI does not. If the transport provider does not have an automatic address generator,
   XTI can return the proper error message if conflicting requests are issued.
- XTI defines protocol-specific information for the TCP/IP and OSI protocols. The Digital UNIX XTI implementation adds support for protocol-specific options for STREAMS-based protocols; TLI does not provide such information.
- XTI provides additional events to manage flow control, such as T\_GODATA and T\_GOEXDATA; in TLI, you keep sending until successful.
- XTI provides additional error messages to convey more precise error information to applications. All functions that change the state of an endpoint use the TOUTSTATE error to indicate the function was called when the endpoint was in the wrong state. Some XTI functions return the TLOOK error to indicate that an urgent asynchronous event occurred. With TLI, you must call the t\_look function explicitly before the function or set a signal for the TLOOK event, which are less convenient. The TBADQLEN error, returned when there are no queued connection requests, prevents an application from waiting forever after issuing the t\_listen function. See the XTI reference pages for more information on error messages.

To make a TLI application a true XTI application, do the following:

1. Include the XTI header file instead of the TLI header file in your source code:

#include <xti.h>

- 2. Make any changes or extensions to your program resulting from the differences between TLI and XTI.
- 3. Recompile your application using the following command syntax:

cc -o name name.c -lxti

## 3.5.3 Rewriting a Socket Application to Use XTI

This section explains the differences between the socket interface and XTI. It assumes that your applications use the standard 4.3BSD socket interface and does not account for any extensions or changes you have made to the socket interface. See Appendix B for examples of both sockets and XTI servers and clients.

Because it was designed eventually to replace the socket interface, XTI shares many common functions with the socket interface. However, you should be aware of any differences between it and your current socket interface when rewriting an application for use with XTI.

XTI provides 25 functions. Of the 13 socket functions that map onto corresponding XTI functions, 5 have subtle differences. Table 3-14 lists each XTI function, its corresponding socket function (if one exists), and whether the two functions share common semantics. Generally, socket calls pass parameters by value, while most XTI functions pass pointers to structures containing a combination of input and output parameters.

Table 3-14: Comparison of XTI and Socket Functions

XTI Function	Socket Function	<b>Shared Semantics</b>
t_accept	accept	No
t_alloc	_	_
t_bind	bind	No
t_close	close	Yes
t_connect	connect	Yes
t_error	_	_
t_free	_	_
t_getinfo	_	_
t_getstate	_	_
t_listen	listen, accept	Yesa
t_look	select	No
t_open	socket	Yes
t_optmgmt	setsockopt, getsockopt	No
t_rcv	recv	Yes
t_rcvconnect	_	_
t_rcvdis	_	_

Table 3-14: (continued)

XTI Function	Socket Function	Shared Semantics
t_rcvrel	_	_
t_rcvudata	recvfrom	Yes
t_rcvuderr	_	_
t_snd	send	Yes
t_snddis	shutdown	No
t_sndrel	_	_
t_sndudata	sendto	Yes
t_sync	_	_
t_unbind	_	_

#### Table notes:

a. In XTI, the t\_listen function specifies the queue length parameter as well as waiting for the incoming connection. In sockets, the listen function only specifies the queue length parameter.

The XTI functions that do not share all semantics with their socket counterparts have the following differences:

### t\_accept

The t\_accept function takes the user-specified resfd argument and establishes a connection with the remote endpoint. In contrast, the accept call from sockets asks the system to select the file descriptor to which the connection will be established. Additionally, the t\_accept function is issued after a connection indication is received; therefore, it does not block. Conversely, the accept call is issued in anticipation of a connect request and therefore may block until the connect request occurs.

### t bind

XTI can bind one protocol address to many endpoints, while the socket interface permits one address to be bound with only one socket.

## t\_look

The t\_look function returns the current event, which can be one of nine possible events: T\_LISTEN, T\_CONNECT, T\_DATA, T\_EXDATA, T\_DISCONNECT, T\_UDERR, T\_OREREL, T\_GODATA, T\_GOEXDATA. The poll function can be used to monitor incoming events on a transport endpoint. The select call can be used to see if a single descriptor is ready for read or write, or if an exceptional condition is pending.

### t\_snddis

The t\_snddis function initiates an abortive release on an established connection or rejects a connection request. After an XTI program issues the t\_snddis functions it can continue to listen for requests with the t\_listen function or re-establish a connection with the t\_connect function. In sockets, once you shut down a connection with the shutdown and close calls, the system automatically frees all local resources that are allocated for this connection. Therefore, in order to continue to listen for connections or establish a connection, the program needs to reissue the socket and bind calls.

XTI and sockets both use a series of states to control the appropriate sequence of calls, but each uses a different set of states. XTI states and socket states do not share similar semantics. For example, XTI states are mutually exclusive; socket states are not.

Few error messages are common among sockets and XTI. Table 3-15 lists the socket error messages that have comparable XTI error messages.

Table 3-15: Comparison of Socket and XTI Messages

Socket Error	XTI Error	Description
EBADF	TBADF	You specified an invalid file descriptor.
EOPNOTSUPP	TNOTSUPPORT	You issued a function the underlying transport provider does not support.
EADDRINUSE	TADDRBUSY	You specified an address that is already in use.
EACCES	TACCES	You do not have permission to use the specified address.

### Note

XTI and TLI are implemented using STREAMS. You should use the poll function instead of the select call on any STREAMS file descriptors.

## 3.6 Differences Between XPG3 and XPG4

This section provides information on the differences between the XPG3 and XPG4 implementation of XTI.

In earlier versions of Digital UNIX, the XTI implementation conformed to X/Open's XPG3 specification. The current implementation conforms to

SPEC1170's XTI (part of Networking Services' specification) as well as X/Open's XPG4 specification for XTI.

There are some changes in the specification of which you, as a programmer, should be aware. This section outlines these differences and the related programming issues.

Note that the implementation of Digital UNIX converges both XPG3 and XPG4 versions of XTI in a single subset. This section also provides details about the usage of the appropriate level of functionality.

In this manual, the terms SPEC1170 or SPEC1170 XTI are used to refer to the implementation of XTI available in this version of Digital UNIX. The terms XPG3 XTI refer to the implementation of XTI that conforms to X/Open's XPG3 specification. Note that the latter can be available in the current versions of Digital UNIX due to binary compatibility or source migration features.

## 3.6.1 Major Differences

Most of the changes between the two specifications are upwardly compatible, with the exception of the t\_optmgmt function.

The following is a quick summary of the basic changes in the XTI from XPG3 to SPEC1170:

- Optional functions were made mandatory. This does not affect the Digital UNIX implementation of XTI, because Digital UNIX implemented all the optional functions in its XPG3 version of XTI.
- Many aspects of the XPG3 specification were clarified, which makes XTI applications more portable.
- Some new error codes were added, ensuring better programmatic behavior.
- Options and management structures were revised to provide more control over various aspects of communications.

The changes to the t\_optmgmt function are extensive and incompatible with the XPG3 specification. In general, an application that uses the XPG3 implementation of the t\_optmgmt function cannot use the t\_optmgmt function on a system running the XPG4 specification, without making some modifications to the source.

## 3.6.2 Source Code Migration

If you have an application that was developed for XPG3 XTI, you have the following choices to support it under Digital UNIX:

- Use the older binaries of the application; see Section 3.6.3.
- Recompile the unaltered sources.
- Make changes to the sources to comply with SPEC1170 XTI.

Which option you choose will depend on your situation. The following sections describe these conditions in details.

## 3.6.2.1 Use the Older Binaries of your Application

This choice is appropriate if the sources and features of your application are not going to change. It is useful to provide continued coverage by ensuring that older releases of your products are still functional.

#### 3.6.2.2 Unaltered Sources

This situation arises from minor changes due to correcting minor problems. Therefore, there are no changes to the structure or features or the application. In this case, you might want to compile the sources in the same manner as XPG3 development environment. In that case, compile your source code with the -DXPG3 compiler switch. This ensures that the headers automatically define the older features for you.

## 3.6.2.3 SPEC1170 Compliant Application

If you need to use the new features supported by SPEC1170 XTI, you will have to make changes in your source code. You cannot combine the features from the XPG3 and SPEC1170 XTI. Therefore, if you have large applications consisting of multiple files, you will need to recompile all files with the new features, rather than just the few you might have changed.

You need to compile your source code with the <code>-DXOPEN\_SOURCE</code> compiler switch. Additionally, you must ensure that the names of the transport protocols (as provided through the streams device special files as in <code>/dev/streams/xtiso/tcp</code>) are updated to reflect the naming convention used in SPEC1170 XTI. For example, the names for TCP and UDP are <code>/dev/streams/xtiso/tcp+</code> and <code>/dev/streams/xtiso/udp+</code>. Check the reference manual for the names for the other protocols.

## 3.6.3 Binary Compatibility

Application binaries developed with XPG3 XTI will run on systems running the current version of Digital UNIX. However, there are certain conditions of which you should be aware.

Under unusual circumstances, the errors in XPG3 programs may have been masked due to the way in which the programs or libraries were compiled and

linked. It is feasible that the new implementation is able to flag such conditions as errors. Since the error manifested is a programming error in the application, you will have to correct it. The common programming errors that may cause these errors are pointer overruns and uninitialized variables.

Another issue to consider is the availability of SPEC1170 features through STREAMS special files. This is significant if your application accepts command line input for the specifying transport protocol or imports the protocol names from some configuration files. Since the system configured with XTI will have the file names for SPEC1170-compliant protocols as well, it is important to warn users and administrators that those special names should not be used with applications running with binary-compatibility mode. The results of such an action are undefined.

If you are planning to run an old applications without recompiling them, check them for binary compatibility to avoid these problems.

## 3.6.4 Packaging

Systems running the current version of Digital UNIX and configured to run XTI support both XPG3 and SPEC1170-compliant functionality. You cannot run the XPG3 and SPEC1170 functionality separately. Therefore, you only need to ensure that XTI subsystem is configured.

## 3.6.5 Interoperability

You can use the XPG3 and SPEC1170 versions of XTI on the same network. If you are using compatible versions of your application, then the operation should be transparent to users.

It is possible to of convert your application in simple steps, so that you have some pieces that are XPG3 XTI compatible and some pieces that are SPEC1170 compatible. The only thing you have to ensure is that application-level protocol remains the same. Apart from that there will be no issue for interoperability of these components. Therefore, if you have client and server components of an application, you can choose to upgrade the server component for SPEC1170 compliance, while the client component is still operational in binary compatibility mode. Later on, once the server functionality is updated satisfactorily, you can choose to update the client software.

## 3.6.6 Using XTI Options

This section provides information on using XTI options in XPG4 and XPG3.

## 3.6.6.1 Using XTI Options in XPG4

This section provides the following information on using XTI options:

- General information on using options
- Format of options
- Elements of negotiation
- Option management of transport endpoint

#### **General Information**

The following functions contain an *opt* argument of the type struct netbuf as an input or output parameter. This argument is used to convey options between the transport user and the transport provider:

- t\_accept
- t connect
- t\_listen
- t\_optmgmt
- t rcvconnect
- t\_rcvudata
- t rcvuderr
- t\_sndudata

There is no general definition about the possible contents of options. There are general XTI options and those that are specific for each transport provider. Some options allow you to tailor your communication needs; for instance, by asking for high throughput or low delay. Others allow the fine-tuning of the protocol behavior so that communication with unusual characteristics can be handled more effectively. Other options are for debugging purposes.

All options have default values. Their values have meaning to and are defined by the protocol level in which they apply. However, their values can be negotiated by a transport user. This includes the simple case where the transport user can enforce its use. Often, the transport provider or even the remote transport user can have the right to negotiate a value of lesser quality than the proposed one, that is, a delay can become longer, or a throughput may become lower.

It is useful to differentiate between options that are association-related and those that are not. (Association-related means a pair of communication transport users.) Association-related options are intimately related to the particular transport connection or datagram transmission. If the calling user

specifies such an option, some ancillary information is transferred across the network in most cases. The interpretation and further processing of this information is protocol-dependent. For instance, in an ISO connection-oriented communication, the calling user can specify quality-of-service parameters on connection establishment. These are first processed and possibly lowered by the local transport provider, then sent to the remote transport provider that may degrade them again, and finally conveyed to the called user that makes the final selection and transmits the selected values back to the caller.

Options that are not association-related do not contain information destined for the remote transport user. Some have purely local relevance; for example, an option that enables debugging. Others influence the transmission; for instance, the option that sets the IP time-to-live field or TCP\_NODELAY. (See the xti\_internet(7) reference page.) Local options are negotiated solely between the transport user and the local transport provider. The distinction between these two categories of options is visible in XTI through the following relationship: on output, the t\_listen and t\_rcvudata functions return association-related options only. The t\_rcvconnect and t\_rcvuderr functions may return options of both categories. On input, options of both categories may be specified with the t\_accept and t\_sndudata functions. The t\_connect and t\_optmgmt functions can process and return both categories of options.

The transport provider has a default value for each option it supports. These defaults are sufficient for the majority of communication relations. Therefore, a transport user should only request options actually needed to perform the task and leave all others at their default value.

This section describes the general framework for the use of options. This framework is obligatory for transport providers. The t\_optmgmt reference page provides information on general XTI options. The xti\_internet reference page provides information on the specific options that are legal with the TCP and UDP transport providers.

### **Format of Options**

Options are conveyed through an *opt* argument of struct netbuf. Each option in the buffer specified is of the form struct t\_opthdr possibly followed by an option value.

A transport provider embodies a stack of protocols. The <code>level</code> field of <code>struct t\_opthdr</code> identifies the XTI level or a protocol of the transport provider as TCP or ISO 8073:1986. The <code>name</code> field identifies the option within the level and the <code>len</code> field contains the total length; that is the length of the option header <code>t\_ophdr</code> plus the length of the option value. The <code>status</code> field is used by the XTI level or the transport provider to indicate success or failure of a negotiation.

Several options can be concatenated; however, The transport user has to ensure that each option starts at a long-word boundary. The macro OPT\_NEXTHDR (pbuf, buflen, poptons) can be used for that purpose. The parameter pbuf denotes a pointer to an option buffer opt.buf and buflen is its length. The parameter poption points to the current options in the option buffer. OPT\_NEXTHDR returns a pointer to the position of the next option or returns a null pointer if the option buffer is exhausted. The macro is helpful for writing and reading the option list.

## **Elements of Negotiation**

This section describes the general rules governing the passing and retrieving of options and the error conditions that can occur. Unless explicitly restricted, these rules apply to all functions that allow the exchange of options.

## **Multiple Options and Options Levels**

When multiple options are specified in an option buffer on input, different rules apply to the levels that may be specified, depending on the function call. Multiple options specified on input to t\_optmgmt must address the same option level. Options specified on input to t\_connect, t\_accept, and t sndudata can address different levels.

## **Illegal Options**

Only legal options can be negotiated; illegal options can cause failure. An option is illegal if the following applies:

- The length specified in the t\_opthdr.len parameter exceeds the remaining size of the option buffer (counted from the beginning of the option).
- The option value is illegal. The legal values are defined for each option. See the t\_optmgmt(3) and xti\_internet(7) reference pages.

If and illegal option is passed to XTI, the following will happen:

- A call to the toptmgmt function fails with a TBADOPT error.
- The t\_accept or t\_connect functions fail with a TBADOPT error or the connection establishment aborts, depending on the implementation and the time the illegal option is detected. If the connection aborts, a T\_DISCONNECT event occurs and a synchronous call to t\_connect fails with a TLOOK error. It depends on timing and implementation conditions whether a t\_accept function can still succeed or fail with a TLOOK error in that case.

 A call to the t\_sndudata function either fails with a TBADOPT error or it successfully returns; but a T\_UDERR event occurs to indicate that the datagram was not sent.

If the transport user passes multiple options in one call and one of them is illegal, the call fails as described previously. It is, however, possible that some or even all of the submitted legal options were successfully negotiated. The transport user can check the current status by a call to the t\_optmgmt function with the T\_CURRENT flag set. See the t\_optmgmt(3) and xti\_internet(7) reference pages.

Specifying an option level unknown to or not supported by the protocol selected by the option level does not cause failure. The option is discarded in calls to the t\_connect, t\_accept, or t\_sndudata functions. The t\_opmgmt function returns T\_NOTSULPORT in the <code>level</code> field of the option.

## **Initiating an Option Negotiation**

A transport user initiates an option negotiation when calling the t\_connect, t\_sndudata, or t\_optmgmt functions with the T\_NEGOTIATE flag set.

The negotiation rules for these functions depend on whether an option request is an absolute requirement. This is explicitly defined for each option. See the t\_optmgmt(3) and xti\_internet(7) reference pages. In the case of an ISO transport provider, for example, the option that requests use of expedited data is not an absolute requirement. On the other hand, the option that requests protection could be an absolute requirement.

### **Note**

The term **absolute requirement** originates from the quality-of-service parameters in the ISO 8072:1986 specification. Its use is extended here to all options.

If the proposed option value is an absolute requirement, there are three possible outcomes:

- The negotiated value is the same as the proposed one. When the result of the negotiation is retrieved, the status field in t\_opthdr is set to T SUCCESS.
- The negotiation is rejected if the option is supported but the proposed value cannot be negotiated. This leads to the following:
  - The t\_optmgmt function successfully returns; but the returned option has its status field set to T\_FAILURE.

- Any attempt to establish a connection aborts; a T\_DISCONNECT event occurs and a synchronous call to the t\_connect function fails with a TLOOK error.
- The t\_sndudata function fails with a TLOOK error or successfully returns; but a T\_UDERR event occurs to indicate that the datagram was not sent.

If multiple options are submitted in one call and one of them is rejected, XTI behaves as just described. Although the connection establishment or the datagram transmission fails, options successfully negotiated before some option was rejected retain their negotiated values. There is no roll-back mechanism. See the Option Management of a Transport Endpoint section for more information.

The t\_optmgmt function attempts to negotiate each option. The *status* fields of the returned options indicate success (T\_SUCCESS) or failure (T\_FAILURE).

• If the local transport provider does not support the option at all, the t\_optmgmt function reports T\_NOTSULPORT in the status field. The t\_connect and t\_sndudata functions ignore this option.

If the proposed option value is not an absolute requirement, the following outcomes are possible:

- The negotiated value is of equal or lesser quality than the proposed one; for example, a delay may become longer.
  - When the result of the negotiation is retrieved, the *status* field in t\_opthdr is set to T\_SUCCESS if the negotiated value equals the proposed one; otherwise, it is set to T PARTSUCCESS.
- If the local transport provider does not support the option at all, t\_optmgmt reports T\_NOTSULPORT in the status field. The t\_connect and t\_sndudata functions ignore this option.

Unsupported options do not cause functions to fail or a connection to abort, since different vendors possibly implement different subsets of options. Furthermore, future enhancements of XTI might encompass additional options that are unknown to earlier implementations of transport providers. The decision whether or not the missing support of an option is acceptable for the communication is left to the transport user.

The transport provider does not check for multiple occurrences of the same options, possibly with different option values. It simply processes the options in the option buffer sequentially. However, the user should not make any assumption about the order of processing.

Not all options are independent of one another. A requested option value might conflict with the value of another option that was specified in the same

call or is currently effective. See the Option Management of a Transport Endpoint section for more information. These conflicts may not be detected at once, but they might later lead to unpredictable results. If detected at negotiation time, these conflicts are resolved within the rules stated above. The outcomes may thus be quite different and depend on whether absolute or nonabsolute requests are involved in the conflict.

Conflicts are usually detected at the time a connection is established or a datagram is sent. If options are negotiated with the t\_optmgmt function, conflicts are usually not detected at this time, since independent processing of the requested options must allow for temporal inconsistencies.

When called, the t\_connect, and t\_sndudata functions initiate a negotiation of all association-related options according to the rules of this section. Options not explicitly specified in the function calls themselves are taken from an internal option buffer that contains the values of a previous negotiation. See the Option Management of a Transport Endpoint section for more information.

## **Responding to a Negotiation Proposal**

In connection-oriented communication, some protocols give the peer transport users the opportunity to negotiate characteristics of the transport connection to be established. These characteristics are association-related options. With the connect indication, the called user receives (through the t\_listen function) a proposal about the option values that should be effective for this connection. The called user can accept this proposal or weaken it by choosing values of lower quality; for example, longer delays than proposed. The called user can, of course, refuse the connection establishment altogether.

The called user responds to a negotiation proposal using the t\_accept function. If the called transport user tries to negotiate an option of higher quality than proposed, the outcome depends on the protocol to which that option applies. Some protocols may reject the option, some protocols take other appropriate action described in protocol-specific reference pages. If an option is rejected, the following error occurs:

The connection fails; a T\_DISCONNECT event occurs. In that case, whether a t\_accept function can still succeed or fail with a TLOOK error depends on timing and implementation conditions.

If multiple options are submitted with the t\_accept function and one of them is rejected, the connection fails as described previously. Options that could be successfully negotiated before the erroneous option was processed retain their negotiated value. There is no rollback mechanism. See the Option Management of a Transport Endpoint section for more information.

The response options can either be specified with the t\_accept call or can be preset for the responding endpoint (not the listening endpoint) resfd in a t\_optmgmt call (action T\_NEGOTIATE) prior to the t\_accept call. (See the Option Management of a Transport Endpoint section for more information.) Note that the response to a negotiation proposal is activated when the t\_accept function is called. A t\_optmgmt function call with erroneous option values as described previously will succeed; the connection aborts at the time the t\_accept function is called.

The connection also fails if the selected option values lead to contradictions.

The t\_accept function does not check for multiple specification of an option. (See the Initiating an Option Negotiation section.) Unsupported options are ignored.

## **Retrieving Information About Options**

This section describes how a transport user can retrieve information about options.

A transport user must be able to:

- Know the result of a negotiation; for example, at the end of a connection establishment.
- Know the proposed option values under negotiation during connection establishment.
- Retrieve option values sent by the remote transport user for notification only; for example, IP options.
- Check option values currently in effect for the transport endpoint.

To this end, the following function take an output argument opt of the struct netbuf:

- t connect
- t\_listen
- t\_optmgmt
- t\_rcvconnect
- t\_rcvudata
- t rcvuderr

The transport user has to supply a buffer to which the options will be written; the opt.buf parameter must point to this buffer and the opt.maxlen parameter must contain the buffer's size. The transport user can set the opt.maxlen parameter to zero to indicate that no options are to be retrieved.

Which options are returned depend on the function call involved:

• t\_connect in synchronous mode and t\_rcvconnect

The functions return the values of all association-related options that were received with the connection response and the negotiated values of those nonassociation-related options that had been specified on input. However, options specified on input in the t\_connect call that are not supported or refer to an unknown option level are discarded and not returned on output.

The *status* field of each option returned with the t\_connect or t\_rcvconnect function indicates if the proposed value (T\_SUCCESS) or a degraded value (T\_PARTSUCCESS) has been negotiated. The *status* field of received ancillary information (for example, IP options) that is not subject to negotiation is always set to T\_SUCCESS.

### • t\_listen

The received association-related options are related to the incoming connection (identified by the sequence number), not to the listening endpoint. (However, the option values currently in effect for the listening endpoint can affect the values retrieved by the t\_listen function, since the transport provider might also be involved in the negotiation process.) Therefore, if the same options are specified in a call to the t\_optmgmt function with action T\_CURRENT, they will usually not return the same values.

The number of received options may vary for subsequent connect indications, since many association-related options are only transmitted on explicit demand by the calling user; for example, IP options or ISO 8072:1986 throughput. It is even possible that no options at all are returned.

The status field is irrelevant.

### • t rcvudata

The received association-related options are related to the incoming datagram, not to the transport endpoint fd. Therefore, if the same options are specified in a call to the t\_optmgmt function with action T\_CURRENT, the t\_optmgmt function will usually not return the same values.

The number of options received may vary from call to call.

The status field is irrelevant.

#### • t rcvuderr

The returned options are related to the options input of the previous t sndudata call that produced the error. Which options are returned

and which values they have depend on the specific error condition. The status field is irrelevant.

### • t\_optmgmt

This call can process and return both categories of options. It acts on options related to the specified transport endpoint, not on options related to a connect indication or an incoming datagram. For more information, see the t\_optmgmt(3) reference page.

## **Privileged and Read-Only Options**

Only privileged users can request privileged options, or option values. The meaning of privilege is hereby implementation-defined.

Read-only options serve for information purposes only. The transport user may be allowed to read the option value but not to change it. For instance, to select the value of a protocol timer or the maximum length of a protocol data unit may be too subtle to leave to the transport user, though the knowledge about this value might be of some interest. An option might be read-only for all users or solely for nonprivileged users. A privileged option might be inaccessible or read-only for nonprivileged users.

An option might be negotiable in some XTI states and read-only in other XTI states. For instance, the ISO quality-of-service options are negotiable in the T\_IDLE and T\_INCON states, and read-only in all other states (except T\_UNINIT).

If a transport user requests negotiation of a read-only option, or a nonprivileged user requests illegal access to a privileged option, the following outcomes are possible:

- The t\_optmgmt function successfully returns, but the returned option
  has its status field set to T\_NOTSULPORT if a privileged option was
  requested illegally, and to T\_READONLY if modification of a read-only
  option was requested.
- If negotiation of a read-only option is requested, the t\_accept or t\_connect functions fail with TACCES or the connection establishment aborts and a T\_DISCONNECT event occurs. If the connection aborts, a synchronous call to t\_connect fails with TLOOK. If a privileged option is illegally requested, the option is quietly ignored. A nonprivileged user is not able to select an option that is privileged or unsupported. Timing and implementation conditions determine whether a t accept call still succeeds or fails with TLOOK.
- If negotiation of a read-only option is requested, the t\_sndudata function may return TLOOK or successfully return, but a T\_UDERR event occurs to indicate that the datagram was not sent. If a privileged option is illegally requested, the option is quietly ignored. A

nonprivileged user cannot select an option that is privileged or unsupported.

If multiple options are submitted to the t\_connect, t\_accept, or t\_sndudata functions and a read-only option is rejected, the connection or the datagram transmission fails as described. Options that could be successfully negotiated before the erroneous option was processed retain their negotiated values. There is no rollback mechanism. See the Option Management of a Transport Endpoint section for more information.

## **Option Management of a Transport Endpoint**

This section describes how option management works during the lifetime of a transport endpoint.

Each transport endpoint is (logically) associated with an internal option buffer. When a transport endpoint is created, this buffer is filled with a system default value for each supported option. Depending on the option, the default may be OPTION ENABLED, OPTION DISABLED, or denote a time span, and so on. These default settings are appropriate for most uses. Whenever an option value is modified in the course of an option negotiation, the modified value is written to this buffer and overwrites the previous one. At any time, the buffer contains all option values that are currently effective for this transport endpoint.

The current value of an option can be retrieved at any time by calling the t\_optmgmt function with the T\_CURRENT flag set. Calling the t\_optmgmt function with the T\_DEFAULT flag set yields the system default for the specified option.

A transport user can negotiate new option values by calling the t\_optmgmt function with the T\_NEGOTIATE flag set. The negotiation follows the rules described in the Elements of Negotiation section.

Some options may be modified only in specific XTI states and are read-only in other XTI states. Many association-related options, for instance, may not be changed in the T\_DATAXFER state, and an attempt to do so fails; see the Privileged and Read-Only Options section. The legal states for each option are specified with its definition.

As usual, association-related options take effect at the time a connection is established or a datagram is transmitted. This is the case if they contain information that is transmitted across the network or determine specific transmission characteristics. If such an option is modified by a call to the t\_optmgmt function, the transport provider checks whether the option is supported and negotiates a value according to its current knowledge. This value is written to the internal option buffer.

The final negotiation takes place if the connection is established or the datagram is transmitted. This can result in a degradation of the option value or even in a negotiation failure. The negotiated values are written to the internal option buffer.

Some options can be change in the T\_DATAXFER state; for example, those specifying buffer sizes. Such changes might affect the transmission characteristics and lead to unexpected side effects; for example, data loss if a buffer size was shortened.

The transport user can explicitly specify both categories of options on input when calling the t\_connect, t\_accept, or t\_sndudata functions. The options are at first locally negotiated option by option and the resulting values written to the internal option buffer. The modified option buffer is then used if a further negotiation step across the network is required; for example, in connection-oriented ISO communication. The newly negotiated values are then written to the internal option buffer.

At any stage, a negotiation failure can cause the transmission to abort. If a transmission aborts, the option buffer preserves the content it had at the time the failure occurred. Options that could be negotiated before the error occurred are written back to the option buffer, whether the XTI call fails or succeeds.

It is up to the transport user to decide which option it explicitly specifies on input when calling the t\_connect, t\_accept, or t\_sndudata functions. The transport user need not pass options at all by setting the *len* field of the function's input *opt* argument to zero (0). The current content of the internal option buffer is then used for negotiation without prior modification.

The negotiation procedure for options at the time of a t\_connect, t\_accept, or t\_sndudata call always obeys the rules in the Initiating an Option Negotiation section whether the options were explicitly specified during the call or implicitly taken from the internal option buffer.

The transport user should not make assumptions about the order in which options are processed during negotiation.

A value in the option buffer is only modified as a result of a successful negotiation of this option. It is, in particular, not changed by a connection release. There is no history mechanism that would restore the buffer state existing prior to the connection establishment of the datagram transmission. The transport user must be aware that a connection establishment or a datagram transmission may change the internal option buffer, even if each option was originally initialized to its default value.

## The Option Value T UNSPEC

Some options may not always have a fully specified value. An ISO transport provider, for instance, that supports several protocol classes might not have a preselected preferred class before a connection establishment is initiated. At the time of the connection request, the transport provider may conclude from the destination address, quality-of-service parameters, and other locally available information which preferred class it should use. A transport user asking for the default value of the preferred class option in the T\_IDLE state would get the value T\_UNSPEC. This value indicates that the transport provider did not yet select a value. The transport user could negotiate another value as the preferred class; for example, T\_CLASS2. The transport provider would then be forced to initiate a connect request with class 2 as the preferred class.

An XTI implementation may also return the T\_UNSPEC value if it currently cannot access the option value. This can happen in the T\_UNBND state in systems where the protocol stacks reside on separate controller cards and not in the host. The implementation may never return T\_UNSPEC if the option is not supported at all.

If T\_UNSPEC is a legal value for a specific option, it can be used on input, as well. It is used to indicate that it is left to the provider to choose an appropriate value. This is especially useful in complex options as ISO throughput, where the option value has an internal structure. The transport user can leave some fields unspecified by selecting this value. If the user proposes T\_UNSPEC, the transport provider is free to select an appropriate value. This might be the default value, some other explicit value, or T\_UNSPEC.

For each option, it is specified whether T\_UNSPEC is a legal value for negotiation purposes.

## The info Argument

The t\_open and t\_getinfo functions return values representing characteristics of the transport provider in the *info* argument. The value of *info->options* is used by the t\_alloc function to allocate storage for an option buffer to be used in an XTI call. The value is sufficient for all uses

In general, *info->options* also includes the size of privileged options; even if these are not read-only for nonprivileged users. Alternatively, an implementation can choose to return different values in *info->options* for privileged and nonprivileged users.

The values in info->etsdu, info->connect, and info->discon possibly diminish as soon as the T DATAXFER state is entered. Calling the

t\_optmgmt function does not influence these values. For more information, see the t\_optmgmt(3) reference page.

## **Portability Issues**

An application programmer who writes XTI programs has the following portability issues across the following:

- Protocol profiles
- Different system platforms

Options are intrinsically coupled with a definite protocol or protocol profile. Therefore, explicit use of options degrades portability across protocol profiles.

Different vendors might offer transport providers different option support. This is due to different implementation and product policies. The lists of options on the t\_optmgmt(3) reference page and in the protocol-specific reference pages are maximal sets, but do not necessarily reflect common implementation practice. Vendors implement subsets that suit their needs. Therefore, making careless use of options endangers portability across different system platforms.

Every implementation of a protocol profile accessible by XTI can be used with the default values of options. Applications can thus be written that do not care about options at all.

An application program that processes options retrieved from an XTI function should discard options it does not know to lessen its dependence from different system platforms and future XTI releases with possibly increased option support.

### 3.6.6.2 Negotiating Protocol Options in XPG3

The Digital UNIX XPG3 implementation of XTI provides an optional function, t\_optmgmt, for retrieving, verifying, and negotiating protocol options with transport providers. After you create an endpoint with t\_open and bind an address to it, you can verify or negotiate options with the transport provider. To do so, issue the t\_optmgmt function, with the following syntax:

t\_optmgmt (fd,req,ret);

In the preceding statement:

fd

Identifies the file descriptor for the endpoint, which is returned by the topen function.

req

Points to a t\_optmgmt structure that sends protocol options to the transport provider and requests actions of the transport provider.

ret

Points to a t\_optmgmt structure that returns the valid protocol options and the actions taken by the transport provider.

Both the req and ret arguments point to a t\_optmgmt structure.

### Note

Although other transport providers may support the t\_optmgmt function, the Digital UNIX TCP transport provider does not. See the transport provider documentation for information about option management.

See t\_optmgmt(3) for more information.

The t\_optmgmt function returns a value of 0 upon successful completion; otherwise, it returns a value of -1, and  $t_{errno}$  is set to one of the values described in Section 3.7. (For multithreaded applications,  $t_{errno}$  is thread specific.)

## 3.7 XTI Errors

XTI returns library errors and system errors. When an XTI function encounters an error, it returns a value of -1, and can do one of the following:

- Check the external variable *t\_errno* to get the specific error. (For multithreaded applications, *t\_errno* is thread specific.)
- Call the t\_error function to print the text of the message associated with the error stored in *t errno*.
- Check the state of the transport endpoint with the t\_getstate function. Some errors change the state of the endpoint.

#### Note

Since a successful call to an XTI function does not clear the contents of  $t\_errno$ , check  $t\_errno$  only after an error occurs.

The <xti.h> header file defines the  $t_errno$  variable as a macro as follows:

```
#define t_errno(*_t_errno())
```

For more information on errors, see the individual XTI reference pages.

# 3.8 Configuring XTI Transport Providers

Use the xtiso kernel configuration option to configure XTI transport providers. You can configure the xtiso option into your system at installation time or you can add it to your system using the doconfig command. See the *Installation Guide*.

You can use the doconfig command in one of the following ways:

- Use the doconfig command without options if you have not customized your kernel. Without options the doconfig command creates a new kernel configuration file for your system.
- Use the doconfig -c command if you have customized your kernel and you do not want to recustomize it. The doconfig -c command allows you to add information to the existing kernel configuration file.

To use the doconfig command without any options, do the following:

- 1. Enter the /usr/sbin/doconfig command at the superuser prompt (#).
- 2. Enter a name for the kernel configuration file. It should be the name of your system in all uppercase letters, and will probably be the default provided in square brackets ([]). For example:

```
Enter a name for the kernel configuration file. [HOST1]:
RETURN
```

3. Enter y when the system asks whether you want to replace the system configuration file. For example:

```
A configuration file with the name 'HOST1' already exists. Do you want to replace it? (y/n) [n]: {\bf y} Saving /sys/conf/HOST1 as /sys/conf/HOST1.bck
```

```
*** KERNEL CONFIGURATION AND BUILD PROCEDURE ***
```

4. Select the X/Open Transport Interface (XTISO, TIMOD, TIRDWR) option from the Kernel Option Selection menu. Confirm your choice at the prompt.

#### For example:

\*\*\* KERNEL OPTION SELECTION \*\*\*

```
Selection Kernel Option
______
            1 System V Devices
2 NTP V3 Kernel Phase Lock Loop (NTP_TIME)
3 Kernel Breakpoint Debugger (KDEBUG)
4 Packetfilter driver (PACKETFILTER)
5 Point-to-Point Protocol (PPP)
6 STREAMS pckt module (PCKT)
7 X/Open Transport Interface (XTISO, TIMOD, TIRDWR)
8 File on File File System (FFM)
9 ISO 9660 Compact Disc File System (CDFS)
10 Audit Subsystem
11 ACL Subsystem
12 Logical Storage Manager (LSM)
13 Advanced File System (ADVFS)
              13
                           Advanced File System (ADVFS)
             14
                           All of the above
              15
                           None of the above
             16
                          Help
Enter the selection number for each kernel option you want.
For example, 1 3 [15]: 7
Enter the selection number for each kernel option you want.
For example, 1 3 : 7
```

You selected the following kernel options:

X/Open Transport Interface (XTISO, TIMOD, TIRDWR) Is that correct? (y/n) [y]: y

Configuration file complete.

5. Enter n when the doconfig command asks whether you want to edit the configuration file.

The doconfig command then creates device special files, indicates where a log of the files it created is located, and builds the new kernel. After the new kernel is built, you must move it from the directory where doconfig places it to the root directory ( /) and reboot your system.

When you reboot, the strsetup -i command runs automatically, creating the device special files for any new STREAMS modules.

6. Enter the strsetup -c command to verify that the device is configured properly.

The following example shows the output from the strsetup -c command:

#### # /usr/sbin/strsetup -c

STREAMS Configuration Information...Fri Nov 3 14:23:36 1995

Name	Type	Major	Module ID
		32	0
clone	a 2		
dlb	device	52	5010
kinfo	device	53	5020
log	device	54	44
nuls	device	55	5001
echo	device	56	5000
sad	device	57	45
pipe	device	58	5304
xtisoUDP	device	59	5010
xtisoTCP	device	60	5010
xtisoUDP+	device	61	5010
xtisoTCP+	device	62	5010
ptm	device	63	7609
pts	device	6	7608
bba	device	64	24880
lat	device	5	5
pppif	module		6002
pppasync	module		6000
pppcomp	module		6001
bufcall	module		0
null	module		5002
pass	module		5003
errm	module		5003
ptem	module		5003
spass	module		5007
rspass	module		5008
pipemod	module		5303
timod	module		5006
tirdwr	module		0
ldtty	module		7701
-			

Configured devices = 15, modules = 14

To use the doconfig -c command to add the XTISO option to the kernel configuration file, do the following:

1. Enter the doconfig -c *HOSTNAME* command from the superuser prompt (#). *HOSTNAME* is the name of your system in all uppercase

letters. For example, for a system called host1 you would enter: # doconfig -c HOST1

2. Add XTISO to the options section of the kernel configuration file.

Enter y at the prompt to edit the kernel configuration file. The doconfig command allows you to edit the configuration file with the ed editor. For information about using the ed editor, see ed(1).

The following ed editing session shows how to add the XTISO option to the kernel configuration file for host1. The number of the line after which you append the new line can differ between kernel configuration files:

```
*** KERNEL CONFIGURATION AND BUILD PROCEDURE ***

Saving /sys/conf/HOST1 as /sys/conf/HOST1.bck

Do you want to edit the configuration file? (y/n) [n]: y

Using ed to edit the configuration file. Press return when ready, or type 'quit' to skip the editing session: 2153

48a options XTISO

. 1,$w
2185
q

**** PERFORMING KERNEL BUILD ***
```

3. After the new kernel is built you must move it from the directory where doconfig places it to the root directory ( / ) and reboot your system.

When you reboot, the strsetup -i command is run automatically, creating the device special files for any new STREAMS modules.

4. Run the strsetup -c command to verify that the device is configured properly.

The following example shows the output from the strsetup -c command:

#### # /usr/sbin/strsetup -c

STREAMS Configuration Information...Fri Nov 3 14:23:36 1995

Name	Type	Major	Module ID
clone		32	0
dlb	device	52	5010
kinfo	device	53	5020
log	device	54	44
nuls	device	55	5001
echo	device	56	5000
sad	device	57	45
pipe	device	58	5304
xtisoUDP	device	59	5010
xtisoTCP	device	60	5010
xtisoUDP+	device	61	5010
xtisoTCP+	device	62	5010
ptm	device	63	7609
pts	device	6	7608
bba	device	64	24880
lat	device	5	5
pppif	module		6002
pppasync	module		6000
pppcomp	module		6001
bufcall	module		0
null	module		5002
pass	module		5003
errm	module		5003
ptem	module		5003
spass	module		5007
rspass	module		5008
pipemod	module		5303
timod	module		5006
tirdwr	module		0
ldtty	module		7701

Configured devices = 15, modules = 14

For detailed information on reconfiguring your kernel or the doconfig command see the *System Administration* manual.

# Sockets 4

The Digital UNIX sockets programming interface supports the XPG4 standard and the Berkeley Software Distribution (BSD) socket programming interface.

In Digital UNIX, sockets provide an interface to the Internet Protocol suite (TCP/IP) and to the UNIX domain for interprocess communication on the same system. However, you can use sockets to build network-based applications that are independent of the underlying networking protocols and hardware.

To use the XPG4 standard implementation in your program, you must compile your program using the c89 compiler command. The examples in this chapter are based on the XPG4 standard. See Section 4.4 for information on the differences between the XPG4 and the BSD interfaces.

This chapter contains the following information:

- Overview of the sockets framework
- Description of the application interface to sockets
- Information on how to use sockets
- Information on the BSD socket interfaces
- Explanation of common socket error messages
- Information about advanced topics

Figure 4-1 highlights the sockets framework and shows its relationship to the rest of the network programming environment:

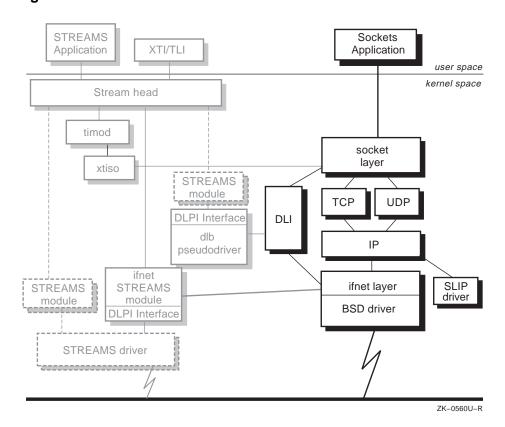


Figure 4-1: The Sockets Framework

4.1 Overview of the Sockets Framework

The sockets framework consists of:

- A set of abstractions, such as communication domains and socket types, that defines socket communication properties
- A programming interface, or set of system and library calls, used by application programs to access the socket framework
- Kernel resources, including networking protocols, that application programs access using system and library calls
  - Digital UNIX implements the Internet Protocol suite and UNIX domain using sockets to achieve interprocess communication. It also implements BSD-based device drivers that are accessed using sockets system calls.

# 4.1.1 Communication Properties of Sockets

This section describes the abstractions and definitions that underlie sockets communication properties.

#### 4.1.1.1 Socket Abstraction

Sockets function as endpoints of communication. A single socket is one endpoint; a pair of sockets constitutes a two-way communication channel that enables unrelated processes to exchange data locally and over networks.

Application programs request the operating system to create a socket when one is needed. The operating system returns a socket descriptor that the program uses to reference the newly created socket for further operations.

Sockets have the following characteristics:

- Exist only as long as some process holds a descriptor referencing it.
- Are referenced by descriptors and have qualities similar to those of a character special device. Read, write, and select operations are performed on sockets by using the appropriate system calls.
- Can be created in pairs or given names and used to rendezvous with other sockets in a communications domain, accepting connections from these sockets or exchanging messages with them.

Sockets are typed according to their communication properties. See Section 4.1.1.3 for a description of the available socket types.

# 4.1.1.2 Communication Domains

Communication domains define the semantics of communication between systems whose hardware and software differ. Communication domains specify the following:

- A set of protocols called the protocol family
- A set of rules for manipulating and interpreting names
- A collection of related socket address formats (an address family)

The socket address for the Internet communication domain contains an Internet address and a port number. The socket address for the UNIX communication domain contains a local pathname.

See Section 4.2.3.4 for more information on socket-related data structures.

Digital UNIX provides default support for the following socket domains<sup>1</sup>:

#### UNIX domain

Digital UNIX provides socket communication between processes running on the same system when a domain of AF\_UNIX is specified. In the UNIX communication domain, sockets are named with UNIX pathnames, such as /dev/printer.

#### · Internet domain

Digital UNIX provides socket communication between a process running locally and one running on a remote host when a domain of AF\_INET is specified. This domain requires that TCP/IP be configured and running on your system.

Table 4-1 summarizes the characteristics of the UNIX and Internet domains.

Table 4-1: Characteristics of the UNIX and Internet Communication Domains

	UNIX	Internet
Socket Types	SOCK_STREAM, SOCK_DGRAM	SOCK_STREAM, SOCK_DGRAM, SOCK_RAW.
Naming	String of ASCII characters, for example, /dev/printer.	32-bit Internet address plus 16-bit port number.
Security	Process connecting to a pathname must have write access to it.	Not applicable.
Raw Access	Not applicable.	Privileged process can access the raw facilities of IP. Raw socket is associated with one IP protocol number, and receives all traffic received for that protocol.

 $<sup>^1\,</sup>$  Digital UNIX can also be configured to support the AF\_DLI domain. For information about the Data Link Interface and using the AF\_DLI domain, see Appendix E.

## 4.1.1.3 Socket Types

Each socket has an associated abstract type which describes the semantics of communications using that socket type. Properties such as reliability, ordering, and prevention of duplication of messages are determined by the socket type. The basic set of socket types is defined in the <sys/socket.h> header file.

#### Note

Typically, header file names are enclosed in angle brackets (< >). To obtain the absolute path to the header file, prepend /usr/include/ to the information enclosed in the angle brackets. In the case of <sys/socket.h>, socket.h is located in the /usr/include/sys directory.

Within the UNIX and Internet domains you can use the following socket types:

SOCK\_DGRAM

Provides datagrams that are connectionless messages of a fixed maximum length where each message can be addressed individually. This type of socket is generally used for short messages because the order and reliability of message delivery is not guaranteed. An important characteristic of a datagram socket is that record boundaries in data are preserved, so individual datagrams are kept separate when they are read.

Often datagrams are used for requests that require a response or responses from the recipient, such as with the finger program. If the recipient does not respond in a specified period of time sending application can repeat the request. The time period varies with the communication domain.

In the UNIX domain, SOCK\_DGRAM is similar to a message queue. In the Internet domain, SOCK\_DGRAM is implemented using the User Datagram Protocol (UDP).

SOCK\_STREAM

Provides sequenced, two-way byte streams across a connection with a transmission mechanism for out-of-band data. The data is transmitted on a reliable basis, in order.

In the UNIX domain, SOCK\_STREAM is like a full-duplex pipe. In the Internet domain, SOCK\_STREAM is implemented using the Transmission Control Protocol (TCP).

SOCK RAW

Provides access to network protocols and interfaces. Raw sockets are only available to privileged processes.

A raw socket allows an application to have direct access to lower-level communications protocols. Raw sockets are intended for advanced users who want to employ protocol features not directly accessible through a normal interface, or who want to build new protocols using existing lower-level protocols. You can also use SOCK\_RAW to communicate with hardware interfaces.

Raw sockets are normally datagram-oriented, though their exact characteristics depend on the interface provided by the protocol. They are available only within the Internet domain.

#### 4.1.1.4 Socket Names

Sockets can be named, which allows unrelated processes on a system or network to locate a specific socket and to exchange data with it. The bound name is a variable-length byte string that is interpreted by the supporting protocol or protocols. Its interpretation varies from communication domain to communication domain. In the Internet domain, names contain an Internet address and port number, and the family is AF\_INET. In the UNIX domain, names contain a pathname and the family is AF\_UNIX.

Communicating processes are bound by an association. In the Internet domain, an association comprises a protocol, local and foreign addresses, and local and foreign ports. When a name is bound to a socket in the Internet domain, the local address and port are specified.

In the UNIX domain, an association comprises local pathnames. Binding a name to a socket in the UNIX domain means specifying a pathname.

In most domains, associations must be unique.

# 4.2 Application Interface to Sockets

The kernel implementation of sockets separates the networking subsystem into the following three interacting layers:

- The socket layer which supplies the interface between the application program and the lower layers, such as the Transmission Control Protocol (TCP) or the User Datagram Protocol (UDP) and IP.
- The protocol layer which consists of transport layer protocols (TCP and UDP) and network layer protocols (IP).

 The device layer which consists of the ifnet layer and the device driver.

In addition to the abstractions described in Section 4.1.1, the socket interface is comprises system and library calls, library functions, and data structures that enable you to manipulate sockets and send and receive data.

Additionally, the kernel provides ancillary services to the sockets framework, such as buffer management, message routing, standardized interfaces to the protocols, and interfaces to the network interface drivers for use by the various network protocols.

#### 4.2.1 Modes of Communication

The sockets framework supports connection-oriented and connectionless modes of communication. Connection-oriented communication means that the application specifies a socket type in a communication domain that supports a connection-oriented protocol. For example, an application could open a SOCK\_STREAM socket in the AF\_INET domain. SOCK\_STREAM sockets in the AF\_INET domain are supported by the TCP protocol, which is a connection-oriented protocol.

Connectionless communication means that the application specifies a socket type in a communication domain that supports a connectionless protocol. For example, a SOCK\_DGRAM socket in the AF\_INET communication domain is supported by the UDP protocol, which is a connectionless protocol.

#### 4.2.1.1 Connection-Oriented Communication

TCP is the connection-oriented protocol implemented on Digital UNIX. TCP is a reliable end-to-end transport protocol that provides for recovery of lost data, transmission errors, and failures of intervening gateways. TCP ensures accurate delivery of data by requiring that two processes be connected before communicating. TCP/IP connections are often compared to telephone connections. Data passed through a SOCK\_STREAM socket in the AF\_INET domain is divided into segments and identified by sequence numbers. The remote process acknowledges receipt of data by including sequence numbers in the acknowledgement. If data is lost enroute, it is resent; thus ensuring that data arrives in the correct sequence to the application.

For applications where large amounts of data are exchanged and the sequence in which the data arrives is important, connection-oriented communication is preferable. File transfer programs are a good example of applications that benefit from the connection-oriented mode of communication offered by TCP.

#### 4.2.1.2 Connectionless Communication

UDP is the connectionless protocol implemented on Digital UNIX. UDP functions as follows:

- Delivers messages based on the messages' address information.
- Requires no connection between communicating processes
- Does not use acknowledgements to ensure that data arrives
- Does not order incoming messages
- Provides no feedback to control the rate at which data is exchanged between hosts.

UDP messages can be lost, duplicated, or arrive out of order. UDP/IP connections are often compared to the postal service.

Where small amounts of data are exchanged and sequencing is not vital, connectionless communication works well. A good example of a program that uses connectionless communication is the rwhod daemon, which periodically broadcasts UDP packets containing system information to the network. It matters little whether or in what sequence those packets are delivered.

UDP is also appropriate for applications that use IP multicast for delivery of datagrams to a subset of hosts on a local area network.

# 4.2.2 Client/Server Paradigm

The most commonly used paradigm in constructing distributed applications is the client/server model. A server process offers services to a network; a client process uses those services. The client and server require a well-known set of conventions before services is rendered and accepted. This set of conventions a protocol comprises that must be implemented at both ends of a connection. Depending on the situation, the protocol can be connection-oriented (asymmetric) or connectionless (symmetric).

In a connection-oriented protocol, such as TCP, one side is always recognized as the server and the other as the client. The server binds a socket to a well-known address associated with the service and then passively listens on its socket. The client requests services from the server by initiating a connection to the server's socket. The server accepts the connection and then server and client can exchange data. An example of a connection-oriented protocol application is Telnet.

In a connectionless protocol, such as UDP, either side can play the server or client role. The client does not establish a connection with the server; instead, it sends a datagram to the server's address. Similarly, the server does not accept a connection from a client. Rather, it issues a recvfrom system call that waits until data arrives from a client. (See Section 4.3.6.)

# 4.2.3 System Calls, Library Calls, Header Files, and Data Structures

This section lists the system and library calls that the socket layer comprises. It also lists the header files that define socket-related constants and structures, and describes some of the most important data structures contained in those header files.

# 4.2.3.1 Socket System Calls

Table 4-2 lists the socket system calls and briefly describes their function. Note that each call has an associated reference page by the same name.

**Table 4-2: Socket System Calls** 

System Call	Description
accept	Accepts a connection on a socket to create a new socket.
bind	Binds a name to a socket.
connect	Initiates a connection on a socket.
getpeername	Gets the name of the connected peer.
getsockname	Gets the socket name.
getsockopt	Gets options on sockets.
listen	Listens for socket connections and specifies the maximum number of queued requests.
recv	Receives messages, peeks at incoming data, and receives out-of-band data.
recvfrom	Receives messages. Has all of the functions of the recv call, plus supplies the address of the peer process.
recvmsg	Receives messages. Has all of the functions of the recv and recvfrom calls, plus receives specially interpreted data (access rights), and performs scatter I/O operations on message buffers.
send	Sends messages. Also sends out-of-band data and normal data without network routing.
sendmsg	Sends messages. Has all of the functions of the send and sendto calls, plus transmits specially interpreted data (access rights), and performs gather I/O operations on message buffers.

## Table 4-2: (continued)

System Call	Description
sendto	Sends messages. Has all of the functions of the send call, plus supplies the address of the peer process.
setsockopt	Sets socket options.
shutdown	Shuts down all socket send and receive operations.
socket	Creates an endpoint for communication and returns a descriptor.
socketpair	Creates a pair of connected sockets.

# 4.2.3.2 Socket Library Calls

Application programs use socket library calls to construct network addresses for use by the interprocess communications facilities in a distributed environment.

Network library subroutines map the following items:

- Host names to network addresses
- Network names to network numbers
- Protocol names to protocol numbers
- Service names to port numbers

Additional socket library calls exist to simplify manipulation of names and addresses.

An application program must include the <netdb.h> header file when using any of the socket library calls.

#### **Host Names**

Application programs use the following network library routines to map Internet host names to addresses:

- gethostbyname
- gethostbyaddr

The gethostbyname routine takes an Internet host name and returns a hostent structure, while the gethostbyaddr routine maps Internet host

addresses into a hostent structure. The hostent structure consists of the following components:

The gethostbyaddr and gethostbyname subroutines return the official name of the host and its public aliases, along with the address family and a null terminated list of variable-length addresses. This list of addresses is required because it is possible for a host to have many addresses with the same name.

The database for these calls is the /etc/hosts file. If the named name server is running, the hosts database is maintained on a designated server on the network. Because of the differences in the databases and their access protocols, the information returned can differ. When using the /etc/hosts version of gethostbyname, only one address is returned, but all listed aliases are included. The named version can return alternate addresses, but does not provide any aliases other than one given as a parameter value.

#### **Network Names**

Application programs use the following network library routines to map network names to numbers and network numbers to names:

- getnetbyaddr
- getnetbyname
- getnetent

The getnetbyaddr, getnetbyname, and getnetent routines extract their information from the /etc/networks file and return a netent structure, as follows:

#### **Protocol Names**

Application programs use the following network library routines to map protocol names to protocol numbers:

- getprotobynumber
- getprotobyname
- getprotoent

The getprotobynumber, getprotobyname, and getprotoent subroutines extract their information from the /etc/protocols file and return the protoent entry, as follows:

## **Service Names**

Application programs use the following network library routines to map service names to port numbers:

- getservbyname
- getservbyport
- getservent

A service is expected to reside at a specific port and employ a particular communication protocol. This view is consistent with the Internet domain, but inconsistent with other network architectures. Further, a service can reside on multiple ports. If this occurs, the higher-level library routines must be bypassed or extended. Services available are contained in the /etc/services file. A service mapping is described by the servent structure, as follows:

The getservbyname routine maps service names to a servent structure by specifying a service name and, optionally, a qualifying protocol. Thus, the following call returns the service specification for a Telnet server by using any protocol:

```
sp = getservbyname("telnet", (char *) NULL);
```

In contrast, the following call returns only the Telnet server that uses the TCP protocol:

```
sp = getservbyname("telnet", "tcp");
```

The getservbyport and getservent routines are also provided. The getservbyport routine has an interface similar to that provided by getservbyname; an optional protocol name can be specified to qualify lookups.

#### **Network Byte Order Translation**

When you have to create or interpret Internet Protocol (IP) suite data in your program, standard methods exist for conversion. The IP suite ensures consistency by requiring particular data formats. Digital UNIX provides functions that let a program convert data to and from those formats. Additionally, the Internet Protocol suite assumes that the most significant byte is in the lowest address, a format known as big-endian. Functions are available to convert from network-byte order to host-byte order and vice versa.

Four functions ensure that data passed by your program is interpreted correctly by the network and vice versa:

- htonl
- htons
- ntohl
- ntohs

Application programs use the following related network library routines to manipulate Internet address strings and 32-bit address quantities:

- inet addr
- inet\_lnaof
- inet makeaddr
- inet netof
- inet\_network
- inet ntoa

Table 4-3 lists and briefly describes the socket library calls. Note that each call has an associated reference page by the same name. The socket library calls are part of libc, so there is no need to link in a special library.

Table 4-3: Socket Library Calls

Name	Description
endhostent	Ends a series of host entry lookups.
endnetent	Ends a series of network entry lookups.
endprotoent	Ends a series of protocol entry lookups.
endservent	Ends a series of service entry lookups.
gethostbyaddr	Given the address of a host, retrieves the host entry from either the name server (named) or the /etc/hosts file.
gethostbyname	Given the name of a host, retrieves the host entry from either the name server (named) or the /etc/hosts file.
gethostent	Retrieves the next host entry from either the name server (named) or the /etc/hosts file, opening this file if necessary.
getnetbyaddr	Given the address of a network, retrieves the network entry from the /etc/networks file.
getnetbyname	Given the name of a network, retrieves the network entry from the /etc/networks file.
getnetent	Retrieves the next network entry from the /etc/networks file, opening this file if necessary.
getprotobyname	Given the protocol name, retrieves the protocol entry from the /etc/protocols file.
getprotobynumber	Given the protocol number, retrieves the protocol entry from the /etc/protocols file.
getprotoent	Retrieves the next protocol entry from the /etc/protocols file, opening this file if necessary.
getservbyname	Given the name of a service, retrieves the service entry from the /etc/services file.
getservbyport	Given the port number of a service, retrieves the service entry from the /etc/services file.
getservent	Retrieves the next service entry from the /etc/services file, opening this file if necessary.
htonl	Converts a 32 bit integer from host-byte order to Internet network-byte order.

Table 4-3: (continued)

Name	Description
htons	Converts an unsigned short integer from host-byte order to Internet network-byte order.
inet_addr	Breaks apart a character string representing numbers expressed in the Internet standard dot (.) notation, and returns an Internet address.
inet_lnaof	Breaks apart an Internet host address and returns the local network address.
inet_makeaddr	Constructs an Internet address from an Internet network number and a local network address.
inet_ntoa	Translates an Internet integer address into a dot- formatted character string.
inet_netof	Breaks apart an Internet host address and returns the network number.
inet_network	Breaks apart a character string representing numbers expressed in the Internet standard dot (.) notation, and returns an Internet network number.
ntohl	Converts a 32 bit integer from Internet network standard-byte order to host-byte order.
ntohs	Converts an unsigned short integer from Internet network-byte order to host-byte order.
sethostent	Begins a series of host entry lookups.
setnetent	Begins a series of network entry lookups.
setprotoent	Begins a series of protocol entry lookups.
setservent	Begins a series of service entry lookups.

## 4.2.3.3 Header Files

Socket header files contain data definitions, structures, constants, macros, and options used by the socket system calls and subroutines. An application program must include the appropriate header file to make use of structures or other information a particular socket system call or subroutine requires. Table 4-4 lists commonly used socket header files.

Table 4-4: Header Files for the Socket Interface

File Name	Description
<sys socket.h=""></sys>	Contains data definitions and socket structures. You need to include this file in all socket applications.
<sys types.h=""></sys>	Contains data type definitions. You need to include this file in all socket applications. This header file is included in <sys socket.h="">.</sys>
<sys un.h=""></sys>	Defines structures for the UNIX domain. You need to include this file in your application if you plan to use UNIX domain sockets.
<netinet in.h=""></netinet>	Defines constants and structures for the Internet domain. You need to include this file in your application if you plan to use TCP/IP in the Internet domain.
<netdb.h></netdb.h>	Contains data definitions for socket subroutines. You need to include this file in your application if you plan to use TCP/IP and need to look up host entries, network entries, protocol entries, or service entries.

#### 4.2.3.4 Socket Related Data Structures

This section describes the following data structures:

- sockaddr
- sockaddr in
- sockaddr\_un
- msghdr

The sockaddr structures contain information about a socket's address format. Because the communication domain in which an application creates a socket determines its address format, it also determines its data structure.

Socket address data structures are defined in the header files described in Section 4.2.3.3. Which header file is appropriate depends on the type of socket you are creating. The possible types of socket address data structures are as follows:

#### struct sockaddr

Defines the generic version of the socket address structure. These sockets are limited to 14 bytes of direct addressing. The

<sys/socket.h> file contains the sockaddr structure, which
contains the following elements:

The sa\_len parameter defines the total length. The sa\_family parameter defines the socket address family or domain, which is AF\_UNIX for the UNIX domain or AF\_INET for the Internet domain. The contents of sa\_data depend on the protocol in use, but generally a socket name consists of a machine-name part and a port-name or service-name part.

```
struct sockaddr_un
```

Defines UNIX domain sockets used for communications between processes on the same machine. These sockets require the specification of a full pathname. The <sys/un.h> header file contains the sockaddr\_un structure. The sockaddr\_un structure contains the following elements:

UNIX domain protocols (AF\_UNIX) have socket addresses up to PATH\_MAX plus 2 bytes long. The PATH\_MAX parameter defines the maximum number of bytes of the pathname.

```
struct sockaddr_in
```

Defines Internet domain sockets used for machine-to-machine communication across a network and local interprocess communication. The <netinet/in.h> file contains the sockaddr\_in structure. The sockaddr in structure contains the following elements:

```
unsigned char sin_len;
sa_family_t sin_family;
in_port_t sin_port;
struct in addr sin addr;
```

The Internet networking routines only support 16-byte structures. Sockets created in the Internet domain (AF\_INET), therefore, have socket addresses that do not exceed 16 bytes.

The msghdr data structure, which is defined in the <sys/socket.h> header file, allows applications to pass access rights to system-maintained objects (such as files, devices, or sockets) using the sendmsg and recvmsg system calls. (See Section 4.3.6 for information on the sendmsg and recvmsg system calls.) The processes transmitting data must be connected with a UNIX domain socket.

The data structure also allows AF\_INET sockets to receive certain data. See the descriptions of the IP\_RECVDSTADDR and IP\_RECVOPTS options in the ip(7) reference page.

The msghdr data structure consists of the following components:

In addition to the XPG4 msghdr data structure, Digital UNIX also supports the 4.3BSD and the 4.4BSD versions of this data structure. The BSD versions of the msghdr data structure are described in greater detail in Section 4.4.

# 4.3 Using Sockets

This section outlines the steps required to create and use sockets. Connection-oriented and connectionless modes of communication are described in the following sections:

Creating sockets

Describes how to create a socket with the socket and socketpair system calls.

Binding names and addresses

Describes how to bind a name and address to a socket with the bind system call.

• Establishing connections (clients)

Describes how to use the connect system call on a client to connect to a server.

Accepting connections (servers)

Describes how to use the listen and accept system calls to connect a server to a client.

• Setting and getting socket options

Describes how to use the setsockopt and getsockopt system calls to set and retrieve the values of socket characteristics.

· Transferring data

Describes how to use the read and write system calls, as well as the send and recv related system calls to transmit data.

- Shutting down sockets
  - Describes how to use the shutdown system call to shut down a socket.
- Closing sockets

Describes how to use the close system call to close a socket.

# 4.3.1 Creating Sockets

The first step in using sockets is creating a socket. Sockets are opened, or created, with the socket or socketpair system calls.

The syntax of the socket system call is as follows:

```
s = socket (domain, type, protocol);
```

In the preceding statement:

```
domain
```

Specifies the communication domain; for example AF\_UNIX or AF INET.

type

Specifies the socket type as SOCK\_STREAM, SOCK\_DGRAM, or SOCK\_RAW.

```
protocol
```

Specifies the transport protocol, such as TCP or UDP. If protocol is specified as zero (0), the system selects an appropriate protocol from those protocols that the communication domain comprises and that can be used to support the requested socket type.

See socket(2) for more information.

The socket call returns a socket descriptor, s, which is an a nonnegative integer that the application program uses to reference the newly created socket in subsequent system calls. The socket descriptor returned is the lowest unused number available in the calling process for such descriptors and is an index into the kernel descriptor table.

For example, to create a stream socket in the Internet domain, you can use the following call:

```
if ((s = socket(AF_INET, SOCK_STREAM,0)) == -1 ) {
            fprintf(file1,"socket() failed\n");
            local_flag = FAILED;
     }
```

This call results in the creation of a stream socket with the TCP protocol providing the underlying communication support. To create a datagram

socket in the UNIX domain, you can use the following call:

```
if ((s = socket(AF_UNIX, SOCK_DGRAM,0)) == -1 ) {
          fprintf(file1, "socket() failed\n");
          local_flag = FAILED;
}
```

This call results in the creation of a datagram socket with a UNIX domain protocol providing the underlying communication support.

The socketpair system call can also be used to create sockets. The socketpair system call creates an unnamed pair of sockets that are already connected. The syntax of the socketpair system call is as follows:

socketpair (domain, type, protocol, socket\_vector[2]);

In the preceding statement:

```
domain
```

Specifies the communication domain. An application using the socketpair system call must specify AF\_UNIX.

type

Specifies the socket type. Can be SOCK\_DGRAM or SOCK\_STREAM.

protocol

Specifies the optional identifier used to define the transport protocol. The value of this variable is always zero (0).

```
socket vector[2]
```

Specifies a two-integer array used to define the file descriptors of the socket pair.

See socketpair(2) for more information.

The socketpair system call returns a pair of socket descriptors, which are a nonnegative integers, that the application uses to reference the newly created socket pair in subsequent system calls.

The following example shows how to create a socket pair:

```
{
    :
    int sv[2];
    :
```

```
if ((s = socketpair (AF_UNIX, SOCK_STREAM, 0, sv)) < 0) {
          local_flag=FAILED;
          fprintf(file1, "socketpair() failed\n");
     }
...
...</pre>
```

## 4.3.1.1 Setting Modes of Execution

Sockets can be set to blocking or nonblocking I/O mode. The O\_NONBLOCK fcntl operation is used to determine this mode. When O\_NONBLOCK is clear (not set), which is the default, the socket is in blocking mode. In blocking mode, when the socket tries to do a read and the data is not available, it waits for the data to become available.

When O\_NONBLOCK is set, the socket is in nonblocking mode. In nonblocking mode, when the calling process tries to do a read and the data is not available, the socket returns immediately with the EWOULDBLOCK error code. It does not wait for the data to become available. Similarly, during writing, when a socket has O\_NONBLOCK set and the output queue is full, an attempt by the socket to write causes the process to return immediately with an error code of EWOULDBLOCK.

The following example shows how to mark a socket as nonblocking:

```
#include <fcntl.h>
    .
    .
    int s;
    .
    if (fcntl(s, F_SETFL, O_NONBLOCK) < 0)
        perror("fcntl F_SETFL, O_NONBLOCK");
        exit(1);
}
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
    .
```

When performing nonblocking I/O on sockets, a program must check for the EWOULDBLOCK error, which is stored in the global value errno. The EWOULDBLOCK error occurs when an operation normally blocks, but the socket on which it was performed is marked as nonblocking. The following socket system calls all return the EWOULDBLOCK error code:

- accept
- connect
- send

- sendto
- sendmsg
- recv
- recvfrom
- recvmsg
- read
- write

Processes that use these system calls on nonblocking sockets must be prepared to deal with the EWOULDBLOCK return codes.

When an operation, such as a send, cannot be completed but partial writes are permissible (for example, when using a SOCK\_STREAM socket), the data that can be sent immediately is processed, and the return value indicates the amount of data actually sent.

# 4.3.2 Binding Names and Addresses

The bind system call associates an address with a socket. The domain for the socket is established with the socket system call. Regardless of the domain in which the bind system call is used, it allows the local process to fill in information about itself, for example, the local port or local pathname. This information allows the server application to be located by a client application.

The syntax for the bind system call is as follows:

```
bind(s, *address, address_len);
```

In the preceding statement:

S

Specifies the socket descriptor.

\*address

Specifies a pointer to a protocol-specific address structure.

address\_len

Specifies the size of the address in \*address.

The following example shows how to use the bind system call on a

#### SOCK\_STREAM socket created in the Internet domain:

See Section 4.6.2 and bind(2) for more information.

# 4.3.3 Establishing Connections

Sockets are created in the unconnected state. Client processes use the connect system call to connect to a server process or to store a server's address locally, depending on whether the communication is connection-oriented or connectionless. For the Internet domain, the connect system call typically causes the local address, local port, foreign address, and foreign port of an association to be assigned.

The syntax of the connect system call depends on the communication domain. The syntax of the connect system call is as follows:

```
connect(s, *address, address_len);
```

In the preceding statement:

s

Specifies the socket descriptor.

\*address

Specifies the server's address to which the client wants to connect.

address\_len

Specifies the size, in bytes, of the address of the server.

An error is returned if the connection was unsuccessful; any name automatically bound by the system remains, however. Common errors associated with sockets are listed in Table 4-5 in Section 4.5. If the connection is successful, the socket is associated with the server and data transfer begins.

See connect(2) for more information.

Selecting a connection-oriented protocol in the Internet domain means choosing TCP. In such cases, the connect system call builds a TCP connection with the destination, or returns an error if it cannot. Client processes using TCP must call the connect system call to establish a connection before they can transfer data through a reliable stream socket (SOCK STREAM).

Selecting a connectionless protocol in the Internet domain means choosing UDP. Client processes using connectionless protocols do not have to be connected before they are used. If connect is used under these circumstances, it stores the destination (or server) address locally so that the client process does not need to specify the server's address each time a message is sent. Any data sent on this socket is automatically addressed to the connected server process and only data received from that server process is delivered.

Only one connected address is permitted at any time for each socket; a second connect system call changes the destination address and a connect system call to a null address (address INADDR\_ANY) causes a disconnect. The connect system call on a connectionless protocol returns immediately, since it results in the operating system recording the server's socket's address (as compared to a connection-oriented protocol, where a connect request initiates establishment of an end-to-end connection).

While a socket using a connectionless protocol is connected, errors from recent send system calls can be returned asynchronously. These errors can be reported on subsequent operations on the socket. A special socket option, SO\_ERROR (used with the getsockopt system call), can be used to query the error status. A select system call, issued to determine when more data can be sent or received, will return true when a process has received an error indication.

In any case, the next operation will return the error and clear the error status.

The syntax of the select system call is as follows:

select (nfds, \*readfds, \*writefds, \*exceptfds, \*timeout);

In the preceding statement:

nfds

Specifies the number of bits that represent open object file descriptors ready for reading or writing, or that have an exception pending.

\*readfds and \*writefds

Point to an I/O descriptor set consisting of file descriptors of objects open for reading or writing.

\*exceptfds

Points to an I/O descriptor set consisting of file descriptors for objects opened for reading or writing that have an exception pending.

\*timeout

Points to a type timeval structure that specifies the time to wait for a response to a select function.

See select(2) for more information.

The following is an example of the select system call:

```
if ( (ret_val = select(20,&read_mask,NULL,NULL,&tp)) != i )
```

#### 4.3.4 Accepting Connections

A connection-oriented server process normally listens at a well-known address for service requests. That is, the server process remains dormant until a connection is requested by a client's connection to the server's address. Then, the server process wakes up and services the client by performing the actions the client requests.

Connection-oriented servers use the listen and accept system calls to prepare for and then accept connections with client processes.

The listen system call is usually called after the socket and bind system calls. It indicates that the server is ready to receive connect requests from clients.

The syntax of the listen system call is as follows:

listen (s, backlog);

In the preceding statement:

S

Specifies the socket descriptor.

backlog

Specifies the maximum number of outstanding connection requests that this server can queue. If the queue is full, the server rejects the connect request and the client must try again.

Servers that process a small number of connections can specify a small backlog. Servers that process a high volume of connections can specify a

larger value. The kernel imposes an upper limit on the backlog, which is determined by the SOMAXCONN parameter.

See listen(2) for more information.

The server accepts a connection to a client by using the accept system call.

The syntax of the accept system call is as follows:

```
accept (s, *address, *address_len);
```

In the preceding statement:

ې

Specifies the socket descriptor.

\*address

Specifies a pointer to a protocol-specific address structure. On return this contains the address of the connecting entity.

```
*address_len
```

Specifies the size of the address structure.

See accept(2) for more information.

An accept call blocks the server until a client requests service. This call returns a failure status if the call is interrupted by a signal such as SIGCHLD. Therefore, the return value from accept is checked to ensure that a connection was established.

When the connection is made, the server normally forks a child process which creates another socket with the same properties as socket s (the socket on which it is listening). Note in the following example how the socket s, used by the parent for queuing connection requests, is closed in the child while the socket g, which is created as a result of the accept call, is closed in the parent. The address of the client is also handed to the doit routine because it is required for authenticating clients. After the accept system call creates the new socket, it allows the new socket to service the client's connection request while it continues listening on the original socket; for example:

```
for (;;) {
  int g, len = sizeof (from);

  g = accept(s, (struct sockaddr *)&from, &len);
  if (g < 0) {
    if (errno != EINTR)
        syslog(LOG_ERR, "rlogind: accept: %m");
    continue;</pre>
```

```
}
if (fork() == 0) {    /* Child */
    close(s);
    doit(g, &from);
}
close(g);    /* Parent */
```

Connectionless servers use the bind system call but, instead of using the accept system call, they use a recvfrom system call and then wait for client requests. No connection is established between the connectionless server and client during the process of exchanging data.

# 4.3.5 Setting and Getting Socket Options

In addition to binding a socket to a local address or connecting to a destination address, application programs must be able to control the socket. For example, with protocols that use time-out and retransmission, the application program may want to obtain or set the time-out parameters. It may also want to control the allocation of buffer space, determine if the socket allows transmission of a broadcast, or control processing of out-of-band data.

The getsockopt and setsockopt system calls provide the application program with the means to control socket operations. The setsockopt system call allows an application program to set a socket option by using the same set of values obtained with the getsockopt system call.

The syntax of the setsockopt system call is as follows:

setsockopt(s, level, optname, \*optval, optlen);

In the preceding statement:

~

Specifies the socket file descriptor.

level

Specifies what portion of code in the system interprets the *optname* parameter; for example, general socket layer or protocol transport layer. To set a socket-level option, set *level* to SOL\_SOCKET, which is defined in the <sys/socket.h> header file. To set a TCP level option, set *level* to IPPROTO\_TCP, which is defined in the <netinet/in.h> header file.

optname

Specifies the name of the option to set, for example, SO\_SNDBUF. Socket options are defined in the <sys/socket.h> header file.

```
*optval
```

Points to a buffer containing data specific to the option being set. This data may specify a Boolean, integer, or some other value, including values in structures.

```
optlen
```

Specifies the size of the buffer to which the optval parameter points.

See setsockopt(2) for more information.

The following example shows how to set the SO\_SNDBUF option on a socket in the Internet communication domain:

The getsockopt system call allows an application program to request information about the socket options that are set with the setsockopt system call.

The syntax of the getsockopt system call is as follows:

```
getsockopt(s, level, optname, *optval, *optlen);
```

The parameters are the same as for the setsockopt system call, with the exception of the *optlen* parameter, which is a pointer to the size of the buffer.

The following example shows how the getsockopt system call can be used to determine the size of the SO\_SNDBUF on an existing socket:

The SOL\_SOCKET parameter indicates that the general socket level code is to interpret the SO\_SNDBUF parameter. The SO\_SNDBUF parameter

indicates the size of the send socket buffer in use on the socket.

Not all socket options apply to all sockets. The options that can be set depend on the address family and protocol the socket uses.

# 4.3.6 Transferring Data

Most of the work performed by the socket layer is in sending and receiving data. The socket layer itself does not impose any structure on data transmitted or received through sockets. Any data interpretation or structuring is logically isolated in the implementation of the communication domain.

The following are the system calls that an application uses to send and receive data:

- read
- write
- send
- sendto
- recv
- recvfrom
- sendmsq
- recvmsg

# 4.3.6.1 Using the read System Call

The read system call allows a process to receive data on a socket without receiving the sender's address.

The syntax for the read system call is as follows:

read (s, \*buf, nbytes);

In the preceding statement:

S

Specifies the socket descriptor.

\*buf

Points to the buffer to receive data.

nbytes

Specifies the size of buf in bytes.

See read(2) for more information.

# 4.3.6.2 Using the write System Call

The write system call is used on sockets in the connected state. The destination of data transferred with the write system call is implicitly specified by the connection.

The syntax for the write system call is as follows:

write (s, \*buf, nbytes);

In the preceding statement:

S

Specifies the socket descriptor.

\*buf

Points to the buffer containing data to be written.

nbytes

Specifies the size of *buf* in bytes.

See write(2) for more information.

## 4.3.6.3 Using the send, sendto, recv and recvfrom System Calls

The send, sendto, recv, and recvfrom system calls are similar to the read and write system calls, sharing the first three parameters with them; however, additional flags are required. The flags, defined in the <sys/socket.h> header file, can be defined as a nonzero value if the application program requires one or more of the following:

Flag	Description
MSG_OOB	Send or receive out-of-band data.
MSG_PEEK	Look at data without reading. Valid for the recy and recyfrom
MSG_DONTROUTE	calls.  Send data without routing packets. Valid for the send and send to calls.

The MSG\_OOB flag signifies out-of-band data, or urgent data, and is specific to stream sockets (SOCK\_STREAM). See Section 4.6.3 for more information about out-of-band data.

The MSG\_PEEK flag allows an application to preview the data that is available to be read, without having the system discard it after the recv or recvfrom call returns. When the MSG\_PEEK flag is specified with a recv system call, any data present is returned to the user but treated as still unread. That is, the next read or recv system call applied to the socket returns the data previously previewed.

The MSG\_DONTROUTE flag is currently used only by the routing table management process and is not discussed further.

#### send

The send system call is used on sockets in the connected state. The send and write system calls function almost identically; the only difference is that send supports the flags described at the beginning of this section.

The syntax for the send system call is as follows:

```
send (s, *message, len, flags);
```

In the preceding statement:

s

Specifies the socket descriptor.

\*message

Points to the buffer containing data to send.

1en

Specifies the length of message in bytes.

flags

Allows the sender to control message transmission. Can be one of the three flags described at the beginning of this section.

See send(2) for more information.

#### sendto

The sendto system call is used on connected or unconnected sockets. It allows the process explicitly to specify the destination for a message.

The syntax for the sendto system call is as follows:

sendto(s, \*message, len, flags, \*dest\_addr, dest\_len);

In the preceding statement:

s

Specifies the socket descriptor.

\*message

Points to the buffer containing the message to be sent.

1en

Specifies the size of the buffer to which the message parameter points.

flag

Allows the sender to control message transmission. Can be one of the three flags described at the beginning of this section.

\*dest addr

Points to the buffer containing the address of the message's intended recipient. The \*dest\_addr parameter is ignored for SOCK STREAM sockets.

dest\_len

Specifies the size of the address in dest\_addr.

See sendto(2) for more information.

#### recv

The recv system call allows a process to receive data on a socket without receiving the sender's address. The read and recv system calls function almost identically; the only difference is that recv supports the flags described at the beginning of this section.

The syntax for the recv system call is as follows:

recv (s, \*message, len, flags);

In the preceding statement:

s

Specifies the socket descriptor.

\*message

Points to a buffer where the message should be placed.

len

Specifies the size of the buffer to which the *message* parameter points.

flags

Allows the receiver to control message reception. Can be one of the three flags described at the beginning of this section.

See recv(2) for more information.

#### recvfrom

The recvfrom system call can be used on connected or unconnected sockets. The recvfrom system call has similar functionality to the recv call but it additionally allows an application to receive the address of a peer with whom it is communicating.

The syntax for the recvfrom system call is as follows:

```
recvfrom (s, *buf, len, flags, *src addr, *src len);
```

In the preceding statement:

S

Specifies the socket descriptor.

\*buf

Points to the buffer to receive data.

1en

Specifies the size of the buffer in bytes.

flags

Allows the receiver to control message reception. Can be one of the three flags described at the beginning of this section.

```
*src addr
```

Points to a buffer to receive the address of the peer (sender). The \*src\_addr parameter is ignored for SOCK\_STREAM sockets.

\*src len

Specifies the length, in bytes, of the buffer pointed to by \*src addr.

See recvfrom(2) for more information.

## 4.3.6.4 Using the sendmsg and recvmsg System Calls

The sendmsg and recvmsg system calls are distinguished from the other send and receive related system calls in that they allow unrelated processes on the local machine to pass file descriptors to each other. These two system calls are the only ones that support the concept of access rights, which means that the system has granted a process the right to access a system-maintained object. Using the sendmsg and recvmsg system calls they can pass that right to another process.

To pass access rights, the sendmsg and recvmsg system calls use the msghdr data structure. The msghdr data structure defines two parameters, the msg\_control and msg\_controllen that deal with the passing and

receiving of access rights between processes. For more information on the msghdr data structure, see Section 4.2.3.4 and Section 4.4.2.

Although the sendmsg and recvmsg system calls can be used on connection-oriented or connectionless protocols and in the Internet or UNIX domains, for processes to pass descriptors they must be connected with a UNIX domain socket.

## sendmsg

The sendmsg system call is used on connected or unconnected sockets. It transfers data using the msghdr data structure. For more information on the msghdr data structure, see Section 4.2.3.4 and Section 4.4.2.

The syntax for the sendmsg system call is as follows:

```
sendmsg(s, *message, flags);
```

In the preceding statement:

s

Specifies the socket descriptor.

\*message

Points to a msghdr structure. For more information on the msghdr structure, see Section 4.2.3.4 and Section 4.4.2.

flags

Contains the size and address of the buffer of control data.

See sendmsg(2) for more information.

The following is an example of the sendmsg system call:

```
struct msghdr send;
struct iovec saiov;
struct sockaddr destAddress;
char sendbuf[BUFSIZE];
send.msg_name = (void *)&destAddress;
send.msg_namelen = sizeof(destAddress);
send.msg_iov = &saiov;
send.msg_iovlen = 1;
saiov.iov base = sendbuf;
saiov.iov_len = sizeof(sendbuf);
send.msg_control = NULL;
send.msg_controllen = 0;
send.msg_flags = 0;
if ((i = sendmsg(s, \&send, 0)) < 0) {
       fprintf(file1, "sendmsg() failed\n");
        exit(1);
}
```

#### recvmsg

The recvmsg system call is used on connected or unconnected sockets. It transfers data using the msghdr data structure. For more information on the msghdr data structure, see Section 4.2.3.4 and Section 4.4.2.

The syntax of the recvmsg system call is as follows:

```
recvmsg(s, *message, flags);
```

In the preceding statement:

S

Specifies the socket descriptor.

\*message

Points to a msghdr structure. For more information on the msghdr structure, see Section 4.2.3.4 and Section 4.4.2.

flags

Allows the sender to control the message transmission. Can be one of the flags described at the beginning of Section 4.3.6.3.

See recvmsg(2) for more information.

The following is an example of the recvmsg system call:

```
struct msghdr recv;
struct iovec recviov;
struct sockaddr_in recvaddress;
char recvbuf[BUFSIZE];
recv.msg_name = (void *) &recvaddress;
recv.msg_namelen = sizeof(recvaddress);
recv.msg_iov = &recviov;
recv.msg_iovlen = 1;
recviov.iov_base = recvbuf;
recviov.iov_len = sizeof(recvbuf);
recv.msg_control = NULL;
recv.msg_controllen = 0
recv.msg_flags = 0
if ((i = recvmsg(r, \&recv, 0)) < 0) {
              fprintf(file1, "recvmsg() failed\n");
            exit(1);
```

## 4.3.7 Shutting Down Sockets

If an application program has no use for any pending data, it can use the shutdown system call on the socket prior to closing it. The syntax of the shutdown system call is as follows:

shutdown (s, how);

In the preceding statement:

S

Specifies the socket descriptor.

how

Specifies the type of shutdown.

See shutdown(2) for more information.

## 4.3.8 Closing Sockets

The close system call is used to close sockets. The syntax of the close system call is as follows:

close(s);

In the preceding statement:

s

Specifies the socket descriptor.

See close(2) for more information.

Closing a socket and reclaiming its resources can be complicated. For example, a close system call is never expected to fail when a process exits. However, when a socket that is promising reliable delivery of data closes with data still queued for transmission or awaiting acknowledgment of reception, the socket must attempt to transmit the data. When the socket discards the queued data to allow the close call to complete successfully, it violates its promise to deliver data reliably. Discarding data can cause naive processes that depend on the implicit semantics of the close call to work unreliably in a network environment.

However, if sockets block until all data is transmitted successfully, a close system call may never complete in some communication domains.

The socket layer compromises in an effort to address the completion problem and still maintain the semantics of the close system call. In normal operation, closing a socket causes any queued but unaccepted connections to be discarded. If the socket is in a connected state, a disconnect is initiated. The socket is marked to indicate that a descriptor is no longer referencing it,

and the close operation returns successfully. When the disconnect request completes, the network support notifies the socket layer, and the socket resources are reclaimed. The network layer attempts to transmit any data queued in the socket's send buffer, but there is no guarantee that it will succeed.

Alternatively, a socket can be marked explicitly to force the application program to linger when closing until pending data is flushed and the connection shuts down. This option is marked in the socket data structure by using the setsockopt system call with the SO\_LINGER option.

#### Note

The setsockopt system call, using the linger option, takes a linger structure, which is defined in the <sys/socket.h> header file.

When an application program indicates that a socket is to linger, it also specifies a duration for the lingering period. If the lingering period expires before the disconnect is completed, the socket layer forcibly shuts down the socket, discarding any data that is still pending.

#### 4.4 BSD Socket Interface

In addition to the XPG4 socket interface, Digital UNIX also supports the 4.3BSD and 4.4BSD socket interfaces. The 4.4BSD socket interface provides a number of changes to 4.3BSD sockets. Most of the changes between the 4.3BSD and 4.4BSD socket interfaces were designed to facilitate the implementation of International Standards Organization (ISO) protocol suites under the sockets framework. The XPG4 socket interface provides a standard version of the socket interface.

#### Note

The availability of the 4.4BSD socket interface does not mean that your site supports ISO protocols. Check with the appropriate personnel at your site.

To use the 4.4BSD socket interface, you must add the following line to your program or makefile:

#define \_SOCKADDR\_LEN

The 4.4BSD socket interface includes the following changes from the 4.3BSD interface for application programs:

 A sockaddr structure for supporting variable-length (long) network addresses • A msghdr structure to allow receipt of protocol information and status with data

The following sections describe these features.

## 4.4.1 Variable-Length Network Addresses

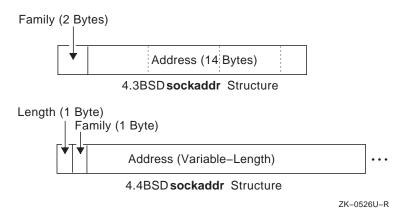
The 4.4BSD version of the sockaddr structure supports variable-length network addresses. The structure adds a length field and is defined as follows:

The 4.3BSD sockaddr structure contains the following fields:

```
u_short sa_family;
char sa_data[14];
```

Figure 4-2 compares the 4.3BSD and 4.4BSD sockaddr structures.

## Figure 4-2: 4.3BSD and 4.4BSD sockaddr Structures



## 4.4.2 Receiving Protocol Data with User Data

The 4.3BSD version of the msghdr structure (which is the default if you use the cc command) provides the parameters needed for using the optional functions of the sendmsg and recvmsg system calls.

#### The 4.3BSD msghdr structure is as follows:

```
/* 4.3BSD msghdr Structure */
struct msghdr {
       caddr_t msg_name;
                                    /* optional address */
             msg_namelen;
                                    /* size of address */
       int.
       struct iovec *msg_iov;
                                    /* scatter/gather array */
       int
            msg_iovlen;
                                    /* # elements in msg_iov */
       caddr_t msg_accrights;
                                    /* access rights sent/re-
                                     /* ceived */
       int.
              msg_accrightslen;
};
```

The  $msg\_name$  and  $msg\_namelen$  parameters are used when the socket is not connected. The  $msg\_iov$  and  $msg\_iovlen$  parameters are used for scatter (read) and gather (write) operations. As stated previously, the  $msg\_accrights$  and  $msg\_accrightslen$  parameters allow the sending process to pass its access rights to the receiving process.

The 4.4BSD structure has additional fields that permit application programs to include protocol information along with user data in messages.

To support the receipt of protocol data together with user data, Digital UNIX provides the msghdr structure from the 4.4BSD socket interface. The structure adds a pointer to control data, a length field for the length of the control data, and a flags field, as follows:

The XPG4 msghdr data structure has the same fields as 4.4BSD. However, the size of the msg\_namelen and msg\_controllen fields are 8 bytes long in the XPG4 msghdr data structure, as opposed to 4 bytes long in the 4.4BSD msghdr data structure. Figure 4-3 shows the 4.3BSD, 4.4BSD, and XPG4 msghdr structures.

Figure 4-3: 4.3BSD, 4.4BSD, and XPG4 msghdr Structures

	msg_name	msg_namelen	_namelen msg_		v msg_iovlen		msg_accrights			
4.3BSD msghdr Structure										
	msg_name	msg_namelen	msg_iov		msg_iovlen	msg_control		msg_ controllen	msg_flags	
4.4BSD <b>msghdr</b> Structure										
	msg_name	msg_nar	amelen ma		_iov	msg_iovlen	msg	_control	msg_ controllen	msg_flags

XPG4 msghdr Structure

ZK-0527U-R

In the 4.3BSD version of the msghdr data structure, the msg\_accrights and msg\_accrightslen fields permit the sending process to pass its access rights to a system-maintained object, in this case a socket, to the receiving process. In the 4.4BSD and XPG4 versions, this is done using the msg\_control and msg\_controllen fields.

## 4.5 Common Socket Errors

Table 4-5 lists some common socket error messages the problems they indicate:

Table 4-5: Common Errors and Diagnostics

Error	Diagnostics
[EAFNOSUPPORT]	The protocol family does not support the addresses in the specified address family.
[EBADF]	The socket parameter is not valid.
[ECONNREFUSED]	The attempt to connect was rejected.
[EFAULT]	A pointer does not point to a valid part of user address space.
[EHOSTDOWN]	The host is down.
[EHOSTUNREACH]	The host is unreachable.
[EINVAL]	An invalid argument was used.
[EMFILE]	The current process has too many open file descriptors

Table 4-5: (continued)

Error	Diagnostics
[ENETDOWN]	The network is down.
[ENETUNREACH]	The network is unreachable. No route to the network is present.
[ENOMEM]	The system was unable to allocate kernel memory to increase the process descriptor table.
[ENOTSOCK]	The socket parameter refers to a file, not a socket.
[EOPNOTSUPP]	The specified protocol does not permit creation of socket pairs.
[EOPNOTSUPP]	The referenced socket can not accept connections.
[EPROTONOSUPPORT]	This system does not support the specified protocol.
[EPROTOTYPE]	The socket type does not support the specified protocol.
[ETIMEDOUT]	The connection timed out without a response from the remote application.
[EWOULDBLOCK]	The socket is marked nonblocking and the operation could not complete.

## 4.6 Advanced Topics

This section contains the following information, which is of interest to developers writing complex applications for sockets:

- Selecting specific protocols
- Binding names and addresses
- Out-of-band data
- IP Multicasting
- Broadcasting and determining network configuration
- The inetd daemon
- Input/output multiplexing
- Interrupt-driven socket I/O
- Signals and process groups
- Pseudoterminals

## 4.6.1 Selecting Specific Protocols

The syntax of the socket system call is described in Section 4.3.1. If the third argument to the socket call, the *protocol* argument, is zero (0), the socket call selects a default protocol to use with the returned socket descriptor. The default protocol is usually correct and alternate choices are not usually available. However, when using raw sockets to communicate directly with lower-level protocols or hardware interfaces, the protocol argument can be important for setting up demultiplexing.

For example, raw sockets in the Internet family can be used to implement a new protocol above IP and the socket receives packets only for the protocol specified. To obtain a particular protocol, you must determine the protocol number as defined within the communication domain. For the Internet domain, you can use one of the library routines described in Section 4.2.3.2.

The following code shows how to use the getprotobyname library call to select the protocol newtop for a SOCK\_STREAM socket opened in the Internet domain:

```
#include <sys/types.h>
#include <sys/socket.h>
#include <netinet/in.h>
#include <netdb.h>

.
.
struct protent *pp;
.
.
pp = getprotobyname("newtcp");
s = socket(AF_INET, SOCK_STREAM, pp->p_proto);
```

## 4.6.2 Binding Names and Addresses

The bind system call associates an address with a socket.

#### 4.6.2.1 Binding to the Wildcard Address

The local machine address for a socket can be any valid network address of the machine. Because one system can have several valid network addresses, binding addresses to sockets in the Internet domain can be complicated. To simplify local address binding, the constant INADDR\_ANY, a wildcard address, is provided. The INADDR\_ANY address tells the system that this server process will accept a connection on any of its Internet interfaces, if it has more than one.

The following example shows how to bind the wildcard value INADDR ANY to a local socket:

```
#include <sys/types.h>
#include <sys/socket.h>
#include <netinet/in.h>
#include <stdio.h>
main()
{
   int s, length;
   struct sockaddr_in name;
   char buf[1024];
   /* Create name with wildcards. */
   name.sin_family = AF_INET;
  name.sin_len = sizeof(name);
  name.sin_addr.s_addr = INADDR_ANY;
   name.sin port = 0;
   if (bind(s, (struct sockaddr *)&name, sizeof(name))== -1) {
      perror("binding datagram socket");
      exit(1);
```

Sockets with wildcard local addresses can receive messages directed to the specified port number, and send to any of the possible addresses assigned to that host. Note that the socket uses a wildcard value for its local address; a process sending messages to the named socket must specify a valid network address. A process can be willing to receive a message from anywhere, but it cannot send a message anywhere.

When a server process on a system with more than one network interface wants to allow hosts to connect to only one of its interface addresses, the server process binds the address of the appropriate interface. For example, if a system has two addresses 130.180.123.45 and 131.185.67.89, a server process can bind the address 130.180.123.45. Binding that address ensures that only connections addressed to 130.180.123.45 can connect to the server process.

Similarly, a local port can be left as unspecified (specified as zero), in which case the system selects a port number for it.

### 4.6.2.2 Binding in the UNIX Domain

Processes that communicate in the UNIX domain (AF\_UNIX) are bound by an association that local and foreign pathnames comprises. UNIX domain sockets do not have to be bound to a name but, when bound, there can never be duplicate bindings of a protocol, local pathname, or foreign pathname. The pathnames cannot refer to files existing on the system. The process that binds the name to the socket must have write permission on the directory where the bound socket will reside.

The following example shows how to bind the name socket to a socket created in the UNIX domain:

```
#include <sys/types.h>
#include <sys/socket.h>
#include <sys/un.h>
#include <stdio.h>
#define NAME "socket"
main()
   int s, length;
   struct sockaddr_un name;
   char buf[1024];
   /* Create name. */
   name.sun_len = sizeof(name.sun_len) +
             sizeof(name.sun_family) +
             strlen(NAME);
   name.sun_family = AF_UNIX;
   strcpy(name.sun_path, NAME);
   if (bind(s, (struct sockaddr *) &name, sizeof(name))==-1) {
     perror("binding name to datagram socket");
      exit(1);
```

#### 4.6.3 Out-of-Band Data

Out-of-band data is a logically independent transmission channel associated with each pair of connected stream sockets. Out-of-band data can be delivered to the socket independently of the normal receive queue or within the receive queue, depending on the status of the SO\_OOBINLINE option, set with the setsockopt system call.

The stream socket abstraction specifies that the out-of-band data facilities must support the reliable delivery of at least one out-of-band message at a time. This message must contain at least one byte of data and at least one message can be pending delivery to the user at any one time.

The socket layer supports marks in the data stream that indicate the end of urgent data or out-of-band processing. The socket mechanism does not return data from both sides of a mark in a single system call.

You can use MSG\_PEEK to peek at out-of-band data. If the socket has a process group, a SIGURG signal is generated when the protocol is notified of its existence. A process can set the process group or process ID to be informed by the SIGURG signal via the appropriate fcntl call, as described in Section 4.6.8 for SIGIO.

When multiple sockets have out-of-band data awaiting delivery, an application program can use a select call for exceptional conditions to determine which sockets have such data pending. The SIGURG signal or select call notifies the application program that data is pending. The application then must issue the appropriate call actually to receive the data.

In addition to the information passed, a logical mark is placed in the data stream to indicate the point at which the out-of-band data was sent. When a signal flushes any pending output, all data up to the logical mark in the data stream is discarded.

To send an out-of-band message, the MSG\_OOB flag is supplied to a send or a sendto system call. To receive out-of-band data, an application program must set the MSG\_OOB flag when performing a recvfrom or recv system call.

An application program can determine if the read pointer is currently pointing to the mark in the data stream by using the the SIOCATMARK ioctl:

```
ioctl(s, SIOCATMARK, &yes);
```

If yes is a 1 on return, meaning that no out-of-band data arrived, the next read returns data after the mark. If out-of-band data did arrive, the next read provides data sent by the client prior to transmission of the out-of-band signal. The following program shows the routine used in the remote login process to flush output on receipt of an interrupt or quit signal. This program reads the normal data up to the mark (to discard it), then reads the out-of-band byte:

```
#include <sys/ioctl.h>
#include <sys/file.h>

.
.
.
oob()
{
  int out = FWRITE, mark;
  char waste[BUFSIZ];

  /* flush local terminal output */
  ioctl(1, TIOCFLUSH, (char *)&out);
```

Sockets 4-45

```
for (;;) {
    if (ioctl(rem, SIOCATMARK, &mark) < 0) {
        perror("ioctl");
        break;
    }
    if (mark)
        break;
    (void) read(rem, waste, sizeof (waste));
}
if (recv(rem, &mark, 1, MSG_OOB) < 0) {
    perror("recv");
    .
    .
    .
}</pre>
```

A process can also read or peek at the out-of-band data without first reading up to the logical mark. This is difficult when the underlying protocol delivers the urgent in-band data with the normal data and only sends notification of its presence ahead of time; for example, the TCP protocol. With such protocols, when the out-of-band byte has not yet arrived and a recv system call is done with the MSG\_OOB flag, the call returns an EWOULDBLOCK error. There can be enough in-band data in the input buffer so that normal flow control prevents the peer from sending the urgent data until the buffer is cleared. The process must then read enough of the queued data so that the urgent data can be delivered.

#### Note

Certain programs that use multiple bytes of urgent data and must handle multiple urgent signals need to retain the position of urgent data within the stream. The socket-level SO\_OOBINLINE option provides this capability and Digital strongly recommends that you use it.

The socket-level SO\_OOBINLINE option retains the position of the urgent data (the logical mark). The urgent data immediately follows the mark within the normal data stream that is returned without the MSG\_OOB flag. Reception of multiple urgent indications causes the mark to move, but no out-of-band data is lost.

## 4.6.4 Internet Protocol Multicasting

Internet Protocol (IP) multicasting provides applications with IP layer access to the multicast capability of Ethernet and Fiber Distribution Data Interface (FDDI) networks. IP multicasting, which delivers datagrams on a best-effort basis, avoids the overhead imposed by IP broadcasting (described in Section 4.6.5) on uninterested hosts; it also avoids consumption of network bandwidth by applications that would otherwise transmit separate packets

with identical data to reach several destinations.

IP multicasting achieves efficient multipoint delivery through use of **host groups**. A host group is a group of zero or more hosts that is identified by a single Class D IP destination address. A Class D address has 1110 in the four high-order bits. In dotted decimal notation, IP multicast addresses range from 224.0.0.0 to 239.255.255.255, with 224.0.0.0 being reserved.

A member of a particular host group receives a copy of all data sent to the IP address representing that host group. Host groups can be permanent or transient. A permanent group has a well-known, administratively assigned IP address. In permanent host groups, it is the address of the group that is permanent, not its membership. The number of group members can fluctuate, even dropping to zero. The all hosts group group is an example of a permanent host group whose assigned address is 224.0.0.1. Digital UNIX systems join the all hosts group to participate in the Internet Group Management Protocol (IGMP). (See Request for Comments 1112: Host Extensions for IP Multicasting for more information about IGMP and IP multicasting.)

IP addresses that are not reserved for permanent host groups are available for dynamic assignment to transient groups. Transient groups exist only as long as they have one or more members.

#### Note

IP multicasting is not supported over connection-oriented transports such as TCP.

IP multicasting is implemented using options to the setsockopt system call, described in the following sections. Definitions required for multicast-related socket options are in the <netinet/in.h> header file. Your application must include this header file if you intend it to receive IP multicast datagrams.

## 4.6.4.1 Sending IP Multicast Datagrams

To send IP multicast datagrams, an application indicates the host group to send to by specifying an IP destination address in the range of 224.0.0.0 to 239.255.255.255 in a sendto system call. The system maps the specified IP destination address to the appropriate Ethernet or FDDI multicast address prior to transmitting the datagram.

An application can explicitly control multicast options with arguments to the setsockopt system call. The following options can be set by an application using the setsockopt system call:

• Time-to-live field (IP MULTICAST TTL)

- Multicast interface ( IP MULTICAST IF )
- Disabling loopback of local delivery (IP MULTICAST LOOP)

#### Note

The syntax for and arguments to the setsockopt system call are described in Section 4.3.5 and the setsockopt(2) reference page. The examples here and in Section 4.6.4.2 illustrate how to use the setsockopt options that apply to IP multicast datagrams only.

The IP\_MULTICAST\_TTL option to the setsockopt system call allows an application to specify a value between 0 and 255 for the time-to-live (TTL) field. Multicast datagrams with a TTL value of 0 restrict distribution of the multicast datagram to applications running on the local host. Multicast datagrams with a TTL value of 1 are forwarded only to hosts on the local subnet. If a multicast datagram has a TTL value greater than 1 and a multicast router is attached to the sending host's network, then multicast datagrams can be forwarded beyond the local subnet. Multicast routers forward the datagram to known networks that have hosts belonging to the specified multicast group. The TTL value is decremented by each multicast router in the path. When the TTL value is decremented to 0, the datagram is not forwarded further.

The following example shows how to use the <code>IP\_MULTICAST\_TTL</code> option to the <code>setsockopt</code> system call:

A datagram addressed to an IP multicast destination is transmitted from the default network interface unless the application specifies that an alternate network interface is associated with the socket. The default interface is determined by the interface associated with the default route in the kernel routing table or by the interface associated with an explicit route, if one exists. Using the IP\_MULTICAST\_IF option to the setsockopt system call, an application can specify a network interface other than that specified by the route in the kernel routing table.

The following example shows how to use the IP\_MULTICAST\_IF option to the setsockopt system call to specify an interface other than the default:

If a multicast datagram is sent to a group of which the sending host is a member, a copy of the datagram is, by default, looped back by the IP layer for local delivery. The IP\_MULTICAST\_LOOP option to the setsockopt system call allows an application to disable this loopback delivery.

The following example shows how to use the IP\_MULTICAST\_LOOP option to the setsockopt system call:

When the value of loop is 0, loopback is disabled. When the value of loop is 1, it is enabled. For performance reasons, Digital recommends disabling the default, unless applications on the same host must receive copies of the datagrams.

#### 4.6.4.2 Receiving IP Multicast Datagrams

Before a host can receive IP multicast datagrams destined for a particular multicast group other than the all hosts group, an application must direct the host to become a member of that multicast group. This section describes how an application can direct a host to add itself to and remove itself from a multicast group.

An application can direct the host it is running on to join a multicast group by using the IP\_ADD\_MEMBERSHIP option to the setsockopt system call as follows:

The *mreq* variable has the following structure:

```
struct ip_mreq{
    struct in_addr imr_multiaddr; /* IP multicast address of group */
    struct in_addr imr_interface; /* local IP address of interface */
};
```

Each multicast group membership is associated with a particular interface. It is possible to join the same group on multiple interfaces. The imr\_interface variable can be specified as INADDR\_ANY, which allows an application to choose the default multicast interface. Alternatively, specifying one of the host's local addresses allows an application to select a particular, multicast-capable interface. The maximum number of memberships that can be added on a single socket is subject to the IP\_MAX\_MEMBERSHIPS value, which is defined in the <netinet/in.h> header file.

To drop membership in a particular multicast group use the IP\_DROP\_MEMBERSHIP option to the setsockopt system call:

The *mreq* variable contains the same structure values used for adding membership.

If multiple sockets request that a host join a particular multicast group, the host remains a member of that multicast group until the last of those sockets is closed.

To receive multicast datagrams sent to a specific UDP port, the receiving socket must have bound to that port using the bind system call. More than one process can receive UDP datagrams destined for the same port if the bind system call (described in Section 4.3.2) is preceded by a setsockopt system call that specifies the SO\_REUSEPORT option. The following example illustrates how to use the SO\_REUSEPORT option to the setsockopt system call:

When the SO\_REUSEPORT option is set, every incoming multicast or broadcast UDP datagram destined for the shared port is delivered to all sockets bound to that port.

Delivery of IP multicast datagrams to SOCK\_RAW sockets is determined by the protocol type of the destination.

## 4.6.5 Broadcasting and Determining Network Configuration

Using a datagram socket, it is possible to send broadcast packets on many networks supported by the system. The network itself must support broadcast; the system provides no simulation of broadcast in the software.

Broadcast messages can place a high load on a network because they force every host on the network to service them. Consequently, the ability to send broadcast packets is limited to sockets that are explicitly marked as allowing broadcasting.

Broadcast is typically used for one of two reasons: to find a resource on a local network without prior knowledge of its address, or to route some information, which requires that information be sent to all accessible neighbors.

#### Note

Broadcasting is not supported over connection-oriented transports such as TCP.

To send a broadcast message, use the following procedure:

1. Create a datagram socket; for example:

```
s = socket(AF INET, SOCK DGRAM, 0);
```

2. Mark the socket for broadcasting; for example:

3. Ensure that at least a port number is bound to the socket; for example:

```
sin.sin_len = sizeof(sin);
sin.sin_family = AF_INET;
sin.sin_addr.s_addr = htonl(INADDR_ANY);
sin.sin_port = htons(MYPORT);
if (bind(s, (struct sockaddr *) &sin, sizeof (sin)) == -1)
    perror("setsockopt");
```

The destination address of the message depends on the network or networks on which the message is to be broadcast. The Internet domain supports a shorthand notation for broadcast on the local network, the address is INADDR\_BROADCAST (as defined in netinet/in.h).

To determine the list of addresses for all reachable neighbors requires knowledge of the networks to which the host is connected. Digital UNIX provides a method of retrieving this information from the system data structures. The SIOCGIFCONF ioctl call returns the interface

configuration of a host in the form of a single ifconf structure. This structure contains a data area that an array of ifreq structures comprises, one for each network interface to which the host is connected. These structures are defined in the <net/if.h> header file, as follows:

```
struct ifconf {
        ifc_len;
                                   /* size of associated buffer */
   int
   union {
     caddr_t ifcu_buf;
      struct ifreq *ifcu_req;
   } ifc_ifcu;
#define ifc_buf ifc_ifcu.ifcu_buf /* buffer address */
#define ifc_req ifc_ifcu.ifcu_req /* array of structures returned */
struct ifreq {
#define IFNAMSIZ
                        16
  char ifr_name[IFNAMSIZ];
                                   /* if name, e.g. "en0" */
   union {
      struct sockaddr ifru_addr;
      struct sockaddr ifru_dstaddr;
      struct sockaddr ifru_broadaddr;
      short ifru_flags;
      int
             ifru_metric;
      caddr_t ifru_data;
   } ifr_ifru;
#define ifr_addr ifr_ifru.ifru_addr /* address */
#define ifr_dstaddr ifr_ifru.ifru_dstaddr /* other end of */
                                              /* p-to-p link */
#define ifr_broadaddr ifr_ifru.ifru_broadaddr /* broadcast address */
#define ifr_flags ifr_ifru.ifru_flags /* flags */
                                              /* metric */
                     ifr_ifru.ifru_metric
#define ifr_metric
                                             /* for use by */
#define ifr_data
                    ifr_ifru.ifru_data
                                              /* interface */
```

The actual call which obtains the interface configuration is as follows:

```
struct ifconf ifc;
char buf[BUFSIZ];

ifc.ifc_len = sizeof (buf);
ifc.ifc_buf = buf;
if (ioctl(s, SIOCGIFCONF, (char *) &ifc) < 0) {
    .
    .
    .
}</pre>
```

After this call, buf contains one ifreq structure for each network to which the host is connected, and ifc.ifc\_len is modified to reflect the number of bytes used by the ifreq structures.

Each structure has a set of interface flags that tells whether the network corresponding to that interface flag is up or down, point-to-point or broadcast, and so on. The SIOCGIFFLAGS ioctl retrieves these flags for

an interface specified by an ifreq structure, as follows:

Once the flags are obtained, the broadcast address must be obtained. In the case of broadcast networks, this is done via the SIOCGIFBRDADDR ioctl; while, for point-to-point networks, the address of the destination host is obtained with SIOCGIFDSTADDR. For example:

```
struct sockaddr dst;

if (ifr->ifr_flags & IFF_POINTOPOINT) {
    if (ioctl(s, SIOCGIFDSTADDR, (char *) ifr) < 0) {
        ...
    }
    bcopy((char *) ifr->ifr_dstaddr, (char *) &dst,
        sizeof (ifr->ifr_dstaddr));
} else if (ifr->ifr_flags & IFF_BROADCAST) {
    if (ioctl(s, SIOCGIFBRDADDR, (char *) ifr) < 0) {
        ...
    }
    bcopy((char *) ifr->ifr_broadaddr, (char *) &dst,
        sizeof (ifr->ifr_broadaddr));
}
```

After the appropriate ioctl calls obtain the broadcast or destination address (now in dst), the sendto call is used; for example:

```
if (sendto(s, buf, buflen, 0, (struct sockaddr *)&dst, sizeof (dst)) < 0)
    perror("sendto");</pre>
```

In the preceding loop, one sendto call occurs for every interface to which the host is connected that supports the notion of broadcast or point-to-point addressing. If a process only wants to send broadcast messages on a given network, code similar to that in the preceding example is used, but the loop needs to find the correct destination address.

#### 4.6.6 The inetd Daemon

Digital UNIX supports the inetd Internet superserver daemon. The inetd daemon, which is invoked at boot time, reads the /etc/inetd.conf file to determine the servers for which it should listen.

#### Note

Only server applications written to run over sockets can use the inetd daemon in Digital UNIX. The inetd daemon in Digital UNIX does not support server applications written to run over STREAMS, XTI, or TLI.

For each server listed in /etc/inetd.conf the inetd daemon does the following:

- 1. Creates a socket and binds the appropriate port number to it.
- 2. Issues a select system call for read availability and waits for a process to request a connection to the service that corresponds to that socket.
- 3. Issues an accept system call, forks, duplicates (with the dup call) the new socket to file descriptors 0 and 1 (stdin and stdout), closes other open file descriptors, and executes (with the exec call) the appropriate server.

Servers that use inetd are simplified because inetd takes care of most of the interprocess communication work required to establish a connection. The server invoked by inetd expects the socket connected to its client on file descriptors 0 and 1, and immediately performs any operations such as read, write, send, or recv.

Servers invoked by the inetd daemon can use buffered I/O as provided by the conventions in the <stdio.h> header file, as long as as they remember to use the fflush call when appropriate. See fflush(3) for more information.

The getpeername call, which returns the address of the peer (process) connected on the other end of the socket, is useful for developers writing server applications that use inetd. The following sample code shows how

to log the Internet address, in dot notation, of a client connected to a server under inetd:

```
struct sockaddr_in name;
size_t namelen = sizeof (name);

.
.
if (getpeername(0, (struct sockaddr *)&name, &namelen) < 0) {
    syslog(LOG_ERR, "getpeername: %m");
    exit(1);
} else
    syslog(LOG_INFO, "Connection from %s", inet_ntoa(name.sin_addr));
.
.</pre>
```

While the getpeername call is especially useful when writing programs to run with inetd, it can be used under other circumstances.

## 4.6.7 Input/Output Multiplexing

Multiplexing is a facility used in applications to transmit and receive I/O requests among multiple sockets. This can be done by using the select call, as follows:

The select call takes as arguments pointers to three sets:

- 1. The set of socket descriptors for which the calling application wants to read data.
- 2. The socket descriptors to which data is to be written.
- 3. Exceptional conditions which are pending.

The corresponding argument to the select call must be a null pointer, if the application is not interested in certain conditions; for example, read, write, or exceptions.

#### Note

Because XTI and TLI are implemented using STREAMS, you should use the poll system call instead of the select system call on any STREAMS file descriptors.

Each set is actually a structure that contains an array of integer bit masks. The size of the array is set by the FD\_SETSIZE definition. The array is long enough to hold one bit for each of the FD\_SETSIZE file descriptors.

The FD\_SET (fd, &mask) and FD\_CLR (fd, &mask) macros are provided to add and remove the fd file descriptor in the mask set. The set needs to be zeroed before use and the FD\_ZERO (&mask) macro is provided to clear the mask set.

The *nfds* parameter in the select call specifies the range of file descriptors (for example, one plus the value of the largest descriptor) to be examined in a set.

A time-out value can be specified when the selection will not last more than a predetermined period of time. If the fields in timeout are set to zero (0), the selection takes the form of a poll, returning immediately. If the last parameter is a null pointer, the selection blocks indefinitely. Specifically, a return takes place only when a descriptor is selectable or when a signal is received by the caller, interrupting the system call.

The select call normally returns the number of file descriptors selected; if the select call returns because the time-out expired, then the value 0 is returned. If the select call terminates because of an error or interruption, a -1 is returned with the error number in errno and with the file descriptor masks unchanged.

Assuming a successful return, the three sets indicate which file descriptors are ready to be read from, written to, or have exceptional conditions pending. The status of a file descriptor in a select mask can be tested with the FD\_ISSET (fd, &mask) macro, which returns a nonzero value if fd is a member of the mask set or 0 if it is not.

To determine whether there are connections waiting on a socket to be used with an accept call, the select call is used, followed by a FD\_ISSET (fd, &mask) macro to check for read readiness on the appropriate socket. If FD\_ISSET returns a nonzero value, indicating data to read, then a connection is pending on the socket.

#### Note

In 4.2BSD, the arguments to the select call were pointers to integers instead of pointers to fd\_set. This type of call works as long as the number of file descriptors being examined is less than the number of bits in an integer; however, the method shown in the following code is recommended.

The following example shows how an application reads data as it becomes available from sockets \$1 and \$2 with a 1-second time-out:

```
#include <sys/time.h>
#include <sys/types.h>
fd_set read_template;
struct timeval wait;
for (;;) {
  wait.tv_sec = 1;  /* one second */
  wait.tv_usec = 0;
  FD_ZERO(&read_template);
  FD_SET(s1, &read_template);
  FD_SET(s2, &read_template);
  nb = select(FD_SETSIZE, &read_template, (fd_set *) 0,
     (fd_set *) 0, &wait);
   if (nb <= 0) {
     An error occurred during the select, or
        the select timed out
   if (FD_ISSET(s1, &read_template)) {
      Socket #1 is ready to be read from.
   if (FD_ISSET(s2, &read_template)) {
      Socket #2 is ready to be read from.
}
```

The select call provides a synchronous multiplexing scheme. Asynchronous notification of output completion, input availability, and exceptional conditions is possible through use of the SIGIO and SIGURG signals described in Section 4.6.9.

## 4.6.8 Interrupt Driven Socket I/O

The SIGIO signal allows a process to be notified using a signal when a socket (or more generally, a file descriptor) has data waiting to be read. Using the SIGIO facility requires the following three steps:

- 1. The process must set up a SIGIO signal handler by using the signal or sigvec calls.
- 2. The process must set the process ID or process group ID that is to receive notification of pending input to its own process ID or the process group ID of its process group. (Note that the default process group of a socket is group 0.) This is done by using a fcntl system call.
- 3. The process must enable asynchronous notification of pending I/O requests with another fcntl system call. The following code shows how to allow a particular process to receive information on pending I/O requests as they occur for socket s. With the addition of a handler for SIGURG, this code can also be used to prepare for receipt of SIGURG signals.

```
#include <fcntl.h>
    .
    .
    int io_handler();
    .
    signal(SIGIO, io_handler);

/* Set the process receiving SIGIO/SIGURG signals to us */
if (fcntl(s, F_SETOWN, getpid()) < 0) {
    perror("fcntl F_SETOWN");
    exit(1);
}

/* Allow receipt of asynchronous I/O signals */
if (fcntl(s, F_SETFL, FASYNC) < 0) {
    perror("fcntl F_SETFL, FASYNC");
    exit(1);
}</pre>
```

## 4.6.9 Signals and Process Groups

Each socket has an associated process number, the value of which is initialized to zero (0). This number must be redefined with the F\_SETOWN parameter to the fcntl system call, as was done in Section 4.6.8, to enable SIGURG and SIGIO signals to be caught. To set the socket's process ID for signals, positive arguments must be given to the fcntl call. To set the

socket's process group for signals, negative arguments must be passed to the fcntl call. Note that the process number indicates the associated process ID or the associated process group; it is impossible to specify both simultaneously.

The F\_GETOWN parameter to the fcntl call allows a process to determine the current process number of a socket.

The SIGCHLD signal is also useful when constructing server processes. This signal is delivered to a process when any child processes change state. Typically, servers use the SIGCHLD signal to reap child processes that exited, without explicitly awaiting their termination or periodic polling for exit status. If the parent server process fails to reap its children, a large number of zombie processes may be created. The following code shows how to use the SIGCHLD signal:

```
int reaper();
signal(SIGCHLD, reaper);
listen(f, 5);
for (;;) {
   int g;
   size_t len = sizeof (from);
   g = accept(f, (struct sockaddr *)&from, &len,);
   if (q < 0) {
      if (errno != EINTR)
         syslog(LOG_ERR, "rlogind: accept: %m");
      continue;
}
#include <wait.h>
reaper()
   union wait status;
   while (wait3(&status, WNOHANG, 0) > 0)
}
```

## 4.6.10 Pseudoterminals

Many programs cannot function properly without a terminal for standard input and output. Since sockets do not provide the semantics of terminals, it is often necessary to have a process communicating over the network do so

through a pseudoterminal (pty). A pseudoterminal is a pair of devices, master and slave, that allow a process to serve as an active agent in communication between applications and users.

Data written on the slave side of a pseudoterminal is used as input to a process reading from the master side, while data written on the master side is processed as terminal input for the slave. In this way, the process manipulating the master side of the pseudoterminal controls the information read and written on the slave side as if it were manipulating the keyboard and reading the screen on a real terminal. The purpose of the pseudoterminal abstraction is to preserve terminal semantics over a network connection; that is, the slave side appears as a normal terminal to any process reading from or writing to it.

For example, rlogind, the remote login server uses pseudoterminals for remote login sessions. A user logging in to a machine across the network is provided a shell with a slave pseudoterminal as standard input, standard output, and standard error. The server process then handles the communication between the programs invoked by the remote shell and the user's local client process. When a user sends a character that generates an interrupt on the remote machine that flushes terminal output, the pseudoterminal generates a control message for the server process. The server then sends an out-of-band message to the client process to signal a flush of data at the real terminal and on the intervening data buffered in the network.

In Digital UNIX, the slave side of a pseudoterminal has a name of the form /dev/ttyxy, where x is any single letter, except d, and is uppercase or lowercase. The y is a hexadecimal digit, meaning it is a single character in the range of 0 to 9 or a to f. The master side of a pseudoterminal has a name of the form /dev/ptyxy, where x and y correspond to x and y on the slave side of the pseudoterminal.

The openpty and forkpty functions were added to the libc.a library to make allocating pseudoterminals easier. These functions use the clone open call to avoid performing multiple open calls.

The following is the syntax for the openpty and forkpty functions:

```
#include <termios.h>
#include <ioctl.h>
:
:
:
int openpty(
    int *master,
    int *slave,
    char *name,
    struct termios *termp,
    struct winsize *winp,);
```

```
pid_t forkpty(
    int *master,
    char *name,
    struct termios *termp,
    struct winsize *winp,);
```

The first two arguments of the openpty function are pointers to integers which, upon successful completion, hold the value of the master and slave file descriptors respectively.

The last three arguments are optional; you should specify them as NULL if they are not used. If they are used, they do the following:

- The third argument is a pointer to a character string which is the pathname of the slave device.
- The fourth argument is a pointer to a termios structure and is used to set the slave's terminal characteristics.
- The fifth argument is pointer to a winsize structure which sets the window size of the slave.

The forkpty function allocates a pseudoterminal. Additionally, it forks a child process and makes the slave pseudoterminal the controlling terminal for the child. The forkpty function takes four arguments instead of five, because the slave file descriptor is not passed back to the calling process. Instead, the slave file descriptor is duplicated in the newly created child process as stdin, stdout, and stderr. The other four arguments are identical to those of the openpty function.

Both the openpty and forkpty functions return -1 to signify an error condition. The openpty function returns a zero (0) upon successful completion, while the forkpty returns the pid of the child process. See the openpty(3) reference page for more information.

The openpty function works as follows:

- 1. Upon successful completion, the slave side of the pseudoterminal is set to the proper terminal modes. At the time the master and slave sides of the pseudoterminal are opened, Digital UNIX performs the necessary security checks.
- 2. The process then forks; the child closes the master side of the pseudoterminal and executes (with the exec call) the appropriate program.
- 3. The parent closes the slave side of the pseudoterminal and begins reading and writing from the master side.

The following example makes use of pseudoterminal. The code in this example makes the following assumptions:

- A connection on a socket already exists.
- The socket is connected to a peer that wants a service of some kind.
- The process disassociated itself from any previous controlling terminal.

```
if (openpty(&mast,&slave,NULL,NULL,NULL) {
   syslog(LOG_ERR, "All network ports in use");
   exit(1);
ioctl(slave, TIOCGETA, &term); /* get default slave termios struct */
term.c_iflag |= ICRNL;
term.c_oflag |= OCRNL;
ioctl(slave, TIOCSETA, &term); /* set slave characteristics
                                                                        * /
i = fork();
if (i < 0) {
   syslog(LOG_ERR, "fork: %m");
   exit(1);
                   /* Parent */
} else if (i) {
   close(slave);
             /* Child */
} else {
   (void) close(s);
   (void) close(master);
   dup2(slave, 0);
   dup2(slave, 1);
dup2(slave, 2);
   if (slave > 2)
      (void) close(slave);
```

See Section 4.3 for information about using sockets.

# Digital UNIX STREAMS 5

Digital UNIX provides a STREAMS framework as specified by AT&T's System V, Version 4.0 release of STREAMS. This framework, which provides an alternative to traditional UNIX character input/output (I/O), allows you to implement I/O functions in a modular fashion. Modularly developed I/O functions allow applications to build and reconfigure communications services easily.

Note that STREAMS refers to the entire framework whereas Stream refers to the entity created by an application program with the open system call.

This chapter contains the following information:

- Overview of the STREAMS framework
- Description of the application interface to STREAMS
- Description of the kernel-level functions
- Instructions on how to configure modules or drivers
- Description of the Digital UNIX synchronization mechanism
- Information on how to create device special files
- Description of error and event logging
- Information about STREAMS reference pages

This chapter provides detailed information about areas where the Digital UNIX implementation of STREAMS differs from that of AT&T System V, Version 4.0. Where the Digital UNIX implementation does not differ significantly from that of AT&T, it provides pointers to the appropriate AT&T documentation.

Note that this chapter does not explain how to program using the STREAMS framework. For detailed programming information you should refer to the *Programmer's Guide: STREAMS*.

#### 5.1 Overview of the STREAMS Framework

The STREAMS framework consists of:

• A programming interface, or set of system calls, used by application programs to access the STREAMS framework

- Kernel resources, such as the Stream head, and queue data structures used by the Stream
- Kernel utilities that handle tasks such as Stream queue scheduling and flow control, memory allocation, and error logging

Figure 5-1 highlights the STREAMS framework and shows its place in the network programming environment.

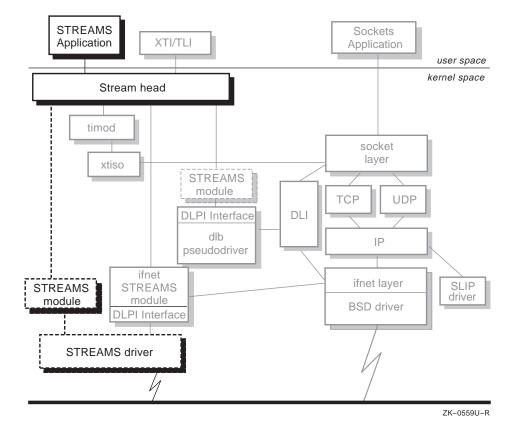


Figure 5-1: The STREAMS Framework

## 5.1.1 A Review of STREAMS Components

To communicate using Digital UNIX STREAMS, an application creates a Stream, which is a full-duplex communication path between a user process and a device driver. The Stream itself is a kernel device and is represented to the application as a character special file. Like any other character special file, the Stream must be opened and otherwise manipulated with system calls.

Every Stream has at least a Stream head at the top and a Stream end at the bottom. Additional modules, which consist of linked pairs of queues, can be inserted between the Stream head and Stream end if they are required for processing the data being passed along the Stream. Data is passed between modules in messages.

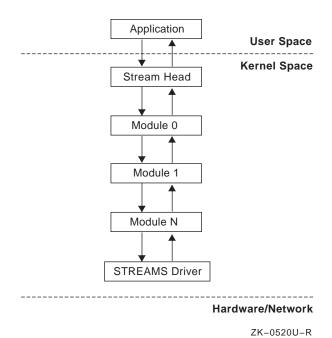
This section briefly describes the following STREAMS components:

- Stream head
- Stream end
- Modules

It also describes messages and their role in the STREAMS framework.

Figure 5-2 illustrates a typical stream. Note that data traveling from the Stream head to the Stream end (STREAMS driver in Figure 5-2) is said to be traveling downstream, or in the write direction. Data traveling from the Stream end to the Stream head is said to be traveling upstream, or in the read direction.

Figure 5-2: Example of a Stream



The Stream head is a set of routines and data structures that provides an interface between user processes and the Streams in the kernel. It is created

when your application issues an open system call. The following are the major tasks that the Stream head performs:

- 1. Interprets a standard subset of STREAMS system calls, such as write and putmsq.
- 2. Translates them from user space into a standard range of STREAMS messages (such as M\_PROTO and M\_DATA) which consist of both data and control information.
- 3. Sends the messages downstream to the next module. Eventually the messages reach the Stream end, or driver.
- 4. Receives messages sent upstream from the driver and transforms the STREAMS message from kernel space to a format appropriate to the system call (such as getmsg or read) made by the application. The format varies depending on the system call.

The Stream end is a special form of STREAMS module and can be either a hardware or pseudodevice driver. If a hardware device driver, the Stream end provides communication between the kernel and an external communication device. If a pseudodevice driver, the Stream end is implemented in software and is not related to an external device. Regardless of whether it is a hardware device driver or a pseudodevice driver, the Stream end receives messages sent by the module above it, interprets them, and performs the requested operations. It then returns data and control information to the application by creating a message of the appropriate type which it sends upstream toward the Stream head.

Drivers are like any other STREAMS modules except for the following:

- They can handle interrupts (although they do not have to).
   Device drivers can have one or more interrupt routines. Interrupt routines should queue data on the read side service routine for later processing.
- They can be connected to multiple Streams.
   A driver can be implemented as a multiplexor, meaning that it is connected to multiple Streams in either the upstream or downstream direction. See the *Programmer's Guide: STREAMS* for more information.
- They are initialized and deinitialized by the open and close system calls. (Other modules use the I\_PUSH and I\_POP commands of the ioctl system call.)

For detailed information on device drivers and device driver routines, see the *Writing Device Drivers: Tutorial* and the *Programmer's Guide: STREAMS*.

Modules process data as it passes from the Stream head to the Stream end and back. A Stream can have zero or more modules on it, depending on the amount and type of processing that the data requires. If the driver can perform all of the necessary processing on the data, no additional modules are required.

Modules consist of a pair of queues that contain data and pointers to other structures that define what each module does. One queue handles data moving downstream toward the driver and the other handles data moving upstream toward the Stream head and application. Pointers link each module's downstream and upstream queues to the next module's downstream and upstream queues.

Depending on their processing requirements, applications request that particular modules be pushed onto the Stream. The Stream head assembles the modules requested by the application and then routes the messages through the pipeline of modules.

Information is passed from module to module using messages. Several different types of messages are defined within the STREAMS environment. All message types, however, fall into the following categories:

- Normal
- High priority

Normal messages, such as M\_DATA and M\_IOCTL, are processed in the order that they are received, and are subject to STREAMS flow control and queuing mechanisms. Priority messages are passed along the stream in an expedited manner.

For more information on messages and message data structures, see Section 5.3.2

# 5.2 Application Interface to STREAMS

The application interface to the STREAMS framework allows STREAMS messages to be sent and received by applications. The following sections describe the application interface, including pointers to the STREAMS header files and data types, and descriptions of the STREAMS and STREAMS-related system calls.

# 5.2.1 Header Files and Data Types

Definitions for the basic STREAMS data types are included in the following header files:

- The <sys/stream.h> header file must be included for all modules and Streams applications.
- The <stropts.h> header file must be included when an application uses the ioctl system call.

• The <strlog.h> header file must be included when an application uses the STREAMS error logger and trace facility.

# Note

Typically, header file names are enclosed in angle brackets (< >). To obtain the absolute path to the header file, prepend /usr/include/ to the information enclosed in the angle brackets. In the case of <sys/stream.h>, stream.h is located in the /usr/include/sys directory.

# 5.2.2 STREAMS Functions

Your application accesses and manipulates STREAMS kernel resources through the following functions:

- open
- close
- read
- write
- ioctl
- mkfifo
- pipe
- putmsg and putpmsg
- getmsg and getpmsg
- poll
- isastream
- fattach
- fdetach

This section briefly describes these functions. For detailed information about these functions, see the Digital UNIX reference pages and the *Programmer's Guide: STREAMS*.

# 5.2.2.1 The open Function

Use the open function to open a Stream. The following is the syntax for the open function:

```
int open (
     const char *path,
     int oflag[ ,
     mode_t mode] );
```

In the preceding statement:

path

Specifies the device pathname supplied to the open function. The device pathnames are located in the /dev/streams directory. To determine which devices are configured on your system issue the following command as root:

```
# /usr/sbin/strsetup -c
```

oflag

Specifies the type of access, special open processing, type of update, and the initial state of the open file.

mode

Specifies the permissions of the file that open is creating.

See open(2) for more information.

The following example shows how the open function is used:

```
int fd;
fd = open("/dev/streams/echo", O_RDWR);
```

# 5.2.2.2 The close Function

Use the close function to close a Stream.

The following is the syntax for the close function:

```
int close(
    int filedes);
```

In the preceding statement:

filedes

Specifies a valid open file descriptor

See close(2) for more information.

The last close for a stream causes the stream associated with the file descriptor to be dismantled. Dismantling a stream includes popping any modules on the stream and closing the driver.

# 5.2.2.3 The read Function

Use the read function to receive the contents of M\_DATA messages waiting at the Stream head.

The following is the syntax for the read function:

```
int read(
    int filedes,
    char *buffer,
    unsigned int nbytes);
```

In the preceding statement:

filedes

Specifies a file descriptor that identifies the file to be read.

\*buffer

Points to the buffer to receive the data being read.

nbytes

Specifies the number of bytes that can be read from the file associated with *filedes*.

See read(2) for more information.

The read function fails on message types other than M\_DATA, and errno is set to EBADMSG.

#### 5.2.2.4 The write Function

Use the write function to create one or more  $M_DATA$  messages from the data buffer.

The following is the syntax for the write function:

```
int write(
    int filedes,
    char *buffer
    unsigned int nbytes);
```

In the preceding statement:

filedes

Specifies a file descriptor that identifies the file to be read.

\*buffer

Points to the buffer to receive the data being read.

nbytes

Specifies the number of bytes to write from the file associated with *filedes*.

See write(2) for more information.

#### 5.2.2.5 The joctl Function

Use the ioctl function to perform a variety of control functions on Streams.

The following is the syntax of the ioctl function:

```
#include <stropts.h>
```

```
int ioctl ( filedes, command, arg)
int fildes, command;
```

In the preceding statement:

```
filedes
```

Specifies an open file descriptor that refers to a Stream.

command

Determines the control function for the Stream head or module to perform. Many of the valid ioctl commands are handled by the Stream head; others are passed downstream to be handled by the modules and driver.

arg

Specifies additional information. The type depends on the *command* parameter.

See streamio(7) for more information.

The following example shows how the ioctl call is used:

```
int fd;
fd = open("/dev/streams/echo", O_RDWR, 0);
ioctl(fd,I_PUSH,"pass");
```

# 5.2.2.6 The mkfifo Function

Use the STREAMS-based mkfifo function to create a unidirectional STREAMS-based file descriptor.

The following is the syntax of the STREAMS-based mkfifo function:

```
int mkfifo(
     const char *path,
     mode_t mode);
```

In the preceding statement:

path

Specifies the file name supplied to the mkfifo function.

mode

Specifies the type, attributes, and access permissions of the file.

#### Note

The default version of the mkfifo function in the libc library is not STREAMS-based. To use the STREAMS version of the mkfifo function the application must link with the sys5 library. See the mkfifo(2) reference page for more information.

Also note that the mkfifo function requires that the File on File Mount File System (FFM\_FS) kernel option is configured. See the *System Administration* manual for information about configuring kernel options.

# 5.2.2.7 The pipe Function

Use the STREAMS-based pipe function to create a bidirectional, STREAMS-based, communication channel. Non-STREAMS pipes and STREAMS-based pipes differ in the following ways:

- Non-STREAMS pipes are unidirectional
- STREAMS operations (such as streamio and putmsg) can not be performed on them

The following is the syntax of the pipe function:

```
int pipe(
  int filedes[2]);
```

In the preceding statement:

filedes Specifies the address of an array of two integers into which new file descriptors are placed.

#### Note

The default version of the pipe function in the libc library is not STREAMS-based. To use the STREAMS version of the pipe function the application must link with the sys5 library. See the pipe(2) reference page for more information.

# 5.2.2.8 The putmsg and putpmsg Functions

Use the putmsg and putpmsg functions to generate a STREAMS message block by using information from specified buffers.

The following is the syntax of the putmsg function:

```
int putmsg(
    int filedes;
    struct strbuf *ctlbuf;
    struct strbuf *databuf;
    int flags;)
```

In the preceding statement:

filedes

Specifies the file descriptor that references an open Stream.

ctlbuf

Points to a strbuf structure that holds the control part of the message.

databuf

Points to a strbuf structure that holds the data part of the message.

flags

An integer that specifies the type of message the application wants to send.

See putmsg(2) for more information.

Use the putpmsg function to send priority banded data down a Stream.

The following is the syntax of the putpmsg function:

```
int putpmsg(
    int filedes;
    struct strbuf *ctlbuf;
    struct strbuf *databuf;
    int band;
    int flags;)
```

The arguments have the same meaning as for the putmsg function. The band argument specifies the priority band of the message.

See putpmsg(2) for more information.

# 5.2.2.9 The getmsg and getpmsg Functions

Use the getmsg and getpmsg functions to retrieve the contents of a message located at the Stream head read queue and place them into user specified buffer(s).

The following is the syntax of the getmsg function:

```
int getmsg(
    int filedes
    struct strbuf *ctlbuf
    struct strbuf *databuf
    int *flags);
```

In the preceding statement:

filedes

Specifies a file descriptor that references an open Stream.

ctlbuf

Points to a strbuf structure that returns the control part of the message.

databuf

Points to a strbuf structure that returns the data part of the message.

flags

Points to an integer that specifies the type of message the application wants to retrieve.

See getmsg(2) for more information.

Use the getpmsg function to receive priority banded data from a Stream.

The following is the syntax of the getpmsg function:

```
int getpmsg(
    int filedes
    struct strbuf *ctlbuf
    struct strbuf *databuf
    int band;
    int *flags);
```

The arguments have the same meaning as for the getmsg function. The band argument points to an integer that specifies the priority band of the message being received.

See getpmsg(2) for more information.

# 5.2.2.10 The poll Function

Use the poll function to identify the Streams to which a user can send data and from which a user can receive data.

The following is the syntax for the poll function:

```
#include <sys/poll.h>
int poll(
         struct pollfd filedes[],
         unsigned int nfds,
```

In the preceding statement:

int timeout);

```
filedes
```

Points to an array of pollfd structures, one for each file descriptor you are polling. By filling in the pollfd structure, the caller can specify a set of events about which to be notified.

nfds

Specifies the number of pollfd structures in the filedes array.

timeout

Specifies the maximum length of time (in milliseconds) to wait for at least one of the specified events to occur.

See pol1(2) for more information.

#### 5.2.2.11 The isastream Function

Use the isastream function to determine if a file descriptor refers to a STREAMS file.

The following is the syntax for the isastream routine:

```
int isastream(
  int filedes;);
```

In the preceding statement:

filedes

Specifies an open file desciptor.

The following example shows how to use the isastream function to verify that you have opened a STREAMS-based pipe instead of a sockets-based pipe:

See the isastream(3) reference page for more information.

# 5.2.2.12 The fattach Function

Use the fattach function to attach a STREAMS-based file descriptor to an object in the file system name space.

The following is the syntax of the fattach function:

```
int fattach(
   int fd,
   const char *path);
```

In the preceding statement:

fd

Specifies an open STREAMS-based file descriptor.

path

Specifies the pathname of an existing file system object. The pathname must reference a regular file. It can not reference, for example, a directory or pipe.

The following example shows how to use the fattach function to name a STREAMS-based pipe:

```
int fds[2];
pipe(fds);
fattach(fd[0], "/tmp/pipe1");
```

# Note

The fattach function requires that the FFM\_FS kernel option be configured. See the *System Administration* manual for information about configuring kernel options.

See the fattach(3) reference page for more information.

# 5.2.2.13 The fdetach Function

Use the fdetach function to detach a STREAMS-based file descriptor from a file name. A STREAMS-based file descriptor may have been attached by using the fattach function.

The following is the syntax of the fdetach function:

```
int fdetach(
    const char *path);
```

In the preceding statement:

path

Specifies the pathname of a file system object that was previously attached.

# Note

The fdetach function requires that the File on File Mount File System (FFM\_FS) kernel option is configured. See the *System Administration* manual for information about configuring kernel options.

See the fdetach(3) reference page for more information.

Table 5-1 lists and briefly describes the reference pages that contain STREAMS-related information. For further information about each component, refer to the appropriate reference page.

Table 5-1: STREAMS Reference Pages

Reference Page	Description	
autopush(8)	Command that manages the system's database of automatically pushed STREAMS modules.	
clone(7)	STREAMS software driver that finds and opens an unused major/minor device on another STREAMS driver.	
* close(2)	Function that closes the file associated with a designated file descriptor.	
dlb(7)	STREAMS pseduodevice driver that provides a communication path between BSD-style device drivers and STREAMS protocol stacks.	
fattach(8)	Command that attaches a STREAMS-based file descriptor to a node in the file system.	
fdetach(8)	Command that detaches a STREAMS-based file descriptor from a file name.	
fdetach(3)	Function that detaches a STREAMS-based file descriptor from a file name.	
getmsg(2) getpmsg(2)	Functions that reference a message positioned at the Stream head read queue.	

Table 5-1: (continued)

Reference Page	Description		
ifnet(7)	STREAMS-based module that provides a bridge between STREAMS-based device drivers written to the Data Link Provider Interface (DLPI) and sockets.		
isastream(3)	Function that determines if a file descriptor refers to a STREAMS file.		
mkfifo(2)	Function that creates a unidirectional STREAMS-based file descriptor.		
* open(2)	Function that establishes a connection between a file and a file descriptor.		
pipe(2)	Function that creates a bidirectional, STREAMS-based, interprocess communication channel.		
pol1(2)	Function that provides a general mechanism for reporting I/O conditions associated with a set of file descriptors and for waiting until one or more specified conditions becomes true.		
<pre>putmsg(2) putpmsg(2)</pre>	Functions that generate a STREAMS message block.		
* read(2)	Function that reads data from a file into a designated buffer.		
strace(8)	Application that retrieves STREAMS event trace messages from the STREAMS log driver.		
strchg(1)	Command that alters the configuration of a Stream.		
strclean(8)	Command that removes STREAMS error log files.		
strconf(1)	Command that queries about a Stream's configuration.		
streamio(7)	Command that performs a variety of control functions on Streams.		
strerr(8)	Daemon that receives error messages from the STREAMS log driver.		
strlog(7)	Interface that tracks log messages used by STREAMS error logging and event tracing daemons.		
strsetup(8)	Command that creates the appropriate STREAMS pseudodevices and displays the setup of your STREAMS modules.		
timod(7)	Module that converts ioctl calls from a transport user supporting the Transport Interface (TI) into messages that a transport protocol provider supporting TI can consume.		

Table 5-1: (continued)

Reference Page	Description
tirdwr(7)	Module that provides a transport user supporting the TI with an alternate interface to a transport protocol provider supporting TI.
* write(2)	Function that writes data to a file from a designated buffer.

Table Notes: An asterisk (\*) means that the page is not STREAMS specific.

# 5.3 Kernel Level Functions

This section contains information with which the kernel programmer who writes STREAMS modules and drivers must be familiar. It contains information about:

- Module data structures
- Message data structures
- STREAMS processing routines for modules and drivers

#### 5.3.1 Module Data Structures

When a module or driver is configured into the system, it must define its read and write queues and other module information.

The qinit, module\_info, and streamtab data structures, all of which are located in the <sys/stream.h> header file, define read and write queues. STREAMS modules must fill in these structures in their declaration sections. See Appendix A for an example.

The only external data structure a module must provide is streamtab.

The qinit structure, shown in the following example, defines the interface routines for a queue. The read queue and write queue each have their own set of structures.

The module\_info structure, shown in the following example, contains module or driver identification and limit values:

The streamtab structure, shown in the following example, forms the uppermost part of the declaration and is the only part which needs to be visible outside the module or driver:

# 5.3.2 Message Data Structures

Digital UNIX STREAMS messages consist of one or more linked message blocks. Each message block consists of a triplet with the following components:

• A data buffer

The data buffer contains the binary data that makes up the message. STREAMS imposes no alignment rules on the format of data in the data buffer, aside from those imposed by messages processed at the Stream head.

• A mblk\_t control structure

The mblk\_t structure contains information that the message owner can manipulate. Two of its fields are the read and write pointers into the data buffer.

A dblk\_t control structure

The dblk\_t structure contains information about buffer characteristics.

For example, two of its fields point to the limits of the data buffer, while others contain the message type.

The Stream head creates and fills in the message data structures when data is traveling downstream from an application. The Stream end creates and fills in the message data structures when data is traveling upstream, as in the case of data coming from an external communications device.

The mblk\_t and dblk\_t structures, shown in the following examples, are located in the <sys/stream.h> header file:

```
/* message block */
struct msgb {
          struct msgb * b_next; /* next message on queue */
          struct msgb * b_prev; /* previous message on queue */
struct msgb * b_cont; /* next message block of message */
unsigned char * b_rptr; /* first unread data byte in buffer */
          unsigned char * b_wptr; /* first unwritten data byte */
struct datab * b_datap; /* data block */
unsigned char b_band; /* message priority */
           unsigned char b_pad1;
           unsigned short b_flag; /* message flags */
                                 b_pad2;
           long
           MSG KERNEL FIELDS
};
typedef struct msgb
                                mblk_t;
/* data descriptor */
struct datab {
          union {
                      struct datab * freep;
                     struct free_rtn * frtnp;
           } db_f;
           unsigned char * db_base; /* first byte of buffer */
          unsigned char * db_lim; /* last byte+1 of buffer */
unsigned char db_ref; /* count of messages pointing */
very to block */
unsigned char db_type; /* message type */
          unsigned char db_type/ / message type /
unsigned char db_iswhat; /* message status */
unsigned int db_size; /* used internally */
caddr_t db_msgaddr; /* used internally */
                                db_filler;
           long
#define db_freep
                               db_f.freep
#define db_frtnp
                               db_f.frtnp
typedef struct datab
                                 dblk t;
/* Free return structure for esballoc */
typedef struct free_rtn {
          void (*free_func)(char *, char *); /* Routine to free buffer */
           char * free_arg;
                                                               /* Parameter to free_func */
} frtn_t;
```

When a message is on a STREAMS queue, it is part of a list of messages linked by b\_next and b\_prev pointers. The q\_next pointer points to the first message on the queue and the q\_last pointer points to the last message on the queue.

# 5.3.3 STREAMS Processing Routines for Drivers and Modules

A module or driver can perform processing on the Stream that an application requires. To perform the required processing, the STREAMS module or driver must provide special routines whose behavior is specified by the STREAMS framework. This section describes the STREAMS module and driver routines, and the following kinds of processing they provide:

- Open processing
- Close processing
- Configuration processing
- Read side put processing
- Write side put processing
- Read side service processing
- Write side service processing

#### Note

STREAMS modules and drivers must provide open, close, and configuration processing. The other kinds of processing described in this section are optional.

The format used to describe each routine in this section is XX\_routine\_name. Digital recommends that you substitute the name of a user-written STREAMS module or driver for the XX. For example, the open routine for the user-written STREAMS pseudodevice driver echo would be echo\_open.

# 5.3.3.1 open and close Processing

Only the open and close routines provide access to the u\_area of the kernel. They are allowed to sleep only if they catch signals.

# open processing

Modules and drivers must have open routines. The read side qinit structure, st\_rdinit defines the open routine in its  $qi\_qopen$  field. A driver's open routine is called when the application opens a Stream. The Stream head calls the open routine in a module when an application pushes the module onto the Stream.

The open routine has the following format:

The open routine can allocate data structures for internal use by the STREAMS driver or module. A pointer to the data structure is commonly stored in the q\_ptr field of the queue\_t structure. Other parts of the module or driver can access this pointer later.

# close processing

Modules and drivers must have close routines. The read side qinit structure, st\_rdinit, defines the close routine in its  $qi_qclose$  field. A driver calls the close routine when the application that opened the Stream closes it. The Stream head calls the close routine in a module when it pops the module from the stack.

The close routine has the following format:

```
XX_close(q, flag, credp)
    queue_t *q;    /* pointer to read queue */
    int flag;    /* file flag */
    cred_t *credp    /* pointer to credentials structure */
```

The close routine may want to free and clean up internally used data structures.

# 5.3.3.2 Configuration Processing

The configure routine is used to configure a STREAMS module or driver into the kernel. It is specific to Digital UNIX and its use is illustrated in Section 5.4.

The configure routine has the following format:

# 5.3.3.3 Read Side Put and Write Side Put Processing

There are both read side and write side XX\_Xput routines; XX\_wput for write side put processing and XX\_rput for read side put processing.

# Write Side Put Processing

The write side put routine, XX\_wput, is called when the upstream module's write side issues a putnext call. The XX\_wput routine is the only interface for messages to be passed from the upstream module to the current module or driver.

The XX\_wput routine has the following format:

```
XX_wput(q, mp)
    queue_t *q; /* pointer to write queue */
    mblk_t *mp; /* message pointer */
```

# **Read Side Put Processing**

The read side put routine, XX\_rput, is called when the downstream modules read side issues a putnext call. Because there is no downstream module, drivers that are Stream ends do not have read side put routines. The XX\_rput routine is the only interface for messages to be passed from the downstream module to the current module.

The XX\_rput routine has the following format:

```
XX_rput(q, mp)
    queue_t *q;    /* pointer to read queue */
    mblk_t *mp;    /* message pointer */
```

The XX\_Xput routines must do at least one of the following:

- Process the message
- Pass the message to the next queue (using putnext)
- Delay processing of the message by putting the message on the module's service routine (using putq)

The XX\_Xput routine should leave any large amounts of processing to the service routine.

# 5.3.3.4 Read Side Service and Write Side Service Processing

If an XX\_Xput routine receives a message that requires extensive processing, processing it immediately could cause flow control problems. Instead of processing the message immediately, the XX\_rput routine (using the putq call) places the message on its read side message queue and the XX\_wput places the message on its write queue. The STREAMS module notices that there are messages on these queues and schedules the module's

read or write side service routines to process them. If the module's XX\_rput routine never calls putq, then the module does not require a read side service routine. Likewise, if the module's XX\_wput routine never calls putq, then the module does not require a write side service routine.

The code for a basic service routine, either read side or write side, has the following format:

# 5.3.4 Digital UNIX STREAMS Concepts

The following STREAMS concepts are unique to Digital UNIX. This section describes these concepts and how they are implemented in Digital UNIX:

- Synchronization
- Timeout

# 5.3.4.1 Synchronization

Digital UNIX supports the use of more than one kernel STREAMS thread. Exclusive access to STREAMS queues and associated data structures is not guaranteed. Messages can move up and down the same Stream simultaneously, and more than one process can send messages down the same Stream.

To synchronize access to the data structures, each STREAMS module or driver chooses the synchronization level it can tolerate. The synchronization level determines the level of parallel activity allowed in the module or driver. Synchronization levels are defined in the sa.sa\_syn\_level field of the

streamadm data structure which is defined in the module's or driver's configuration routine. The sa.sa\_syn\_level field must have one of the following values:

# SQLVL\_QUEUE

Queue Level Synchronizaton. This allows one thread of execution to access any instance of the module or driver's write queue at the same time another thread of execution can access any instance of the module or driver's read queue. Queue level synchronization can be used when the read and write queues do not share common data. The SQLVL\_QUEUE argument provides the lowest level of synchronization available in the Digital UNIX STREAMS framework.

For example, the q\_ptr field of the read and write queues do not point to the same memory location.

# SQLVL\_QUEUEPAIR

Queue Pair Level Synchronizaion. Only one thread at a time can access the read and write queues for each instance of this module or driver. This synchronization level is common for most modules or drivers which process data and have only per-stream state.

For example, within an instance of a module, the q\_ptr field of the read and write queues points to the same memory location. There is no other shared data within the module.

# SQLVL MODULE

Module Level Synchronization. All code within this module or driver is single threaded. No more than one thread of execution can access all instances of the module or driver. For example, all instances of the module or driver are accessing data.

# SQLVL\_ELSEWHERE

Arbitrary Level Synchronization. The module or driver is synchronized with some other module or driver. This level is used to synchronize a group of modules or drivers that access each other's data. A character string is passed with this option in the sa.sync\_info field of the streamadm structure. The character string is used to associate with a set of modules or drivers. The string is decided by convention among the cooperating modules or drivers.

For example, a networking stack such as a TCP module and an IP module which share data might agree to pass the string tcp/ip. No more than one thread of execution can access all modules or drivers synchronized on this string.

# SQLVL\_GLOBAL

Global Level Synchronization. All modules or drivers under this level are single threaded. Note there may be modules or drivers using other levels not under the same protection. This option is available primarily for debugging.

# 5.3.4.2 Timeout

The Digital UNIX kernel interface to timeout and untimeout is as follows:

```
timeout(func, arg, ticks);
untimeout(func, arg);
```

However, to maintain source compatibility with AT&T System V Release 4 STREAMS, the <sys/stream.h> header file redefines timeout to be the System V interface, which is:

```
id = timeout(func, arg, ticks);
untimeout(id);
```

The *id* variable is defined to be an int.

STREAMS modules and drivers must use the System V interface.

# 5.4 Configuring a User-Written STREAMS-Based Module or Driver in the Digital UNIX Kernel

For your system to access any STREAMS drivers or modules that you have written, you must configure the drivers and modules into your system's kernel.

STREAMS modules or drivers are considered to be configurable kernel subsystems; therefore, follow the guidelines in the *Programmer's Guide* manual for configuring kernel subsystems.

The following sample procedure shows how to add to the kernel a STREAMS-based module (which can be a pushable module or a hardware or pseudodevice driver) called mymod, with it's source files mymodule1.c and mymodule2.c.

1. Declare a configuration routine in your module source file, in this example, /sys/streamsm/mymodulel.c.

Example 5-1 shows a module (mymod\_configure) that can be used by a module. To use the routine with a driver, do the following:

a. Remove the comment signs from the following line:

```
/* sa.sa_flags = STR_IS_DEVICE | STR_SYSV4_OPEN; */
This line follows the following comment line:
/* driver */
```

b. Comment out the following line:

```
sa.sa_flags = STR_IS_MODULE | STR_SYSV4_OPEN;
```

This line follows the following comment line:

```
/* module */
```

# **Example 5-1: Sample Module**

```
Sample mymodule.c
#include <sys/sysconfig.h>
#include <sys/errno.h>
struct streamtab mymodinfo = { &rinit, &winit };
cfq_subsys_attr_t mymod_attributes[] = {
      {"",0,0,0,0,0,0,0}
                      /* required last element */
};
mymod_configure(
 cfg_op_t op;
 caddr_t indata;
             indata_size;
 ulong
 caddr_t
            outdata;
 ulong
             outdata_size)
       dev_t devno = NODEV;
       struct streamadm sa;
       if (op != CFG_OP_CONFIGURE)
                return EINVAL;
                    = OSF_STREAMS_10;
       sa.sa_version
       /* sa.sa_flags = STR_IS_DEVICE | STR_SYSV4_OPEN; */
sa.sa_ttys = NULL;
       sa.sa_sync_level = SQLVL_MODULE;
                                           5
```

# Example 5-1: (continued)

```
sa.sa_sync_info = NULL;
strcpy(sa.sa_name, "mymod");

if ((devno = strmod_add(devno, &mymodinfo, &sa)) == NODEV)
{
         return ENODEV;
}

return ESUCCESS;
}
```

- 1 The subroutine in this example supplies an empty attribute table and no attributes are expected to be passed to the subroutine. If you want to develop attributes for your module, refer to the *Programmer's Guide* manual.
- 2 The first available slot in the cdevsw table is automatically allocated for your module. If you wish to reserve a specific device number, you should define it after examining the cdevsw table in the conf.c program. For more information on the cdevsw table and how to add device driver entries to it, see the *Writing Device Drivers: Tutorial*.
- **3** This example routine only supports the CFG\_OP\_CONFIGURE option. See the *Programmer's Guide* manual for information on other configuration routine options.
- 4 The STR\_SYSV4\_OPEN option specifies to call the module's or device's open and close routines, using the AT&T System V Release 4 calling sequence. If this bit is not specified, the AT&T System V Release 3.2 calling sequence is used.
- **5** Other options for the sa.sync\_level field are described in Section 5.3.4.
- 2. Statically link your module with the kernel.

If you want to make the STREAMS module dynamically loadable, see the *Programmer's Guide* for information on configuring kernel subsystems. If the module you are configuring is a hardware device driver, also see the *Writing Device Drivers: Tutorial*.

To statically link your module with the kernel, put your module's source files (mymodule1.c and mymodule2.c) into the /sys/streamsm directory and add an entry for each file to the /sys/conf/files file.

The following example shows the entries in the /sys/conf/files file for mymodule1.c and mymodule2.c:

```
streamsm/mymodule1.c optional mymod Notbinary streamsm/mymodule2.c optional mymod Notbinary
```

Add the MYMOD option to the kernel configuration file. The default kernel configuration file is /sys/conf/HOSTNAME (where HOSTNAME is the name of your system in uppercase letters.) For example, if your system is named DECOSF, add the following line to the /sys/conf/DECOSF configuration file:

```
options MYMOD
```

If you are configuring a hardware device driver continue with step 3; if not, got to step 4.

3. If you are configuring a hardware device driver, complete steps 3a to 3d.

If you are not configuring a hardware device driver, go to step 4.

If you are configuring a hardware device driver, you should already have an XXprobe and an interrupt routine defined. See the *Writing Device Drivers: Tutorial* for information about defining probe and interrupt routines.

a. Add the following line to the top of the device driver configuration file, which for this example is /sys/streams/mydriver.c:

```
#include <io/common/devdriver.h>
```

b. Define a pointer to a controller structure; for example:

```
struct controller *XXinfo;
```

For information on the controller structure, see the *Writing Device Drivers: Tutorial*.

c. Declare and initialize a driver structure; for example:

```
struct driver XXdriver =
{
    XXprobe, 0, 0, 0, 0, XXstd, 0, 0, "XX", XXinfo
};
```

For information on the driver structure, see the *Writing Device Drivers: Tutorial*.

d. Add the controller line to the kernel configuration file.

The default kernel configuration file is /sys/conf/HOSTNAME (where HOSTNAME is the name of your system in uppercase letters). For example, if your system name is DECOSF, would add a line

similar to the following to the /sys/conf/DECOSF configuration file:

controller XXO at bus vector XXintr

For information about the possible values for the bus keyword, see the *System Administration* manual.

- 4. Reconfigure, rebuild, and boot the new kernel for this system by using the doconfig command. See the doconfig(8) reference page or the *System Administration* manual for information on reconfiguring your kernel.
- 5. Run the strsetup -c command to verify that the device is configured properly:
  - # /usr/sbin/strsetup -c

STREAMS Configuration Information...Wed Jun 2 09:30:11 1994

Name	Type	Major	Minor	Module ID
clone		32	0	
ptm	device	37	0	7609
pts	device	6	0	7608
log	device	36	0	44
nuls	device	38	0	5001
echo	device	39	0	5000
sad	device	40	0	45
pipe	device	41	0	5304
kinfo	device	42	0	5020
xtisoUDP	device	43	0	5010
xtisoTCP	device	44	0	5010
dlb	device	49	0	5010
bufcall	module			0
timod	module			5006
tirdwr	module			0
ifnet	module			5501
ldtty	module			7701
null	module			5003
pass	module			5003
errm	module			5003
spass	module			5007
rspass	module			5008
pipemod	module			5303

Configured devices = 11, modules = 11

# 5.5 Device Special Files

This section describes the STREAMS device special files and how they are created. It also provides an overview of the clone device.

All STREAMS drivers must have a character special file created on the system. These files are usually in the /dev/streams directory and are created at installation, or by running the /usr/sbin/strsetup utility.

A STREAMS driver has a device major number associated with it which is determined when the driver is configured into the system. Drivers other than STREAMS drivers usually have a character special file defined for each major and minor number combination. The following is an example of an entry in the /dev directory:

```
    crw-----
    1 root
    system
    8, 1024 Aug 25 15:38 rrzla

    crw-----
    1 root
    system
    8, 1025 Aug 25 15:38 rrzlb

    crw-----
    1 root
    system
    8, 1026 Aug 25 15:38 rrzlc
```

In this example, rrzla has a major number of 8 and a minor number of 1024. The rrzlb device has a major number of 8 and a minor number of 1025, and rrzlc has a major number of 8 and a minor number 1026.

You can also define character special files for each major and minor number combination for STREAMS drivers. The following is an example of an entry in the /dev/streams directory:

```
crw-rw-rw- 1 root system 32, 0 Jul 13 12:00 /dev/streams/echo0 crw-rw-rw- 1 root system 32, 1 Jul 13 12:00 /dev/streams/echo1
```

In this example, echo0 has a major number of 32 and a minor number of 0, while echo1 has a major number of 32, and a minor number of 1.

For an application to open a unique Stream to a device, it must open a minor version of that device that is not already in use. The first application can do an open on /dev/streams/echo0 while the second application can do an open on /dev/streams/echo1. Since each of these devices has a different minor number, each application acquires a unique Stream to the echo driver. This method requires that each device (in this case, echo) have a character special file for each minor device that can be opened to it. This method also requires that the application determine which character special file it should open; it does not want to open one that is already in use.

The clone device offers an alternative to defining device special files for each minor device that can be opened. When the clone device is used, each driver needs only one character special file and, instead of an application having to determine which minor devices are currently available, clone allows a second (or third) device to be opened using its (clone device's) major number. The minor number is associated with the device being opened (in this case, echo). Each time a device is opened using clone device's major number, the STREAMS driver interprets it as a unique Stream.

The strsetup command sets up the entries in the /dev/streams directory to use the clone device. The following is an example entry in the /dev/streams file:

```
crw-rw-rw- 1 root system 32, 18 Jul 13 12:00 /dev/streams/echo
```

In this example, the system has assigned the major number 32 to the clone device. The number 18 is the major number associated with echo. When an application opens /dev/streams/echo, the clone device intercepts the call. Then, clone calls the open routine for the echo driver. Additionally, clone notifies the echo driver to do a clone open. When the echo driver realizes it is a clone open it will return its major number, 18, and the first available minor number.

#### Note

The character special files the /usr/sbin/strsetup command creates are created by default in the /dev/streams directory with clone as the major number. If you configure into your kernel a STREAMS driver that either does not use clone open, or uses a different name, you must modify the /etc/strsetup.conf file described in the strsetup.conf(4) reference page.

To determine the major number of the clone device on your system, run the strsetup -c command.

# 5.6 Error and Event Logging

STREAMS error and event logging involves the following:

- The error logger daemon
- The trace logger
- The strclean command

The error logger daemon, strerr, logs in a file any error messages sent to the STREAMS error logging and event tracing facility.

The trace logger, strace, writes to standard output trace messages sent to the STREAMS error logging and event tracing facility.

The strclean command can be run to clean up any old log files generated by the strerr daemon.

A STREAMS module or driver can send error messages and event tracing messages to the STREAMS error logging and event tracing facility through the strlog kernel interface. This involves a call to strlog.

The following example shows a STREAMS driver printing its major and minor device numbers to both the STREAMS error logger and the event tracing facility during its open routine:

A user process can also send a message to the STREAMS error logging and event tracing facility by opening a Stream to /dev/streams/log and calling putmsg. The user process must contain code similar to the following to submit a log message to strlog:

```
struct strbuf ctl, dat;
struct log_ctl lc;
char *message = "Last edited by <username> on <date>";

ctl_len = ctl.maxlen = sizeof (lc);
ctl.buf = (char *)&lc;

dat.len = dat.maxlen = strlen(message);
dat.buf = message;
lc.level = 0;
lc.flags = SL_ERROR|SL_NOTIFY;

putmsg (log, &ctl, &dat, 0);
```

# Extensible SNMP Application Programming Interface 6

The Simple Network Management Protocol (SNMP) is an application layer protocol that allows remote management and data collection from networked devices. A networked device can be anything that is connected to the network, such as a router, a bridge, or a host.

A managed networked device contains software that acts as the SNMP agent for the device. It handles the application layer protocol for SNMP and carries out the management commands. These commands consist of getting information and setting of operational parameters.

There are also network management application programs (usually running on a host somewhere on the network) that send SNMP commands to the various managed devices on the network to perform the management tasks. These tasks can consist of configuration management, network traffic monitoring and network trouble shooting.

The Extensible Simple Network Management Protocol (eSNMP) is the SNMP agent architecture for a host machine on the network running Digital UNIX Version 4.0 (or higher). It includes a master-agent process and multiple related processes containing eSNMP subagents. The master-agent performs the SNMP protocol handling and the subagents perform the requested management commands. This section assumes you are familiar with the following:

- SNMP protocol
- Management Information Base (MIB) definitions and Request For Comments (RFCs)
- Object Identifiers (OIDs) and the International Standards Organization (ISO) registration hierarchy (1.3.6.1.2.1, and so on)
- The C programming language

This chapter provides the following information:

- Overview of eSNMP
- Overview of the eSNMP application programming interface (API)
- Detailed information on the eSNMP routines

# 6.1 Overview of eSNMP

This section describes the components and architecture the eSNMP agent for Digital UNIX. It contains information on the following:

- Components of eSNMP
- Architecture
- SNMP Versions

# 6.1.1 Components of eSNMP

The eSNMP components are as follows:

- /usr/sbin/snmpd The master-agent daemon.
- /usr/sbin/os\_mibs The subagent daemon provided by Digital UNIX.
- /usr/sbin/mosy The MIB compiler.
- /usr/sbin/snmpi The object table code generator.
- /usr/shlib/libesnmp.so The eSNMP Library.
- /usr/include/esnmp.h-eSNMP definitions.
- /usr/examples/esnmp/\* Example code.

The Management Information Base (MIB) defines a set of data elements that relate to network management. Many of these are standardized in the RFCs which are produced as a result of the Internet Engineering Task Force (IETF) working group standardization effort of the Internet Society.

The data elements defined in the RFCs are identified using a naming scheme with a hierarchical structure. Each name at each level of the hierarchy has a number associated with it. You can refer to the data elements in the MIB definitions by name or by its corresponding sequence of numbers. This is called the Object Identifier (OID). You can extend an OID for an specific data element further by adding more numbers to identify a specific instance of the data element. The entire collection of managed data elements is called the MIB tree.

Each SNMP agent implements those MIB elements that pertain to the device being managed, plus a few common MIB elements. These are the supported MIB tree elements. An extensible SNMP agent is one that permits its supported MIB tree to be distributed among various processes and change dynamically.

For eSNMP there is a single master-agent and there may be any number of subagents. The master-agent itself does not support (implement) any MIBs, it handles the SNMP protocol and maintains a registry of subagents and the

MIBs they support. The master-agent for eSNMP is the daemon process /usr/sbin/snmpd.

The eSNMP protocol contains one standard subagent that implements the common MIB elements contained under the mib-2 OID name. This is the daemon process /usr/sbin/os\_mibs. Another eSNMP subagent is built into the gated daemon process (/usr/sbin/gated). Additional subagents will be added by Digital and third parties. These subagents communicate with the master-agent and work together to appear to the management application programs as a single SNMP agent for the host.

# 6.1.2 Architecture

The master-agent listens on the preassigned User Datagram Protocol (UDP) port for an incoming SNMP request. When the master-agent receives an SNMP request, it authenticates it against the local security database and handles any authentication or protocol errors. If the request is valid, the snmpd daemon consults its MIB registry. (See the snmpd(8) reference page for more information.) For each MIB object contained in the request it determines which registered MIB could contain that object and which subagent has registered that MIB. The master-agent then builds a series of messages; one for each subagent that will be involved in this SNMP request. These messages do not carry SNMP, but use the more efficient eSNMP protocol<sup>1</sup> for communication between the master-agent and the subagents.

Each subagent program is linked with the shareable library libesnmp.so. This library contains the protocol implementation that enables communication between the master-agent and the subagent. This code parses the master-agent's message and consults its local object table.

The object table is a data structure that is defined and initialized in code emitted by the MIB compiler tools. It contains an entry for each MIB object that is contained in the MIBs implemented in that subagent. One part of an object table entry is the address of a function which services requests for this MIB object. These functions are called method routines.

The eSNMP library code calls into the indicated method routine for each of the MIB variables in the master-agent's message. The eSNMP library code creates a response packet based on the function return values and sends it back to the master-agent.

The master-agent starts a timer and marshals the response packets from all involved subagents. The master-agent may rebuild and resend a new set of

On Digital UNIX this protocol is based on the Distributed Protocol Interface (DPI) V2, RFC 1592. DPI V2 is an experimental protocol and Digital UNIX does not attempt to adhere to the protocol completely. DPI was accounted for in the design so that Digital could more easily implement an eventual standard protocol. Digital UNIX will track and support standards that emerge in the area of extensible SNMP.

subagent messages, depending on the specific request; for example, a GetNext request. When the master-agent has all required data or error responses or has timed out waiting for a response from a subagent, it builds an SNMP response message and sends it to the originating SNMP application. The interaction between the master-agent and subagent is invisible to the requesting SNMP management application.

Subagent programs are linked against libesnmp.so shareable library, which performs all the protocol handling and dispatching. Subagent developers need to code the method routines for their MIB objects.

# 6.1.3 SNMP Versions:

The IETF working group is readdressing SNMPv2 and RFCs have not been published, at the time of this writing.

Extensible SNMP support for SNMPv2 does exist in the following areas. This is based on the original SNMPv2 RFCs that were submitted and withdrawn:

- The MIB tools (the mosy and snmpi programs) support SNMPv2 SMI and textual conventions.
- The eSNMP library API supports SNMPv2, variable binding exceptions, and error codes. It ignores MIB objects with SNMPv2-only data types when processing a SNMPv1 request and does not call associated method routines.
- The master-agent currently supports SNMPv1 only and maps all SNMPv2-specific information from the subagent into SNMPv1-adherent data. In a future release the master-agent will support both SNMPv1 and SNMPv2 in a bilingual manner. This will not require code changes within subagents. Therefore, documented SNMPv2 features (such as GetBulk) are latent.

# 6.2 Overview of the Extensible SNMP Application Programming Interface

The subagent's function is to establish communications with the masteragent, register the MIBs that it is going to handle, and process requests from the master-agent. It must also be able to send SNMP traps on behalf of the host application.

The subagent consist of the following:

- A main function written by the developer
- The eSNMP Library routines which perform the eSNMP protocol work

- The method routines written by the developer that handle specific MIB elements
- The object table structures generated from MIB definition files using the mosy and snmpi programs

The subagent is usually embedded within a host application, such as a router daemon. Here the subagent processing is only a small part of the work performed by the process. The main routine of the host application contains the calls to the eSNMP library to perform the eSNMP protocol. In other cases, the subagent is a standalone daemon process that has its own main routine.

The eSNMP library calls the method routines while processing a packet from the master-agent. Each MIB variable in the object table has a pointer to the method routine that is to handle that variable. Since the object tables are generated by the mosy and snmpi programs, the method routine names are static.

The eSNMP developer's kit provided with Digital UNIX consists of the following:

- /usr/sbin/mosy MIB compiler utility
- /usr/sbin/snmpi Object table code generator utility
- /usr/examples/esnmp/mib-converter.sh MIB text extraction tool
- /usr/shlib/libesnmp.so eSNMP library
- /usr/include/esnmp.h eSNMP definitions file
- /usr/examples/esnmp/\* Subagent example code

The eSNMP library (libesnmp.so) contains the following:

- The master-agent to subagent protocol handling routines
   These routines implement communication with the master-agent on behalf of the subagent; they are:
  - esnmp\_init Initializes the protocol (performs a handshake with the master-agent)
  - esnmp\_register Registers a MIB with the master-agent
  - esnmp poll Processes a packet from the master-agent
  - esnmp\_trap Requests the master-agent to generate an SNMP trap
  - esnmp\_are\_you\_there Pings the master-agent
  - esnmp\_unregister Unregisters a MIB
  - esnmp\_term Ends communication with the master-agent and

terminate extensible SNMP

esnmp\_sysuptime - Time handling and synchronization

# Support routines

These are also resolved in libesnmp.so, and are optional routines for convenience in developing method routines. These include, but are not limited to, the following:

- str2oid Converts an ASCII dot-format string into internal OID format; see Section 6.3.3.5 for more information.
- cmp\_oid Compares the value of two OID structures; see Section
   6.3.3.10 for a complete list.

The esnmp.h header file is associated with the eSNMP library. This file defines all data structures, constants, and function prototyes required to implement subagents to this API.

# 6.2.1 Subtrees

Understanding subtrees is crucial to understanding the eSNMP API and how your subagent will work.

## Note

This section assumes that you understand the OID naming structure used in SNMP. If not, refer to RFC1442 *Structure of Management Information*.

The information in SNMP is structured hierarchically like an inverted tree. Data can be associated with any leaf node in this hierarchy. Each node has a name and a number. Each node can also be identified by an OID, which is an accumulation of the numbers that make up a path from the root down to that node in the tree.

For example, the chess MIB used in the sample code has an element with the name chess. The OID for the element chess is

1.3.6.1.4.1.36.2.15.2.99, which is derived from its position in the hierarchy: (The chess MIB appears in the /usr/examples/esnmp directory.)

```
iso(1)
  org(3)
  dod(6)
   internet(1)
    private(4)
    enterprise(1)
     digital(36)
    ema(2)
```

```
sysobjects(15)
decosf(2)
chess(99)
```

Any node in the MIB hierarchy can define a subtree. All elements within the subtree have an OID that starts with OID of the subtree base. For example, if we define chess to be a subtree base, the elements with the same prefix as the chess OID are all within the subtree:

```
      chess
      1.3.6.1.4.1.36.2.15.2.99

      chessProductID
      1.3.6.1.4.1.36.2.15.2.99.1
      ObjectID

      chessMaxGames
      1.3.6.1.4.1.36.2.15.2.99.2
      Integer32

      chessNumGames
      1.3.6.1.4.1.36.2.15.2.99.3
      Integer32

      gameTable
      1.3.6.1.4.1.36.2.15.2.99.4
      Integer32

      gameEntry
      1.3.6.1.4.1.36.2.15.2.99.4.1
      Integer32

      gameIndex
      1.3.6.1.4.1.36.2.15.2.99.4.1.1
      Integer32

      gameNumMoves
      1.3.6.1.4.1.36.2.15.2.99.4.1.2
      DisplayString

      gameStatus
      1.3.6.1.4.1.36.2.15.2.99.4.1.4
      INTEGER

      moveTable
      1.3.6.1.4.1.36.2.15.2.99.5
      Integer32

      moveEntry
      1.3.6.1.4.1.36.2.15.2.99.5
      Integer32

      moveIndex
      1.3.6.1.4.1.36.2.15.2.99.5
      Integer32

      moveByWhite
      1.3.6.1.4.1.36.2.15.2.99.5
      Integer32

      moveByBlack
      1.3.6.1.4.1.36.2.15.2.99.5
      Integer32

      chessTraps
      1.3.6.1.4.1.36.2.15.2.99.5
      Integer32

      moveTrap
      1.3.6.1.4.1.36.2.15.2.99.6
      Integer32
```

It is this subtree base that is registered with the master-agent to tell it that this subagent handles all requests related to the elements within the subtree.

The master-agent expects a subagent to handle all objects subordinate to the registered subtree. This principle guides your choice of subtrees.

For example, registering a subtree of chess is reasonable because it is realistic to assume that the subagent could handle all requests for elements in this subtree. Registering an entire application-specific MIB usually makes sense because the particular application expects to handle all objects defined in the MIB.

Registering a subtree of transmission (under MIB-2) would be a mistake, because it is unlikely that the subagent is prepared to handle every MIB subordinate to transmission (FDDI, Token Ring, and so on).

A subagent may register as many subtrees as it wants. It can register OIDs that overlap with other registrations by itself or other subagents; however, it cannot register the same OID more than once. The subagents can register and unregister subtrees at any time after it has established communication with the master-agent.

Normally it is the nonterminal nodes that are registered as a subtree with the master-agent. However, terminal nodes (those of one object type), or even specific instances, can be registered as a subtree.

The master-agent distributes requests to the subagent that has the subtree with the highest priority (largest priority number) or the most recent (if priority is equal), matching the OID on the variable bindings of the request.

# 6.2.2 Object Tables

The mosy and snmpi utilities are used to generate the C language code that defines the object tables from the MIBs. The object tables are defined in the emitted files <code>subtree\_tbl.h</code> and <code>subtree\_tbl.c</code>, files that are compiled into your subagent.

These modules are created by the utilities and it is not recommended that they be edited. If the MIBs change or a future version of the eSNMP development utilities require your object tables to be rebuilt, it is easy to rebuild the files and recompile them if you did not edit the files.

#### 6.2.2.1 The subtree tbl.h File

The *subtree\_tbl*.h file contains the following information:

- A declaration of the subtree structure
- Index definitions for each MIB variable in the subtree
- Enumeration definitions for MIB variables with enumerated values
- MIB group data structure definitions
- Method routine function prototypes

The first section is a declaration of the subtree structure. The subtree is automatically initialized by code in the <code>subtree\_tbl.c</code> file. A pointer to this structure is passed to the <code>esnmp\_register</code> routine to register a subtree with the master-agent. All access to the object table for this subtree is through this pointer. The declaration has the following form:

# extern SUBTREE subtree subtree;

The next section contains index definitions for each MIB variable in the SUBTREE of the form:

#define I mib-variable nnnn

These values are unique for each MIB variable within a subtree and are the index into the object table for this MIB variable. These values are also generally used to differentiate between variables that are implemented in the same method routine so they can be used in a switch operation.

The next section contains enumeration definitions for those integer MIB variables that are defined with enumerated values, as follows:

### #define D\_mib-variable\_enumeration-name value

These are useful since they describe the architected value that enumerated integer MIB variables may take on; for example:

```
/* enumerations for gameEntry group */
  #define D_gameStatus_complete 1
  #define D_gameStatus_underway 2
  #define D_gameStatus_delete 3
```

The next section contains the MIB group data structure definitions of the form:

```
typedef struct xxx {
      type mib-variable;
    .
    char mib-variable_mark;
    .
    .
    .
} mib-group_type
```

One of these data structures is emitted for each MIB group within the subtree. Each structure definition contains a field representing each MIB variable within the group. If the MIB variable name is not unique within the pool of MIBs presented to the snmpi program at the time the <code>subtree\_tbl.h</code> file is built, the snmpi program does not qualify the name with the name of its parent variable (group name) to make it unique. In addition to the MIB variable fields, the structure includes a 1-byte <code>mib-variable\_mark</code> field for each variable. You can use these for maintaining status of a MIB variable; for example, the following is the group structure for the chess MIB:

```
typedef struct _chess_type {
   OID    chessProductID;
   int    chessMaxGames;
   int    chessNumGames;
   char chessProductID_mark;
   char chessMaxGames_mark;
   char chessNumGames_mark;
}
```

These MIB group structures are provided for convenience, but are not mandatory. You can use whatever structure is easiest for you in your method routine.

The next section is the method routine function prototypes. Each MIB group

within the subtree has a method routine prototype defined. A MIB group is a collection of MIB variables that are leaf nodes and share a common parent node.

There is always a function prototype for the method routine that handles the Get, GetNext, and GetBulk operations. If the group contains any writable variables, there is also a function prototype for the method routine that handles Set operations. Pointers to these routines appear in the subtree's object table which is initialized in the <code>subtree\_tbl.c</code> module. You must write method routines for each prototype that is defined, as follows:

```
extern int mib-group_get(METHOD *method)
extern int mib-group_set(METHOD *method)
```

# For example:

```
extern int chess_get(METHOD *method);
extern int chess_set(METHOD *method);
```

Method routines are discussed in more detail in Section 6.3.2.3.

# 6.2.2.2 The subtree tbl.c File

The *subtree* tbl.c file contains the following information:

- An array of integers representing the OIDs for each MIB variable
- An array of OBJECT structures. (See esnmp.h.)
- The initialized SUBTREE structure

The first section is the array of integers used for the OIDs of each MIB variable in the subtree, as follows:

```
static unsigned int elems[] = { ...
```

The next section is an array of OBJECT structures. There is one OBJECT for each MIB variable within the subtree. (See esnmp.h.)

An OBJECT represents a MIB variable and has the following fields:

- object\_index The constant I\_mib-variable from the subtree\_tbl.hfile.
- oid The this is the variable's OID (points to a part of elems[]).
- type The variable's data type.
- getfunc The address of method routine to call for Get operations.
- setfunc The address of method routine to call for Set operations.

The master-agent has no knowledge of object tables or MIB variables. It only maintains a registry of subtrees. When a request for a particular MIB variable arrives, it is processed as follows. In the following procedure, the MIB variable is mib var and the subtree is subtree 1:

- 1. The master-agent finds which subagent registered subtree\_1 which contains (for Get or Set requests) or might contain (for GetNext or GetBulk requests) mib var.
- 2. It sends an eSNMP message to the subagent that registered subtree\_1.
- 3. The subagent consults its list of registered subtrees and locates subtree\_1. It searches the object table of subtree\_1 and locates the following:
  - mib\_var (for Get and Set requests)
  - The first object lexicographically after mib\_var (for Next or Bulk requests)
- 4. It calls the appropriate method routine. If the method routine completes successfully, the data is returned to the master-agent. If not, for Get or Set, an error is returned. For Next or Bulk, the libsnmp code keeps trying subsequent objects in the object table of subtree\_1 until a method routine returns success or the table is exhausted; in either case a response is returned.
- 5. If the master-agent detects subtree\_1 could not return data on a Next or Bulk routine, it recursively tries the subtree lexicographically after subtree 1.

The next section is the SUBTREE structure itself. It is a pointer to this structure that is passed to the <code>esnmp\_register</code> eSNMP library routine to register the subtree. It is through this pointer that the library routines find the object structures. The following is an example of the <code>chess</code> subtree structure:

The SUBTREE structure has the following elements:

- name This is the name of the base node of the subtree.
- dots The ASCII string representation of the subtree's OID; it is what actually gets registered.
- oid The OID of the base node of the subtree; it points back to the array of integers.
- *object\_tbl* A pointer to the array of objects in the object table. It is indexed by the I\_xxxx definitions found in the *subtree\_tbl*.hfile.

• *last* – This is the index of the last object in the object\_tbl file. It is used to determine when the end of the table has been reached.

The final section of the *subtree\_tbl.c* contains short routines for allocating and freeing the *mib-group\_type* structures. These are provided as a convenience and are not a required part of the API.

# 6.2.3 Implementing a Subagent.

As a subagent developer, you are usually presented with a UNIX application, daemon, or driver (such as the gated daemon or ATM drivers) and have to implement an SNMP interface. The following steps explain how you do this:

1. Obtain a MIB specification.

MIB development starts with a MIB specification. Usually these are RFCs, written in concise MIB format according to RFC 1212. Designing and specifying a MIB is beyond the scope of this document; it is assumed you have a MIB specification.

The standard RFCs can be obtained from the the InterNIC directory at the following URL:

http://ds.internic.net/ds/dspglintdoc.html

If you have to build your own MIB specification, you can look at a similar MIBs written by another vendor. One source for a listing of these is in the archives section of the Network Management page at the following URL:

http://smurfland.cit.buffalo.edu/NetMan/index.html

You need MIBs for all of the elements you are implementing in the subagent and for any elements referenced by these MIBs (such that all element names resolve to the OID numbers). As a minimum you will need the SMI MIB rfc1442.my and the textual conventions v2-tc.my. These are in the /usr/examples/esnmp directory.

2. Compile your MIBs.

Once you obtain MIB definitions, use them to generate the object tables for your new subagent. The objective is to take the MIB specification text for each of the MIBs, remove the ASN.1 specifications, and compile them into C language modules that contain the local object tables.

Compile your MIBs using the following tools:

mib-converter.sh

The mib-converter.sh is a gawk shell script that extracts the MIB ASN.1 definitions from the RFC text. This step removes the text before and after the MIB definition and removes page headings

and footings.

The mib-converter.sh script may not remove everything that needs to be removed; therefore, you may need to remove some things manually, using a text editor. The following is an example of how to use the mib-converter.sh script:

# /usr/examples/esnmp/mib-converter.sh mib-def.txt > \
mib-def.my

Be careful; some RFCs contain more than one MIB definition. You can only use the mib-converter.sh script shell on RFCs that contain a single MIB definition. The mosy compiler may not handle it either. If you use an RFC that contains more than one MIB definition, make each one into a separate file. The resulting files containing the MIBs should be in the following form:

mib-def.my

- mosy

The Managed Object Syntax (mosy) compiler parses .my files created by the mib-converter.sh script and compiles them into .defs files. The .defs files describe the object hierarchy within the MIB. The .defs files are front-ends to several tools. The following is an example of how to use the mosy compiler:

# mosy mib-def.my

The mosy compiler produces mib-def.defs files.

The mosy program is taken from ISODE 8.0 (distributed with the 4BSD/ISODE SNMPv2 package).

– snmpi

The MIB data initializer creation program (snmpi) reads a concatenation of the .def files compiled by the mosy compiler and generates the C code to define the static structures of the object table for a specified MIB subtree.

#### Note

The snmpi program supplied with Digital UNIX is different from the snmpi program in 4BSD/ISODE SMUX.

Concatenate the .def files the mosy compiler compiles into the objects.defs file. Be sure to include the compiled versions of rfc1442.my and v2-tc.my. The objects.defs file must

contain enough MIBs to resolve all MIB names, even if they are not used by your subtrees. Then generate the object table files using the following command:

### # /usr/sbin/snmpi objects.defs subtree

The snmpi program has a print option that allows you to dump the contents of the entire tree generated as a result of the objects it finds into the objects.defs file. If you are having trouble with the subtrees you may find this to be helpful. Use the following command to generate a listing:

### # /usr/sbin/snmpi -p objects.defs > objects.txt

The snmpi program outputs the <code>subtree\_tbl.c</code> and <code>subtree\_tbl.h</code>. The <code>subtree</code> is the name of the base MIB variable name for a MIB subtree. These two files are C code used to initialize the MIB object table for the specified subtree. (This is the local object table referred to above.) Repeat this process for each MIB subtree being implemented in your subagent. Note that the snmpi program defaults to using MIB groups as the level of granularity for method routines; that is, the assumption is made that all MIB variables within a group should be serviced by the same method routine. (It also provides the <code>mib-group\_type</code> data structure to help do this.)

The <code>mib-group\_type</code> structure is not part of the API; it is provided as a convenience. It is helpful to use the <code>mib-group</code> organization of the object table. This is because, generally, those objects are logically related and usually accessed as a group; for example, <code>ipRoutes</code> are returned more or less complete from the kernel routing tables.

### 3. Code the method routines and the API calls.

Write the code that calls the eSNMP library API to initialize communications with the master-agent (snmpd), and register your MIBs. (See Section 6.2.4.)

Write the code for the required method routines. (See Section 6.3.) Usually you need one Get method routine and one Set method routine for each MIB group within your registered MIB subtree. The <code>subtree\_tbl.h</code> files generated in the previous step define the names and function prototype for each method routine you need.

### 4. Build the subagent.

An example Makefile is provided in the /usr/examples/esnmp directory.

5. Execute and test your subagent.

Run your subagent like any other program or daemon. There are trace facilities built into the eSNMP library routines to assist in the debugging process. Use the set\_debug\_level routine in the main section to enable the trace.

Once the subagent has initialized and successfully registered a MIB subtree, you can send SNMP requests using standard applications. For example, POLYCENTER Netview, HP OPenview, or any MIB browser. If you do not have access to SNMP applications, you can use the snmp\_request and snmp\_traprcv programs to help debug subagents.

Note that if you interactively debug, your subagent will probably cause SNMP requests to timeout.

Normally all error and warning messages are recorded in the system's daemon log. When running the sample chess subagent and the os\_mibs subagent, you specify a trace runtime argument, as follows:

```
os_mibs -trace
```

With the trace option active, the program does not daemonize and all trace output goes to stdout; it displays each message that is processed.

You can use this feature in your own subagents by calling the set\_debug\_level routine and pass it the TRACE parameter.

Anything passed in the debug macro is sent to stdout, as follows:

```
ESNMP LOG ((TRACE, ("message text \n"));
```

To send everything to the daemon log, call the set\_debug\_level routine and pass it the WARNING | DAEMON\_LOG parameter or the set\_debug\_level routine and pass it the ERROR | DAEMON\_LOG parameter to suppress warning messages.

# 6.2.4 Subagent Protocol Operations

The eSNMP API provides for autonomous subagents that are not closely tied to the master agent (snmpd). Subagents can be part of other subsystems or products and have primary functions not related to SNMP. For instance, the gated daemon is primarily concerned with Internet routing; however it also functions as a subagent.

In particular, the snmpd daemon does not start or stop any subagent daemons during its startup or shutdown procedures. It also does not maintain any on-disk configuration information about subagents. Whenever the snmpd daemon starts, it has no knowledge of previously registered subagents or subtrees.

Typically all daemons on a Digital UNIX system are started or stopped together, as the system changes run levels. But subagents should correctly handle situations where they start before the snmpd daemon, or are running while the snmpd daemon is restarted to reload information from its configuration file. In these situations subagents need to restart the eSNMP protocol as described in the following sections.

## 6.2.4.1 Order of Operations

Subagent protocol operations follow the following sequence:

- Initialization (esnmp\_init)
- 2. Registration (esnmp\_register [esnmp\_register ...])
- 3. Data communication

```
The following loop happens continuously:

{
    determine sockets with data pending
    if the eSNMP socket has data pending
        esnmp_poll

    periodically call esnmp_are_you_there as required during periods of inactivity
}
```

4. Termination (esnmp\_term)

Note that is very important that subagents call the <code>esnmp\_term</code> function when they are stopping. This enables eSNMP to free system resources being used by the subagent.

The example subagent in the /usr/examples/esnmp directory shows how to code subagent protocol operations.

# 6.2.4.2 Function Return Values

The eSNMP API function return values indicate to a subagent both the success or failure of the requested operation and the state of the master agent. The following list provides a description of each return value and the indicated subagent actions:

• ESNMP LIB OK

The operation was successful.

• ESNMP LIB NO CONNECTION

The connection between the subagent and the master agent could not be initiated. This value is returned by the esnmp init function.

- Causes The master agent is not running or is not responding.
- Action Restart the protocol by calling the esnmp\_init function again after a suitable delay.

### • ESNMP LIB DUPLICATE

A duplicate subagent identifier has been received by the master agent. This means that another process with the same subagent identifier is connected to the master agent and that this process should terminate. This value is returned by the esnmp\_poll function.

- Causes Typically this means the subagent daemon was started more than once; but it may indicate a different subagent used the same identifier.
- Action This invocation of the subagent process will never be able successfully initialize eSNMP, so the subagent should terminate.

## ESNMP\_LIB\_LOST\_CONNECTION

Lost communications with the master-agent. This value is returned by the esnmp\_register, esnmp\_poll, esnmp\_are\_you\_there, esnmp unregister, and esnmp trap functions.

- Causes An attempt to send a packet to the master agent's socket failed; this is normally due to the master agent terminating abnormally.
- Action Restart the protocol by calling the esnmp\_init function after a suitable delay.

## • ESNMP\_LIB\_BAD\_REG

The attempt to send a registration failed. This value is returned by the esnmp\_register, esnmp\_unregister, and esnmp\_poll. functions.

- Causes are as follows:
  - The esnmp\_init function has not been successfully called prior to calling the esnmp\_register function.
  - The timeout parameter in the esnmp\_register function is invalid.
  - The subtree passed to the esnmp\_register function has already been queued for registration or has been registered by this subagent.
  - A previous registration was failed by the master-agent (when returned by the esnmp\_poll function). See the log file to determine the details regarding why it failed and which subtree was at fault.

- Trying to unregister a subtree that was not registered (esnmp unregister).
- Action Call the esnmp\_register function in the proper sequence and with correct arguments.
- ESNMP LIB CLOSE

The master-agent is stopping. This value is returned by the esnmp\_poll function.

- Causes The master agent is beginning an orderly shutdown.
- Action Restart the protocol with the esnmp\_init function as suited by the subagent.
- ESNMP LIB NOTOK

An eSNMP protocol error occurred and the packet was discarded. This value is returned by the esnmp\_poll, and esnmp\_trap functions.

- Causes This indicates a packet-level protocol error within eSNMP, probably due to lack of memory resources within the subagent.
- Action Continue.

# 6.3 Extensible SNMP Application Programming Interface

This section provides detailed information on the SNMP Application Programming Interface, which consists of the following:

- Calling interface
- Method routine calling interface
- The libesnmp support routines

# 6.3.1 Calling Interface

The calling interface contains the following routines:

- esnmp\_init
- esnmp\_register
- esnmp\_unregister
- esnmp\_poll
- esnmp\_are\_you\_there
- esnmp\_trap
- esnmp\_term

• esnmp\_sysuptime

# 6.3.1.1 The esnmp init Routine

The esnmp\_init routine locally initializes the extensible SNMP subagent, and initiates communication with the master-agent.

This call does not block waiting for a response from the master-agent. After calling the esnmp\_init routine, call the esnmp\_register routine for each subtree that is to be handled by this subagent.

Call this routine during program initialization or to restart the eSNMP protocol. If you are restarting, the esnmp\_init routine clears all registrations so each subtree must be reregistered.

You should attempt to create a unique <code>subagent\_identifier</code>, perhaps using the program name (<code>argv[0]</code>) and additional descriptive text. The master-agent does not open communications with a subagent whose subagent-identifier is already in use.

The syntax for the esnmp init routine is as follows:

int esnmp init (int \*socket, char \*subagent identifier)

The arguments are defined as follows:

socket

The address of the integer that receives the socket descriptor used by eSNMP.

```
subagent_identifier
```

The address of a null-terminated string that uniquely identifies this subagent (usually program name).

The return values are as follows:

# Status ESNMP\_LIB\_NO\_CONNECTION

Could not initialize or communicate with the master-agent. Try again after a delay.

```
ESNMP_LIB_OK
```

Indicates the esnmp\_init routine has completed successfully.

The following is an example of the esnmp\_init routine:

```
#include <esnmp.h>
int socket;
status = esnmp init(&socket, "gated");
```

# 6.3.1.2 The esnmp\_register Routine

The esnmp\_register routine requests registration of a single MIB subtree. Before the master-agent can pass SNMP requests on to the subagent, it must register the willingness to process all messages for MIB variables subordinate to a subtree identifier.

The initialization routine (esnmp\_init) must be called prior to calling the esnmp\_register routine. The esnmp\_register function must be called for each subtree structure corresponding to each subtree that it will be handling. At any time subtrees can be unregistered by calling esnmp\_unregister and then be reregistered by calling the esnmp\_register.

When restarting the eSNMP protocol by calling esnmp\_init, all registrations are cleared. All subtrees must be reregistered.

A subtree is identified by the base MIB name and its corresponding OID number of the node which is the parent of all MIB variables that are contained in the subtree; for example, the MIB-2 tcp subtree has an OID of 1.3.6.1.2.1.6. All elements subordinate to this (those that have the same first 7 digits) are included in the subtree's object table. The subtree can also be a single MIB object (a leaf node) or even a specific instance.

By registering a subtree, the subagent is indicating that it will process SNMP requests for all MIB variables (or OIDs) within that subtree's range. Therefore, a subagent should register the most fully qualified (longest) subtree that still contains its instrumented MIB variables.

For example, the Digital UNIX operating system contains support for MIB-2 implemented as an eSNMP subagent. This subagent does not register MIB-2 (1.3.6.1.2.1); instead, it registers the following MIBs: at, dot5, egp, fddi, icmp, interfaces, IP, snmp, system, tcp, and udp.

The master-agent requires that a subagent cannot register the same subtree more than once. Other than this one restriction, a subagent may register subtrees that overlap the OID range of subtrees that it previously registered or those of subtrees registered by other subagents.

For example, consider the two Digital UNIX daemons, os\_mibs and gated. The os\_mibs daemon registers the ip subtree and the gated daemon registers the ipRouteTable subtree at a higher priority. Requests for operations on MIB objects within ipRouteEntry, such as ipRouteIfIndex, will go to gated because it is a higher priority. Requests for other ip objects, such as ipNetToMediaIfIndex, will be passed to os\_mibs. If the gated process should terminate or unregister the ipRouteEntry subtree, subsequent requests for ipRouteIfIndex will go to os\_mibs because the ip subtree, which includes the ipRouteEntry objects, will now be the highest priority in that range.

When the master-agent receives a SIGUSR1 signal, it puts its MIB registry in to the /var/tmp/snmpd\_dump.log file. See the snmpd(8) reference page for more information.

The syntax for the esnmp\_register routine is as follows:

int esnmp\_register(SUBTREE \*subtree,int timeout,int priority)

The arguments are defined as follows:

#### subtree

A pointer to a SUBTREE structure corresponding to the subtree to be handled. The SUBTREE structures are externally declared and initialized in the code emitted by the mosy and snmpi utilities (xxx\_tbl.c and xxx\_tbl.h, where xxx is the name of the subtree) taken directly from the MIB document.

#### timeout

The number of seconds the master-agent should wait for responses when requesting data in this subtree. This value must be between zero (0) and ten (10). If the value is zero (0), the default timeout is used (3 seconds). Digital recommends you use the default.

#### priority

This is the registration priority. The entry with largest number has the highest priority. The range is 0 to 65535. The subagent that has registered a subtree that has the highest priority over a range of Object Identifiers (OIDs) gets all requests for that range of OIDs.

Subtrees that are registered with the same priority are ranked in order by time of registration. The most recent registration has the highest priority.

The *priority* argument is a mechanism for cooperating subagents to handle different configurations.

The return values are as follows:

### ESNMP LIB OK

Indicates the esnmp register routine has completed successfully.

# ESNMP LIB BAD REG

Indicates the esnmp\_init routine has not been called, the timeout parameter is invalid, or this subtree has already been queued for registration.

#### ESNMP LIB LOST CONNECTION

Indicates the subagent has lost communications with the master-agent.

Note that the status indicates only the initiation of the request. The actual status returned in the master-agent's response will be returned in a

subsequent call to the esnmp\_poll routine.

The following is an example of the esnmp\_register routine:

# 6.3.1.3 The esnmp\_unregister Routine

The esnmp\_unregister routine unregisters a MIB subtree with the master-agent.

This routine can be called by the application code to tell the eSNMP subagent not to process requests for variables in this subtree anymore. You can later reregister a subtree, if needed, by calling the esnmp\_register routine

The syntax for the esnmp\_unregister routine is as follows:

int esnmp\_unregister(SUBTREE \*subtree)

The arguments are as follows:

\*subtree

A pointer to the subtree structure for the subtree to be unregistered.

The return values are as follows:

ESNMP\_LIB\_OK

Indicates the routine completed successfully.

ESNMP\_LIB\_BAD\_REG

Indicates the subtree was not registered.

# ESNMP\_LIB\_LOST\_CONNECTION

Indicates that the request to unregister the subtree could not send. You should restart the protocol.

# 6.3.1.4 The esnmp\_poll Routine

The esnmp\_poll routine processes a pending message that has been sent by the master-agent. This routine is called after the user's select() call has indicated data is ready on the eSNMP socket. (This socket was returned from the call to the esnmp\_init routine). If no message is pending on the socket, the esnmp\_poll routine will block until one is received.

If a received message indicates a problem, an entry is made to the syslog file and an error status is returned.

If the received message is a request for SNMP data, the object table is consulted and the appropriate method routines are called.

The syntax for the esnmp\_poll routine is as follows:

### int esnmp\_poll()

The return values are as follows:

## ESNMP LIB OK

Indicates the esnmp poll routine has completed successfully.

### ESNMP\_LIB\_BAD\_REG

Indicates a previous registration was failed by the master-agent. See the log file.

### ESNMP LIB DUPLICATE

Indicates an esnmp\_init error, a duplicate subagent identifier has already been received by the master-agent.

# ESNMP\_LIB\_NO\_CONNECTION

Indicates an esnmp\_init request was failed by master-agent, restart after a delay. See the log file.

### ESNMP LIB CLOSE

Received a CLOSE message.

### ESNMP LIB NOTOK

Indicates an eSNMP protocol error occurred. The packet was discarded.

# ESNMP LIB LOST CONNECTION

Indicates that communication with master-agent was lost. Restart the connection.

## 6.3.1.5 The esnmp\_are\_you\_there Routine

The esnmp\_are\_you\_there routine requests the master-agent to respond immediately that it is up and functioning. This call does not block waiting for a response. It is intended to cause the master-agent to reply immediately. The response should be processed by calling the esnmp\_poll routine.

If no response is received within the timeout period the application code should restart the eSNMP protocol by calling the esnmp\_init routine. There are no timers maintained by the eSNMP library.

The syntax for the esnmp\_are\_you\_there routine is as follows:

int esnmp\_are\_you\_there()

The return values are as follows:

ESNMP\_LIB\_OK

The request was sent.

ESNMP\_LIB\_LOST\_CONNECTION

Cannot send the request because the master-agent is down.

# 6.3.1.6 The esnmp\_trap Routine

The esnmp\_trap routine sends a trap message to the master-agent. This function can be called at anytime. If the eSNMP protocol has not initialized with the master-agent, traps are queued and sent when communication is possible.

The trap message is actually sent to the master-agent after the master-agent's response to the esnmp\_init call has been processed. This processing happens within any API call, for most cases during subsequent calls to the esnmp\_poll routine. The quickest way actually to send traps to the master-agent is to call the esnmp\_init, esnmp\_poll, and esnmp\_trap routines.

The master-agent formats the trap into an SNMP trap message and sends it to management stations based on its current configuration. For information on configuring the master-agent see the snmpd(8) and snmpd.conf(4) reference pages.

There is no response returned from the master-agent for a trap.

The syntax for the esnmp\_trap routine is as follows:

int esnmp\_trap(int generic\_trap, int specific\_trap, char \*enterprise,
VARBIND \*vb)

The arguments are as follows:

generic\_trap
A generic trap code
specific\_trap
A specific trap code

enterprise

An enterprise OID string in dot notation.

vb

A VARBIND list of data (a NULL pointer indicates no data)

The return values are as follows:

ESNMP LIB OK

Indicates the routine completed successfully.

ESNMP\_LIB\_LOST\_CONNECTION

Indicates it could not send the trap message to master-agent.

ESNMP\_LIB\_NOTOK

Indicates something failed and message could not be generated.

# 6.3.1.7 The esnmp\_term Routine

The esnmp\_term routine sends a close message to the master-agent and shuts down the eSNMP protocol. Subagents should call this routine when terminating, so that the master-agent can update its MIB registry more quickly. It is important that terminating subagents call this routine, so that system resources used by eSNMP on their behalf can be released.

The syntax for the esnmp\_term routine is as follows:

void esnmp term(void)

The return values are:

ESNMP LIB OK

The esnmp\_term routine always returns ESNMP\_LIB\_OK, even if the packet could not be sent.

### 6.3.1.8 The esnmp sysuptime Routine

The esnmp\_sysuptime routine converts UNIX system time obtained from gettimeofday into a value with the same timebase as sysUpTime. This can be used as a TimeTicks data type (the time since the SNMP agent started) in units of 1/100 seconds. The time base is obtained from the master-agent in response to the esnmp\_init routine, so calls to this function before that time will not be accurate.

This provides a general purpose mechanism to convert UNIX timestamps into SNMP TimeTicks. The function returns the value that sysUpTime was when the passed timestamp was now. Passing a null timestamp returns the current value of sysUpTime.

The syntax is as follows:

unsigned int esnmp\_sysuptime( struct timeval \*timestamp)

The arguments are as follows:

```
struct timeval *timestamp
```

Is a pointer to struct *timeval* containing a value obtained from the gettimeofday system call. The structure is defined in include/sys/time.h.

A NULL pointer means return the current sysUpTime.

The following is an example of the esnmp\_sysuptime routine:

```
#include <include/sys/time.h>
#include <esnmp.h>
struct timeval timestamp;
gettimeofday(&timestamp, NULL);
o_integer(vb, object, esnmp_sysuptime(&timestamp));
```

The return is as follows:

0

Indicates an error (gettimeofday failed); otherwise, timestamp contains the time in 1/100ths seconds since the SNMP protocol started.

# 6.3.2 Method Routine Calling Interface

The method routine calling interface contains the following functions:

- \*\_get
- \*\_set

Section 6.3.2.3 provides additional information on method routines.

# 6.3.2.1 The \*\_get Routine

The \*\_get routine is a method routine for the specified MIB item, which is typically a MIB group (for example, system in MIB-2) or a table entry (for example, ifEntry in MIB-2). However, it is up to your discretion. See the snmpi(8) reference page for more information.

The libesnmp routines call whatever routine is specified for Get operations in the object table identified by the registered subtree.

The syntax for the \*\_get routine is as follows:

int mib\_item\_get( METHOD \*method )

# The arguments are:

#### method

A pointer to a METHOD structure, which contains the following fields:

#### action

One of ESNMP\_ACT\_GET, ESNMP\_ACT\_GETNEXT, or ESNMP ACT GETBULK.

#### serial num

An integer number that is unique to this SNMP request. Each method routine called while servicing a single SNMP request will receive the same value of <code>serial\_num</code>. New SNMP requests are indicated by a new value of <code>serial\_num</code>.

### repeat\_cnt

Used for GetBulk only. This value indicates the current iteration number of a repeating *VARBIND*. This number increments from 1 to max\_repetitions, and is 0 for nonrepeating VARBIND structures.

### max\_repetitions

For GetBulk. The maximum number of repetitions to perform. This will be 0 for nonrepeating VARBIND structures. You may be able to optimize subsequent processing by knowing the maximum number repeat calls will be made.

#### varbind

A pointer to the VARBIND structure for which we must fill in the OID and data fields. Upon entry of the method routine, the method->varbind->name is the OID that was requested.

Upon exit of the method routine, the method->varbind contains the requested data, and the method->varbind->name is updated to reflect the actual instance OID for the returned VARBIND.

The libsnmp routines (o\_integer, o\_string, o\_oid, and o\_octet) are generally used to load data. The libsnmp instance2oid routine is used to update the OID in method->varbind->name.

### object

A pointer to the object table entry for the MIB variable being referenced. The method->object->object\_index is this object's unique index within the object table (useful when one method routine services many objects).

The method->object->oid is the OID defined for this object in the MIB. The instance requested is derived by comparing this

OID with the OID in the request found in the method->varbind->name. The oid2instance function is useful for this.

row

Is not used on Get operations.

flags

Is not used on Get operations.

security

Is a pointer to security information (SNMPv2) and is not currently unused.

The return values for the \*\_get method routine are as follows:

ESNMP\_MTHD\_noError

Indicates the routine completed successfully.

ESNMP\_MTHD\_noSuchObject

The requested object cannot be returned or does not exist.

ESNMP\_MTHD\_noSuchInstance

The requested instance cannot be returned or does not exist.

ESNMP\_MTHD\_genErr

Indicates a general processing error.

## 6.3.2.2 The \* set Method Routine

The \*\_set method routine for a specified MIB item, which is typically a MIB group (for example, system in MIB-2) or a table entry (for example, ifEntry in MIB-2). However, it is up to your discretion.

The libesnmp routines call whatever routine is specified for Set operations in the object table identified by the registered subtree.

This function is pointed to by some number of elements of the subagent object table. When a request arrives for an object, its method routine is called. The \*\_set method routine is called in response to a Set SNMP request.

SNMP requests may contain many VariableBindings (encoded MIB variables). The libsnmp code executing in a subagent matches each VariableBinding with an object table entry. The object table's method routine is then called.

Therefore, a method routine is called to service a single MIB variable and the same method routine may be called several times during a single SNMP request.

The syntax for the \*\_set method routine is as follows: int mib\_item\_set( METHOD \*method )

The arguments are as follows:

method

Is a pointer to a METHOD structure, which contains the following fields:

action

The action value can be one of the following: ESNMP\_ACT\_SET, ESNMP\_ACT\_COMMIT, ESNMP\_ACT\_UNDO, or ESNMP\_ACT\_CLEANUP

serial\_num

An integer number that is unique to this SNMP request. Each method routine called while servicing a single SNMP request will receive the same value of <code>serial\_num</code>. New SNMP requests are indicated by a new value of <code>serial\_num</code>.

repeat\_cnt

This argument is not used for Set calls.

max\_repetitions

This argument is not used for Set calls.

varbind

Is a pointer to the VARBIND structure which contains the MIB variable's supplied data value and name (OID). The instance information has already been extracted from the OID and placed in method->row->instance.

object

Is a pointer to the object table entry for the MIB variable being referenced. The method->object->object\_index is this object's unique index within the object table (useful when one method routine services many objects).

The method->object->oid is the OID defined for this object in the MIB.

flags

Is a read-only integer bitmask set by libesnmp. If set, the ESNMP\_FIRST\_IN\_ROW bit indicates that this call is the first object to be set in the row. If set, the ESNMP\_LAST\_IN\_ROW bit indicates that this call is the last object to be set in the row. Only METHOD structures with the ESNMP\_LAST\_IN\_ROW bit set are passed to the method routines for commit, undo, and cleanup phases.

#### row

Is a pointer to a ROW\_CONTEXT structure (defined in the esnmp.h header file). All Set calls to the method routine which refer to the same group and have the same instance number will be presented with the same row structure. The method routines can accumulate information in the row structures during Set calls for use during the omit and undo phases. The accumulated data can be released by the method routines during the cleanup phase.

#### instance

Is an address of an array containing the instance OID for this conceptual row. The libesnmp routine builds this array by subtracting the object oid from the requested variable binding oid.

# instance\_len

Is the size of the method->row->instance.

#### context

Is a pointer to be used privately by the method routine to reference data needed to process this request.

# save

Is a pointer to be used privately by the method routine to reference data needed to potentially undo this request.

#### state

Is an integer to be used privately by the method routine to hold any state information it requires.

# security

Is pointer to security info (SNMPv2) and is not currently used.

The returns for the \*\_set method routine are as follows:

### ESNMP\_MTHD\_noError

Indicates the routine completed successfully.

## ESNMP\_MTHD\_notWritable

Indicates the requested object is not settable or was not implemented.

## ESNMP\_MTHD\_wrongLength

Indicates the requested value is the wrong length.

### ESNMP\_MTHD\_wrongEncoding

Indicates the requested value is represented incorrectly.

### ESNMP\_MTHD\_wrongValue

Indicates the requested value is out of range.

### ESNMP MTHD noCreation

Indicates the requested instance cannot ever be created.

### ESNMP\_MTHD\_inconsistentName

Indicates the requested instance cannot currently be created.

### ESNMP MTHD inconsistentValue

Indicates the requested value is not consistent.

# ESNMP MTHD resourceUnavailable

Indicates a failure due to some resource constraint.

### ESNMP MTHD genErr

Indicates a general processing error.

# ESNMP\_MTHD\_commitFailed

Indicates the commit phase failed.

# ESNMP\_MTHD\_udoFailed

Indicates the undo phase failed.

## Overall Processing of the \*\_set Routine

Every variable binding is parsed and its object is located in the object table. A METHOD structure is created for each VARBIND. These METHOD structures point to a ROW\_CONTEXT structure, which is useful for handling these phases. Objects in the same conceptual row all point to the same ROW\_CONTEXT structure. This determination is made by checking the following:

- The referenced objects are in the same MIB group
- The VARBIND structures have the same instance OIDs.

Each ROW\_CONTEXT structure is loaded with the instance information for that conceptual row. The ROW\_CONTEXT structure context and save fields are set to NULL, and the state field is set to ESNMP\_SET\_UNKNOWN structure.

The method routine for each object is called, being passed its METHOD structure with an action code of ESNMP\_ACT\_SET.

If all method routines return success, a single method routine (the last one called for the row) is called for each row, with method->action == ESNMP ACT COMMIT.

If any row reports failure, all rows that have been successfully committed are told to undo the phase. This is accomplished by calling a single method routine for each row (the same one that was called for the commit phase), with a method->action == ESNMP\_ACT\_UNDO.

Finally, each row is released. The same single method routine for each row is called with a method->action == ESNMP\_ACT\_CLEANUP. This occurs for every row, regardless of the results of previous processing.

### **ESNMP ACT SET**

Each object's method routine is called during the Set phase, until all objects are processed or a method routine returns an error status value. (This is the only phase during which each object's method routine is called.) For variable bindings in the same conceptual row, method->row points to a common ROW CONTEXT.

The method->flags bitmask have the ESNMP\_LAST\_IN\_ROW bit set, if this is the last object being called for this ROW\_CONTEXT. This enables you to do a final consistency check, since you have seen every variable binding for this conceptual row.

The method routine's job in this phase is to determine if the SetRequest will work, return the correct SNMP error code if not, and prepare any context data it needs to actually perform the Set during the commit phase.

The method->row->context is private to the method routine; libesnmp does not use it. A typical use is to store the address of an emitted foo\_type structure that has been loaded with the data from the VARBIND for the conceptual row.

# ESNMP\_ACT\_COMMIT

Even though several variable bindings may be in a conceptual row, only the last one in order of the SetRequest is processed. So, for all the method routines that point to a common row, only the last method routine is called.

This method routine must have available to it all necessary data and context to perform the operation. It must also save a snapshot of current data or whatever it needs to undo the Set if required. The method->row->save is intended to hold a pointer to whatever data is needed to accomplish this. A typical use is to store the address of an emitted foo\_type structure that has been loaded with the current data for the conceptual row.

The method->row->save is also private to the method routine; libesnmp does not use it.

If the set operation succeeds, return ESNMP\_MTHD\_noError; otherwise, back out the commit as best you can and return a value of ESNMP\_MTHD\_commitFailed.

If any errors were returned during the commit phase, libesnmp enters the undo phase; if not, it enters the cleanup phase.

#### Note

The undo phase may occur even if the Set operation in your subagent is successful because the SetRequest spanned subagents and a different subagent failed.

### ESNMP\_ACT\_UNDO

For each conceptual row that was successfully committed, the same method routine is called with method->action == ESNMP\_ACT\_UNDO. The ROW\_CONTEXT structures that have not yet been called for the commit phase are not called for the undo phase; they are called for cleanup phase.

The method routine should attempt to restore conditions to what they were before it executed the commit phase. (This is typically done using the data pointed to by the method->row->save.)

If successful, return ESNMP\_MTHD\_noError; otherwise, return ESNMP\_MTHD\_undoFail.

## ESNMP ACT CLEANUP

Regardless of what else has happened, at this point each ROW\_CONTEXT participates in cleanup phase. The same method routine that was called for commit phase is called with method->action == ESNMP\_ACT\_CLEANUP.

This indicates the end of processing for the SetRequest. The method routine should perform whatever cleanup is required; for instance, freeing dynamic memory that might have been allocated and stored in method->row->context and method->row->save, and so on.

The function return status value is ignored for the cleanup phase.

### 6.3.2.3 Method Routines

You must write the code for the method routines declared in the *subtree\_tbl.*h file. Each method routine has one argument, which is a pointer to the METHOD structure, as follows:

int mib-group\_get(METHOD \*method)
int mib-group\_set(METHOD \*method)

The Get method routines are used to perform Get, GetNext, and GetBulk operations.

Get method routines perform the following tasks:

1. Extract the instance portion of the requested OID. You can do this manually by comparing <code>method->object->oid</code> (the object's base OID) to <code>method->varbind->name</code> (the requested OID). You can use the oid2instance libesnmp routine to do this.

- 2. Determine the instance validity. The instance OID may be null or any length, depending on what was requested and how your object was selected. You may be able to fail the request immediately by checking on the instance OID.
- 3. Extract the data. Based on the instance OID and method->action, determine what data, if any, is to be returned.
- 4. Load the response OID back into the method routine's VARBIND. Set the *method->varbind* with the OID of the actual MIB variable instance you are returning. This is usually accomplished by loading an array of integers with the instance OID you wish to return and calling the instance20ID libesnmp routine.
- 5. Load the response data back into the method routine's VARBIND.

Use one of the libesnmp library routine with the corresponding data type to load the *method->varbind* with the data to return:

- o\_integer
- o\_string
- o octet
- o\_oid

These routines make a copy of the data you specify. The libesnmp function manages any memory associated with copied data. The method routine must manage the original data's memory.

The routine does any necessary conversions to the type defined in the object table for the MIB variable and copies the converted data into method->varbind.

See the Value Representation section for information on data value representation.

- 6. Return the correct status value, as follows:
  - ESNMP\_MTHD\_noError The routine completed successfully or no errors were found.
  - ESNMP\_MTHD\_noSuchInstance

For SNMPV1 – Returned as an error code.

For SNMPV2 – Translated to a noSuchInstance exception.

ESNMP\_MTHD\_noSuchObject

For SNMPV1 – Returned as a noSuchInstance error.

For SNMPv2 – Translated as a noSuchObject exception

 ESNMP\_MTHD\_genErr – An error occurred and the routine did not complete successfully.

# **Value Representation**

The values in a VARBIND for each data type are represented as follows. (Refer to the esnmp.h file for a definition of the OCT and OID structures.)

• ESNMP\_TYPE\_Integer32 (varbind->value.sl)

This is a 32-bit signed integer. Use the o\_integer routine to insert an integer value into the VARBIND. Note that the prototype for the value argument is unsigned long, so you may need to cast this to a signed int.

 ESNMP\_TYPE\_DisplayString, ESNMP\_TYPE\_NsapAddress, ESNMP\_TYPE\_Opaque, ESNMP\_TYPE\_OctetString (varbind->value.oct)

This is an octet string. It is contained in the VARBIND as an OCT structure that contains a length and a pointer to a dynamically allocated character array. Included on the end of the character array is a null terminator that is not included in the length.

The DisplayString is different only in that the character array can be interpreted as ASCII text where the OctetString can be anything.

Use the o\_string routine to insert a value into the VARBIND which is a buffer and a length. New space will be allocated and the buffer copied into the new space.

Use the o\_octet routine to insert a value into the VARBIND, which is a pointer to an OCT structure. New space is allocated and the buffer pointed to by the OCT structure is copied.

• ESNMP\_TYPE\_ObjectId (varbind->value.oid and the varbind->name field)

This is an object identifier. It is contained in the VARBIND as an OID structure which contains the number of elements and a pointer to a dynamically allocated array of unsigned integers, one for each element.

The *varbind->name* field is used to hold the object identifier and instance information that identifies MIB variable. Use the OID2Instance function to extract the instance elements from an incoming OID on a request. Use the Instance20ID function to combine the instance elements with the MIB variable's base OID to set the VARBIND structure's *name* field when building a response.

Use the o\_oid function to insert an object identifier into the VARBIND when the OID value to be returned as data is in the form of a pointer to an OID structure.

Use the o\_string function to insert an object ID into the VARBIND when the OID value to be returned as data is in the form of a pointer to an ASCII string containing the OID in dot format; for example

1.3.6.1.2.1.3.1.1.2.0.

### • ESNMP TYPE NULL

This is the NULL or empty type. This is used to indicate that there is no value. The length is 0 and the value union in the VARBIND is zero-filled.

The incoming VARBIND structures on a Get, GetNext, and GetBulk will have this data type. A method routine should never return such a value. An incoming Set request never has such a value in a VARBIND.

• ESNMP\_TYPE\_IpAddress (varbind->value.oct)

This is an IP address. It is contained in the VARBIND in an OCT structure which has a length of 4 and a pointer to a dynamically allocated buffer containing the 4 bytes of the IP address in network order.

Use the o\_integer function to insert an IP address into the VARBIND when the value is an unsigned integer in network byte order.

Use the o\_string function to insert an IP address into the VARBIND when the value is a byte array (in network byte order). Use a length of 4.

• ESNMP\_TYPE\_UInteger32 ESNMP\_TYPE\_Counter32 ESNMP TYPE Gauge32 (varbind->value.ul)

The 32-bit counter and 32-bit gauge data types are stored in the VARBIND as an unsigned int.

Use the o\_integer function to insert an unsigned value into the VARBIND.

• ESNMP TYPE TimeTicks (varbind->value.ul)

The 32-bit timeticks type values are stored in the VARBIND as an unsigned int, in .01-second units.

Use the o\_integer function to insert an unsigned value into the VARBIND.

• ESNMP\_TYPE\_BitString (varbind->value.oct)

The BitString is contained in the VARBIND as an OCT structure which contains a length equal to the number of bytes needed to contain the value which is ((qty-bits - 1)/8 + 2), and a pointer to a dynamically buffer containing the bits of the bitstring in the form uubbbbb..bb, where the first octet (uu) is 0x00-0x07 and indicates the number of unused bits in the last octet (bb). The bb octets represent the bit string itself, where bit zero (0) comes first and so on.

Use the o\_octet routine to insert a value into the VARBIND which is a pointer to an OCT structure pointing to a buffer containing the bits in the *uubbbbb..bb* form. New space will be allocated and the buffer pointed to by the OCT structure will be copied.

This is not compatible with SNMPv1. It will be returned or set only for

SNMPv2 requests.

• ESNMP\_TYPE\_Counter64 (varbind->value.u164)

The 64-bit counter is stored int a VARBIND as an unsigned long which, on an Alpha machine, has a 64-bit value.

Use the o\_integer function to insert an unsigned long value (64 bits) into the VARBIND.

This is not compatible with SNMPv1. It is returned or set for SNMPv2 requests only.

# 6.3.3 The libsnmp Support Routines

This section provides information on the libsnmp support routines, which consists of the following:

- o\_integer
- o\_octet
- o\_oid
- o\_string
- str2oid
- sprintoid
- instance2oid
- oid2instance
- inst2ip
- cmp\_oid
- cmp\_oid\_prefix
- clone\_oid
- free\_oid
- clone buf
- mem2oct
- cmp\_oct
- clone\_oct
- free\_oct
- free\_varbind\_data
- set\_debug\_level
- is\_debug\_level

• ESNMP LOG

# 6.3.3.1 The o\_integer Routine

The o\_integer routine loads an integer value into the VARBIND with the appropriate type.

The syntax is as follows:

int o\_integer( VARBIND \*vb, OBJECT \*obj, unsigned long value)

The arguments are as follows:

VARBIND \*vb

Is a pointer to the VARBIND structure which is to receive the data. This function does not allocate the VARBIND structure.

OBJECT \*obi

Is a pointer to the OBJECT structure for the MIB variable associated with the OID in the VARBIND.

unsigned long value

The value to be inserted into the VARBIND.

The real type as defined in the object structure must be one of the following; otherwise, an error is returned.

If the real type is IpAddress, then it assumes that the 4-byte integer is in network byte order and will be packaged into one of the following octet strings:

ESNMP TYPE Integer32:

32-bit INTEGER

ESNMP\_TYPE\_Counter32:

32-bit Counter (unsigned)

ESNMP\_TYPE\_Gauge32:

32-bit Gauge (unsigned)

ESNMP TYPE TimeTicks:

32-bit TimeTicks (unsigned)

ESNMP TYPE UInteger32:

32-bit INTEGER (unsigned)

ESNMP\_TYPE\_Counter64:

64-bit Counter (unsigned)

ESNMP TYPE IpAddress:

IMPLICIT OCTET STRING (4)

The following is an example of the o\_integer routine:

The following are the return values:

ESNMP\_MTHD\_noError

The routine completed successfully.

ESNMP\_MTHD\_genErr

An error has occurred.

### 6.3.3.2 The o\_octet Routine

The o\_octet routine loads an octet value into the VARBIND with the appropriate type.

The syntax is as follows:

```
int o_octet(VARBIND *vb, OBJECT *obj, OCT *oct)
```

The arguments are as follows:

```
VARBIND *vb
```

Is a pointer to the VARBIND structure which is to receive the data. This function does not allocate the VARBIND structure.

### Note

If the original value in the varbind *vb* is not NULL, this routine attempts to free it. So if you malloc your own *vb* structure, be sure to fill it with zeros before using it.

```
OBJECT *obj
```

Is a pointer to the OBJECT structure for the MIB variable associated with the OID in the VARBIND.

```
OCT *value
```

Is the value to be inserted into the VARBIND.

The real type as defined in the object structure must be one of the following; otherwise, an error is returned:

```
ESNMP_TYPE_OCTET_STRING
OCTET STRING (ASN.1)
```

ESNMP\_TYPE\_IpAddress

IMPLICIT OCTET STRING (4) – in octet form, network byte order

ESNMP\_TYPE\_DisplayString
DisplayString (Textual Con)

ESNMP\_TYPE\_NsapAddress IMPLICIT OCTET STRING

ESNMP\_TYPE\_Opaque IMPLICIT OCTET STRING

ESNMP\_TYPE\_BIT\_STRING

BIT STRING (ASN.1) – The first byte is the number of unused bits in the last byte.

The following is an example of the o\_octet routine:

The returns are as follows:

ESNMP MTHD noError

Indicates that the routine completed successfully.

ESNMP\_MTHD\_genErr

Indicates that an error condition has occurred.

# 6.3.3.3 The o\_oid Routine

The o\_oid routine loads an OID value into the VARBIND with the appropriate type.

The syntax is as follows:

int o\_oid(VARBIND \*vb, OBJECT \*obj, OID \*oid)

The arguments are as follows:

VARBIND \*vb

Is a pointer to the VARBIND structure that is to receive the data. This

function does not allocate the VARBIND structure.

### Note

If the original value in the varbind *vb* is not NULL, this routine attempts to free it; therefore, if you malloc your own *vb* structure, fill it with zeros (0s) before using it.

```
OBJECT *obj
```

Is a pointer to the OBJECT structure for the MIB variable associated with the oid in the VARBIND.

```
OID *value
```

Is the value to be inserted into the VARBIND structure as data.

The real type as defined in the object structure must be the following; otherwise, an error is returned:

```
ESNMP_TYPE_OBJECT_IDENTIFIER
OBJECT IDENTIFIER (ASN.1)
```

The following is an example of the o oid routine:

The returns are as follows:

ESNMP\_MTHD\_noError

Indicates the routine ended successfully.

ESNMP\_MTHD\_genErr

Indicates an error condition has occurred.

## 6.3.3.4 The o\_string Routine

The o\_string routine loads a string value into the VARBIND with the appropriate type.

The syntax is as follows:

int o\_string( VARBIND \*vb, OBJECT \*obj, unsigned char \*ptr, int len)

The arguments are as follows:

VARBIND \*vb

Is a pointer to the VARBIND structure which is to receive the data. This function does not allocate the VARBIND structure.

### Note

If the original value in the varbind *vb* is not NULL, this routine attempts to free it; therefore, if you malloc your own *vb* structure, fill it with zeros (0s) before using it.

OBJECT \*obi

Is a pointer to the OBJECT structure for the MIB variable associated with the oid in the VARBIND.

unsigned char \*ptr

Is the pointer to the buffer containing data to be inserted into the VARBIND as data.

int len

Is the length of the data in buffer to which ptr points.

The real type as defined in the object structure must be one of the following; otherwise, an error is returned:

ESNMP\_TYPE\_OCTET\_STRING OCTET STRING (ASN.1)

ESNMP\_TYPE\_IpAddress

IMPLICIT OCTET STRING (4) – in octet form, network byte order

ESNMP\_TYPE\_DisplayString
DisplayString (Textual Con)

ESNMP\_TYPE\_NsapAddress IMPLICIT OCTET STRING

ESNMP\_TYPE\_Opaque IMPLICIT OCTET STRING

ESNMP\_TYPE\_BIT\_STRING

BIT STRING (ASN.1) – The binary value of first byte is the number of unused bits in the last byte.

ESNMP\_TYPE\_OBJECT\_IDENTIFIER
OBJECT IDENTIFIER (ASN.1) – in dot notation, 1.3.4.6.3

The following is an example of the o string routine:

The return values are as follows:

```
ESNMP_MTHD_noError
```

Indicates that the routine completed successfully.

```
ESNMP_MTHD_genErr
```

Indicates that an error condition has occurred.

#### 6.3.3.5 The str2oid Routine

The str2oid routine converts a null-terminated OID string (in dot notation) to an OID structure.

It dynamically allocates the elements buffer and inserts its pointer into the OID structure passed in. It is the responsibility of the caller to free this buffer. The OID can have a maximum of 128 elements.

Note that the str20id routine does not allocate an OID structure.

The syntax is as follows:

```
OID * str2oid (OID *oid, char *s)
```

The following is an example of the str20id routine:

The returns are as follows:

**NULL** 

Indicates an error has occurred; otherwise, the pointer to the OID structure (its first argument) is returned.

# 6.3.3.6 The sprintoid Routine

The sprintoid routine converts an OID into a null-terminated string in dot notation. An OID can have up to 128 elements. A full sized OID can require a large buffer.

The syntax is as follows:

```
char *sprintoid ( char *buffer, OID *oid)
```

The following is an example of the sprintoid routine:

```
#include <esnmp.h>
#define SOMETHING_BIG 1024
OID abc;
char buffer[SOMETHING_BIG];
:
: assume abc gets initialized with some value:
printf("dots=%s0, sprintoid(buffer, &abc));
```

The return values are its first argument.

#### 6.3.3.7 The instance2oid Routine

The instance2oid routine makes a copy of the object's base OID and appends a copy of the instance array to make a complete OID for a value. The instance is an array of integers and len is the number of elements. The instance array may be created by oid2instance or constructed from key values as a result of a get\_next search.

It dynamically allocates the elements buffer and inserts its pointer into the OID structure passed in. The caller is responsible for freeing this buffer.

Point to the OID structure that is to receive the new OID values and call this routine. Any previous value in the OID structure is freed (it calls free\_oid first) and the new values are dynamically allocated and inserted. Be sure the initial value of the new OID is all zeros, if you do not want it to be freed.

Note that the instance2oid routine does not allocate an OID structure, only the array containing the elements.

The syntax is as follows:

```
OID * instance2oid ( OID *new, OBJECT *obj, unsigned int *instance, int len)
```

The arguments are as follows:

```
OID *new
```

Is a pointer to the OID that is to receive the new OID value.

```
OBJECT *obj
```

Is a pointer to the object table entry for the MIB variable being obtained. The first part of the new OID is the OID from this MIB object table entry.

```
unsigned int *instance
```

Is a pointer to an array of *instance* values. These values are appended to the base OID obtained from the MIB object table entry to construct the new OID.

int len

Is the number of elements in the *instance* array.

The following is an example of the instance2oid routine:

The returns are as follows:

**NULL** 

Indicates an error has occurred; otherwise, the pointer to the OID (its first argument) is returned.

#### 6.3.3.8 The oid2instance Routine

The oidlinstance routine extracts the instance values from an OID and copies them to the specified array of integers. It then returns the number of elements in the array. The instance is the elements of an OID beyond those elements that identify the MIB variable. They are used as indexes to identify a specific instance of a MIB value.

If there are more elements in the OID than expected (more than specified by the max\_len parameter), the function copies the number of elements specified by max\_len only and returns the total number of elements that would have been copied had there been space.

The syntax is as follows:

int oid2instance ( OID \*oid, OBJECT \*obj, unsigned int \*instance, int max len)

The arguments are as follows:

oid

Is an incoming OID containing an instance or part of an instance.

obj

Is a pointer to the object table entry for the MIB variable.

ingtance

Is a pointer to an array of unsigned integers where the index will be placed.

max len

Is a number of elements available in the instance array.

The N will be in instance[0] and the IP address will be in instance[2], instance[3], instance[4], and instance[5].

The returns are as follows:

- Less than zero indicates that an error, should not be if the object was obtained by looking at this oid.
- Zero indicates there are no instance elements.
- Any positive integer indicates the number of elements in the index. (This could be larger than the max\_len parameter).

#### 6.3.3.9 The inst2ip Routine

The inst2ip routine returns an IP address derived from an OID instance. For evaluation of an instance for Get and Set operations use the EXACT mode. For GetNext and GetBulk operations use the NEXT mode. When using the NEXT mode, this routine's logic assumes that the search for data will be performed using greater than or equal to matches.

The syntax is as follows:

int inst2ip( unsigned int \*inst, int length, unsigned int \*ipAddr, int exact, int carry)

The arguments are as follows:

#### ingt

Is a pointer to an array of unsigned int containing the instance numbers returned by the oid2instance routine to be converted to an IP address.

Each element is in the range 0 to 255. Using the EXACT mode, the routine returns 1 if an element is out of range. Using NEXT mode, a value greater than 255 causes that element to overflow. It is set to 0 and the next most significant element is incremented, so it returns a lexically equivalent value of the next possible ipAddress.

#### length

Is the number of elements in the instance array. Instances beyond the fourth are ignored. If the length is less than 4, the missing values are assumed to be 0. A negative length results in an ipaddr value of 0. For an exact match (such as Get) there must be at exactly four elements.

#### ipAddr

Is a pointer to where to return the IP address value. It is in network byte order; that is, the most significant element is first.

#### exact

Can be either TRUE or FALSE.

TRUE means do an EXACT match. If any element is greater than 255 or if there are not exactly 4 elements, return 1. The carry argument is ignored.

FALSE means do a NEXT match. That is, return the lexically next IP address if the carry is set and the length is at least 4. If there are fewer than 4 elements, assume the missing values are 0. If any one element contains a value greater than 255, then zero the value and increment the next most significant element. Return 1 only in the case where there is a carry from the most significant (the first) value.

#### carry

Is the carry to add to the IP address on a NEXT match. If you are trying to determine the next possible IP address, pass in a 1; otherwise, pass in a 0. A length of less than 4 cancels the carry.

The following are examples of the inst2ip routine.

The following example converts an instance to an IP address for a Get

operation, which is an EXACT match.

The following example shows a GetNext where there is only one key or that the *ipaddr* is the least significant part of the key. This is a NEXT match; therefore, a 1 is passed in for *carry*.

In the following example, if there is more than one part to a search key and you are doing a GetNext, you want to find the next possible value for the search key so you can do a cascaded greater-than or equal-to search.

If you have a search key of a number and two <code>ipAddr</code> values that are represented in the instance part of the OID as <code>N.A.A.A.B.B.B.B</code> with <code>N</code> as single valued integer and <code>A.A.A.A</code> portion making up one IP address and the <code>B.B.B.B</code> portion making up a second IP address and a total length of 9 if all elements are given, you start by converting the least significant part of the key, (that would be the <code>B.B.B.B</code> portion). You do that by calling the <code>inst2ip</code> routine passing in a 1 for the carry and 5 for the length. If the conversion of the <code>B.B.B.B</code> portion generated a carry (returned 1), you will pass it on to the next most significant part of the key. Therefore, convert the <code>A.A.A.A</code> portion by calling the <code>inst2ip</code> routine, passing in 1 for the length and the carry returned from the conversion of the <code>B.B.B.B</code> portion. The most significant element <code>N</code> is a number; therefore, add the carry from the

A conversion to the number. If that also overflows, then this is not a valid search key.

```
#include <esnmp.h>
            *incoming = &method->varbind->name;
OID
OBJECT
           *object = method->object;
int instLength;
unsigned int instance[9];
unsigned int ip_addrA;
unsigned int ip_addrB;
int
           iface;
int
         carry;
-- The instance is N.A.A.A.B.B.B.B --
instLength = oid2instance(incoming, object, instance, 9);
iface = (instLength < 1) ? 0 :(int) instance[0];</pre>
carry = inst2ip(&instance[1],instLength - 1,&ip_addr,FALSE,1);
carry = inst2ip(&instance[5],instLength - 5,&ip_addr,FALSE,carry);
iface += carry;
if (iface > 0xFFFFFFFF)
-- a carry caused an overflow in the most significant element
return ESNMP_MTHD_noSuchInstance;
```

The returns are as follows:

- If the *carry* is 0, the routine completed successfully.
- If the *carry* equals 1, it indicates an error if EXACT match or there was a carry for a NEXT match. If there was a carry, the returned *ipAddr* is 0.

#### 6.3.3.10 The cmp oid Routine

The cmp\_oid routine compares two OID structures. This routine does an element-by-element comparison starting with the most significant element (element 0) and working toward the least significant element. If all other elements are equal, the OID with the fewest elements is considered less.

The syntax is as follows:

```
int cmp_oid( OID *q, OID *p)
```

The returns are as follows:

- +1 Indicates that oid q is greater than oid p.
- 0 Indicates that oid q is in oid p.
- -1 Indicates that oid q is less than oid p.

#### 6.3.3.11 The cmp\_oid\_prefix Routine

The cmp\_oid\_prefix routine compares an OID against a prefix. A prefix could be the OID on an object in the object table. The elements beyond the prefix are the instance information.

This routine does an element-by-element comparison, starting with the most significant element (element 0) and working toward the least significant element. If all elements of the prefix OID match exactly with corresponding elements of OID q, it is considered an even match if OID q contains additional elements. OID q is considered greater than the prefix if the first nonmatching element is larger. It is considered smaller if the first nonmatching element is less.

The syntax is as follows:

```
int cmp_oid_prefix( OID *q, OID *prefix)
```

The following is and example of the cmp\_oid\_prefix routine:

```
#include <esnmp.h>
OID *q;
OBJECT *object;
if (cmp_oid_prefix(q, &object->oid) == 0)
    printf("matches prefix0);
```

The returns are as follows:

- -1 Indicates the oid is less than the prefix.
- 0 Indicates the oid is in the prefix.
- +0 Indicates the oid is greater than the prefix.

#### 6.3.3.12 The clone\_oid Routine

The clone oid routine makes a copy of the OID structure.

Pass in a pointer to the source OID structure to be cloned and a pointer to the new OID structure that is to receive the duplicated OID values.

It dynamically allocates the element's buffer and inserts its pointer into the OID structure passed in.

It is the responsibility of the caller to free this buffer.

Note that any previous elements buffer pointed to by the new OID structure will be freed and pointers to the new, dynamically allocated, buffer will be inserted. Be sure to initialize the new OID structure with zeroes (0), unless it contains an element buffer that can be freed.

Also note that this routine does not allocate an OID structure.

The syntax is as follows:

```
OID *clone_oid ( OID *new, OID *oid)
```

The arguments are as follows:

```
OID *new
```

Is a pointer to the OID structure that is to receive the copy.

```
OID *old
```

Is a pointer to the OID structure where the data is to be obtained.

The following is an example of the clone oid routine:

The returns are as follows:

**NULL** 

Indicates an error; otherwise, the pointer to the OID (its first argument) is returned.

#### 6.3.3.13 The free\_oid Routine

The free oid routine frees an OID structure's elements buffer.

It frees the buffer pointed to by oid->elements then zeros that field and oid->nelem.

Note that this routine does not deallocate the OID structure itself, only the elements buffer attached to it.

The syntax is as follows:

```
void free oid (OID *oid)
```

The following is an example of the free\_oid routine:

```
#include <esnmp.h>
OID oid;
:
: assume oid was assigned a value (perhaps with clone_oid()
: and we are now finished with it.
:
free_oid(&oid);
```

#### 6.3.3.14 The clone\_buf Routine

The clone\_buf routine duplicates a buffer in a dynamically allocated space. One extra byte is always allocated on end and filled with \0. If the length is less than 0, its length is set to 0. There is always a buffer pointer, unless there is a malloc error.

It is the callers responsibility to free the allocated buffer.

The syntax is as follows:

```
char *clone_buf( char *str, int len)
```

The arguments are as follows:

str

Is a pointer to the buffer to be duplicated.

1en

Is a number of bytes to copy.

The following is an example of the clone\_buf routine:

```
#include <esnmp.h>
char *str = "something nice";
char *copy;
copy = clone_buf(str, strlen(str));
```

The returns are as follows:

**NULL** 

Indicates a malloc error; otherwise, the pointer to allocated buffer containing a copy of the original buffer is returned.

#### 6.3.3.15 The mem2oct Routine

The mem2oct routine converts a string, (a buffer and length) to an OCT structure.

It dynamically allocates a new buffer, copies the indicated data into it, and updates the OCT structure with the new buffer's address and length.

It is the responsibility of the caller to free the allocated buffer.

Note this routine does not allocate an OCT structure and that it does not free data previously pointed to in the OCT structure before making the assignment.

The syntax is as follows:

OCT \* mem2oct( OCT \*new, char \*buffer, int len)

The following is an example of the mem2oct routine:

The following are the return values:

NULL

Indicates an error; otherwise, the pointer to the OCT structure (its first argument) is returned.

#### 6.3.3.16 The cmp\_oct Routine

The cmp\_oct routine compares two octets. The two octets are compared byte-by-byte for the length of the shortest octet. If all bytes are equal, the lengths are compared. An octet with a null pointer is considered the same as a zero-length octet.

The syntax is as follows:

```
int cmp oct ( OCT *oct1, OCT *oct2)
```

The following is an example of the cmp\_oct routine:

The returns are as follows:

- -1 The string to which the first octet points is less than the second.
- 0 The string to which the first octet points is equal to the second.
- +1 The string to which the first octet points is greater than the second.

#### 6.3.3.17 The clone oct Routine

The clone\_oct routine makes a copy of the OCT structure.

It passes in a pointer to the source OCT structure to be cloned and a pointer to the new OCT structure that is to receive the duplicated OCT structure's values.

It dynamically allocates the buffer, copies the data, and updates the new OCT structure with the buffer's address and length.

It is the responsibility of the caller to free this buffer.

Note that any previous buffer to which the new OCT structure points is freed and pointers to the new, dynamically allocated buffer are inserted. Be sure to initialize the new OCT structure with zeros (0), unless it contains a buffer that can be freed.

Also note that this routine does not allocate an OCT structure, only the elements buffer pointed to by the OCT structure.

The syntax is as follows:

```
OCT * clone_oct ( OCT *new, OCT *old)
```

The arguments are as follows:

```
OCT *new
```

Is a pointer to the OCT structure that is to receive the copy.

```
OCT *old
```

Is a pointer to the OCT structure where the data is to be obtained.

The following is an example of the routine:

The returms are as follows:

**NULL** 

Indicates an error; otherwise, the pointer to the OCT structure (its first argument) is returned.

#### 6.3.3.18 The free\_oct Routine

The free\_oct routine frees the buffer attached to the OCT structure.

It frees a dynamically allocated buffer to which the OCT structure points, then zeros (0) the pointer and length fields in the OCT structure. If the buffer is already NULL this routine does nothing.

Note that this routine does not deallocate the OCT structure, only the buffer to which it points.

The syntax is as follows:

```
void free_oct ( OCT *oct)
```

The following is an example of the free\_oct routine:

```
#include <esnmp.h>
OCT octet;
:
: assume octet was assigned a value (perhaps with mem2oct()
: and we are now finished with it.
:
free_oct(&octet);
```

#### 6.3.3.19 The free\_varbind\_data Routine

The free\_varbind\_data routine frees the dynamically allocated fields within the VARBIND structure.

The routine performs a free\_oid ( $vb \rightarrow name$ ) operation. If the vb->typefield indicates, it then frees the vb->value data using either the free oct or the free oid routine.

It does not deallocate the VARBIND structure itself; only the name and data buffers to which it points.

The syntax is as follows:

```
void free_varbind_data( VARBIND *vb)
```

The following is an example of the free varbind data routine:

```
#include <esnmp.h>
VARBIND *vb;

vb = (VARBIND*)malloc(sizeof(VARBIND));
clone_oid(&vb->name, oid);
clone_oct(&vb->value.oct, data);
    .
    .
    free_varbind_data(vb);
free(vb);
```

#### 6.3.3.20 The set debug level Routine

The set\_debug\_level routine sets the logging level which dictates what log messages are generated. You should call the routine during program initialization in response to runtime options. If not called, this will be set to WARNING and ERROR messages to stdout as the default.

The following values can be set:

- ERROR For when a bad error occurred, requiring a restart.
- WARNING For when a packet cannot be handled; this also implies ERROR.
- TRACE For when tracing all packets; this also implies ERROR and WARNING.
- DAEMON\_LOG Causes output to go to syslog rather than to standard output.
- EXTERN\_LOG Causes the callback function to be called to output log messages. If this bit is set, you must provide the second argument, which is a pointer to a user supplied external callback function. If DAEMON\_LOG and EXTERN\_LOG are not specified, output goes to standard output.
- callback A user-supplied external callback function:

void callback\_function( int level, char \*message)

The <code>level</code> will be <code>ERROR</code>, <code>WARNING</code>, or <code>TRACE</code>. If the <code>EXTERN\_LOG</code> bit is set in <code>stat</code>, the <code>callback</code> function will be called whenever an <code>ESNMP\_LOG</code> macro is executed and the log level indicates that a log message is to be generated.

This facility allows an implementer to control where eSNMP library functions output log messages. If EXTERN\_LOG bit will not be set, pass in a NULL pointer for the callback function argument.

The syntax is as follows:

void set\_debug\_level(int stat, LOG\_CALLBACK\_ROUTINE callback\_routine)

The following is an example of the set\_debug\_level routine:

```
#include <esnmp.h>
extern void log_handler(int level, char *message);

if (daemonize)
    set_debug_level(EXTERN_LOG | WARNING, log_handler);
else
    set_debug_level(TRACE, NULL);
```

#### 6.3.3.21 The is debug level Routine

The is\_debug\_level routine tests the log level to see if the specified level is set. You can set the levels as follows:

• ERROR – For when a bad error occurred, requiring restart.

- WARNING For when a packet cannot be handled.
- TRACE For when tracing all packets.
- DAEMON\_LOG For output going to syslog.
- EXTERN\_LOG For the *callback* function is to be called to output log messages.

The syntax is as follows:

```
int is_debug_level( int type)
```

The return values are as follows:

TRUE

The requested level is set and the ESNMP\_LOG will generate output, or output will go to the specified destination.

FALSE

The is\_debug\_level routine is not set.

The following is an example of the is\_debug\_level routine:

```
#include <esnmp.h>
if (is_debug_level(TRACE))
    dump_packet();
```

#### 6.3.3.22 The ESNMP\_LOG Routine

The ESNMP\_LOG routine is an error declaration C macro defined in the <esnmp.h> header file. It gathers the information that it can obtain and sends it to the log. If DAEMON\_LOG is set, log messages are sent to the daemon log. If EXTERN\_LOG is set, log messages are sent to the callback function; otherwise, log messages go to standard output.

#### Note

The esnmp\_log routine is called using the ESNMP\_LOG macro, which uses the helper routine esnmp\_logs to format part of the text. Do not use these functions without the ESNMP\_LOG macro.

```
level
    Can be one of the following:
    ERROR
        Declares an error condition.

WARNING
        Declares a warning.

TRACE
        Put in log file if trace is active.

The syntax is as follows:

ESNMP_LOG( level, ( format, ...))

The following is an example of the ESNMP_LOG routine:
#include <esnmp.h>
ESNMP_LOG( ERROR, ("Cannot open file %s\n", file));
```

# Digital UNIX STREAMS/Sockets Coexistence 7

This chapter describes the ifnet STREAMS module and dlb STREAMS pseudodriver communication bridges. Before reading it, you should be familiar with basic STREAMS and sockets concepts and have reviewed the information in Chapter 4 and Chapter 5.

The Digital UNIX network programming environment supports the STREAMS and sockets frameworks for network programming. However, there is no native communication path at the data link layer between the two frameworks. The term **coexistence** refers to the ability to exchange data between the sockets and STREAMS frameworks. The term **communication bridge** refers to the software (ifnet STREAMS module or the dlb STREAMS pseudodriver) that enables the two frameworks to exchange data at the data link layer.

Programs written to sockets and STREAMS must intercommunicate for the following reasons:

- A system cannot have two drivers for the same device.
- Programs may need to access STREAMS-based device drivers from BSD protocol stacks or, conversely, may need to access BSD device drivers from STREAMS-based protocol stacks.

For example, if your system is running a STREAMS device driver and you have an application that uses the TCP/IP implemented on Digital UNIX, which is sockets-based, you need a path by which the data gets from the sockets-based protocols stack to the STREAMS device driver and back again. The ifnet STREAMS module allows an application using TCP/IP to exchange data with a STREAMS device driver. Section 7.1 describes the ifnet STREAMS module.

Conversely, if you have a STREAMS protocol stack implemented on your system but want to use the BSD device driver implemented on Digital UNIX, you need a path by which the data gets from the STREAMS protocol stack to the BSD device driver and back again. The dlb STREAMS pseudodriver allows the STREAMS protocol stack to route its data to the BSD device driver. Section 7.2 describes the dlb STREAMS pseudodriver.

# 7.1 Bridging STREAMS Drivers to Sockets Protocol Stacks

The ifnet STREAMS module is a communication bridge that allows STREAMS network drivers to access sockets-based network protocols. The ifnet STREAMS module functions like any other STREAMS module, being pushed on the Stream above the STREAMS device driver. Once it is on the Stream, it handles all of the translation required between the DLPI interface of the STREAMS driver and the BSD ifnet layer. The ifnet STREAMS module exports both standard STREAMS interfaces as well as ifnet layer interfaces.

Note that STREAMS network drivers can also continue to use STREAMS-based network protocols while using the ifnet STREAMS module.

Figure 7-1 highlights the ifnet STREAMS module and shows its place in the network programming environment.

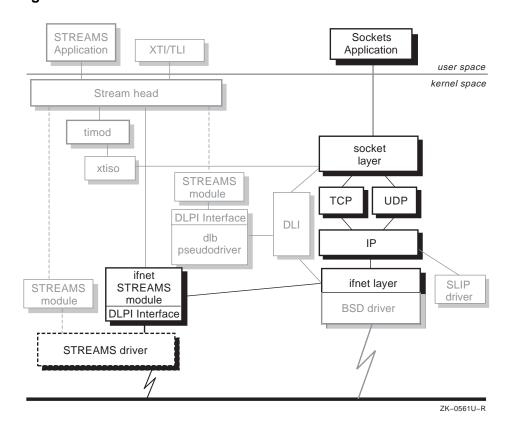


Figure 7-1: The ifnet STREAMS module

#### 7.1.1 The STREAMS Driver

This section describes how to prepare the system running the STREAMS driver to use the ifnet STREAMS module.

#### Note

The ifnet STREAMS module only supports Ethernet STREAMS device drivers.

This section also lists the DLPI primitives that the STREAMS driver must support in order for the ifnet STREAMS module to operate successfully.

#### 7.1.1.1 Using the ifnet STREAMS Module

If your device driver supports the primitives listed in Section 7.1.1.2, no source code changes to either the driver or STREAMS kernel code are needed for you to use the ifnet STREAMS module.

To use the ifnet STREAMS module, the STRIFNET and DLPI options must be configured in your kernel and you must set up STREAMS for the driver.

The STRIFNET and DLPI options may have been configured into your system at installation time. (For information on configuring options during installation, see the *Installation Guide*.) You can check to see if the options are configured, by issuing the following command:

```
# /usr/sbin/strsetup -c
```

If ifnet and dlb appear in the Name column, the options are configured in you kernel. If not, you must add them using the doconfig command.

To configure STRIFNET and DLPI into your kernel, perform the following steps:

- 1. Log in as superuser.
- 2. Enter the /usr/sbin/doconfig command. If you have a customized configuration file, you should use the /usr/sbin/doconfig -c command. For more information, see the doconfig(8) reference page.
- 3. Enter a name for the kernel configuration file. It should be the name of your system in uppercase letters, and will probably be the default provided in square brackets ([]); for example:

```
Enter a name for the kernel configuration file. [HOST1]:
RETURN
```

4. Enter y when the system asks whether you want to replace the system configuration file; for example:

```
A configuration file with the name 'HOST1' already exists. Do you want to replace it? (y/n) [n]: \boldsymbol{y}
```

```
Saving /sys/conf/HOST1 as /sys/conf/HOST1.bck

*** KERNEL CONFIGURATION AND BUILD PROCEDURE ***
```

5. Select the options you want to include in you kernel.

#### Note

The STRIFNET and DLPI options are not available from this menu. To include these options, you must edit the configuration file, as shown in the following step.

6. Add DLPI and STRIFNET to the options section of the kernel configuration file.

Enter y when the system asks whether you want to edit the kernel configuration file. The doconfig command allows you to edit the configuration file with the ed editor. For information about using the ed editor, see ed(1).

The following ed editing session shows how to add the DLPI and STRIFNET options to the kernel configuration file for host1. Note that the number of the line after which you append the new lines can differ between kernel configuration files:

```
Do you want to edit the configuration file? (y/n) [n]: y
```

Using ed to edit the configuration file. Press return when ready, or type 'quit' to skip the editing session: 2153

```
48a
options DLPI
options STRIFNET
.
1,$w
2185
q
**** PERFORMING KERNEL BUILD ***
```

- 7. After the new kernel is built, you must move it from the directory where doconfig places it to the root directory ( / ) and reboot your system.
  - When you reboot, the strsetup -i command runs automatically, and creates the device special files for any new STREAMS modules.
- 8. Run the strsetup -c command to verify that the device is configured properly.

The following example shows the output from the strsetup -c command:

#### # /usr/sbin/strsetup -c

STREAMS Configuration Information...Thu Nov 9 08:38:17 1995

Name	Туре	Major	Module ID
clone		32	0
dlb	device	52	5010
dlpi	device	53	800
kinfo	device	54	5020
log	device	55	44
nuls	device	56	5001
echo	device	57	5000
sad	device	58	45
pipe	device	59	5304
xtisoUDP	device	60	5010
xtisoTCP	device	61	5010
xtisoUDP+	device	62	5010
xtisoTCP+	device	63	5010
ptm	device	64	7609
pts	device	6	7608
bba	device	65	24880
lat	device	5	5
pppif	module		6002
pppasync	module		6000
pppcomp	module		6001
bufcall	module		0
ifnet	module		5501
null	module		5002
pass	module		5003
errm	module		5003
ptem	module		5003
spass	module		5007
rspass	module		5008
pipemod	module		5303
timod	module		5006
tirdwr	module		0
ldtty	module		7701

Configured devices = 16, modules = 15

For more detailed information on reconfiguring your kernel or the doconfig command see the *System Administration* manual and the doconfig(8) reference page.

To set up STREAMS for the driver you must do the following:

#### 1. Write an application program similar to the following:

```
^{\star} Application program to set up the "pifnet" streams for IP
* and ARP. This must be run prior to ifconfig
#include <stdio.h>
#include <fcntl.h>
#include <errno.h>
#include <stropts.h>
#include <sys/ioctl.h>
#include <signal.h>
#include "dlpihdr.h"
#define IP_PROTOCOL
                                0x800
#define ARP_PROTOCOL
                               0x806
#define PIFNET_IOCTL_UNIT
                              1236
main(argc, argv)
       int argc;
       char *argv[];
{
       extern char *getenv();
        char *p;
        short unit = 0;
        char devName[256];
        if (argc != 3) usage();
        strcpy(devName, argv[1]);
       unit = atoi(argv[2]);
        sigignore(SIGHUP);
        setupStream(devName, unit, IP_PROTOCOL);
        setupStream(devName, unit, ARP_PROTOCOL);
        * sleep forever to keep the Streams alive.
        if (fork()) /* detach */
           exit();
       pause();
}
usage()
{
       fprintf(stderr, "usage: pifnetd devname unit-number0);
       exit(1);
}
setupStream(devName, unit, serviceClass)
 char *devName;
 short unit;
 u_long serviceClass;
        int fd, status;
        dl_bind_req_t bindreq;
       dl_bind_ack_t bindack;
        int flags;
        struct strioctl str;
        struct strbuf pstrbufctl, pstrbufdata, gstrbufctl, \
                                               gstrbufdata;
```

```
char ebuf[256];
 * build the stream
fd = open(devName, O_RDWR, 0);
if (fd < 0)
         sprintf(ebuf, " open '%s' failed", devName);
         perror(ebuf);
         exit(1);
if (ioctl(fd, I_PUSH, "ifnet") < 0)</pre>
         sprintf(ebuf, " ioctl I_PUSH failed");
         perror(ebuf);
         exit(1);
}
 \mbox{\scriptsize \star} tell pifnet the unit number for the device
str.ic_cmd = PIFNET_IOCTL_UNIT;
str.ic_timout = 15;
str.ic_len = sizeof (short);
str.ic_dp = (char *) &unit;
status = ioctl(fd, I_STR, &str);
if (status < 0)
         sprintf(ebuf, " %s - ioctl");
         perror(ebuf);
         exit(1);
 * bind the stream to a protocol
bindreq.dl_primitive = DL_BIND_REQ;
bindreq.dl_sap = serviceClass;
bindreq.dl_max_conind = 0;
bindreq.dl_service_mode = DL_CLDLS;
bindreq.dl_conn_mgmt = 0;
bindreq.dl_xidtest_flg = 0;
pstrbufctl.len = sizeof(dl_bind_req_t);
pstrbufctl.buf = (void *)&bindreq;
pstrbufdata.buf = (char *)0;
pstrbufdata.len = -1;
pstrbufdata.maxlen = 0;
status = putmsg(fd, &pstrbufctl, (struct strbuf *)0, 0);
if (status < 0)
{
         perror("putmsg");
         exit(1);
}
 * Check requested binding
```

```
gstrbufctl.buf = (char *)&bindack;
gstrbufctl.maxlen = sizeof(dl_bind_ack_t);
gstrbufctl.len = 0;
status = getmsg(fd, &gstrbufctl, (struct strbuf *)0, &flags);
if (status < 0)
{
         perror("getmsg");
         exit(1);
}

if (bindack.dl_primitive != DL_BIND_ACK)
{
         errno = EPROTO;
         perror(" DL_BIND_ACK");
         exit(1);
}</pre>
```

In this sample application the driver's name is /dev/streams/ln. The application creates two Streams; one for the Internet Protocol (IP) and one for the Address Resolution Protocol (ARP). After setting up the Streams, the application must keep running, using the pause command, in order to keep the Streams alive.

Note that, if the driver is a style-2 driver, you must add a DL\_ATTACH\_REQ to the application program. For more information about the DL\_ATTACH\_REQ primitive or style-2 drivers, see the DLPI specification in /usr/share/doclib/dlpi/dlpi.ps.

- 2. Generate an executable file for the application. Compile, link, and debug the program until it runs without errors.
- 3. Move the executable into a directory that is convenient for you.

The executable can be located in any directory.

4. Add a line invoking the program to the /sbin/init.d/inet file.

Although you can manually start the program each time you reboot, it is easiest to add a line to the /sbin/init.d/inet file to run it automatically when the system reboots. Be sure to add the line before the system's ifconfig lines.

In the following example, each time the system reboots, the /sbin/init.d/inet file runs a program called run\_ifnet, which resides in the /etc directory:

```
# # Enable network # case $1 in echo "Configuring network" /sbin/hostname $HOSTNAME echo "hostname: \c"
```

```
/sbin/hostname
       if [ "$NETDEV_0" != '' ]; then
               echo >/tmp/ifconfig_"$NETDEV_0".tmp
# place command invoking executable BEFORE \fP
      ifconfig lines
               /etc/run_ifnet
               /sbin/ifconfig $NETDEV_0 $IFCONFIG_0 > \
                              /tmp/ifconfig_"$NETDEV_0".tmp 2>&1
              if [ $? != 0 ]; then
                       ERROR='cat /tmp/ifconfig_"$NETDEV_0".tmp'
                       if [ "$ERROR" = "$ERRSTRING" ]; then
                               /sbin/ifconfig $NETDEV_0 up
                       else
                               echo "$0: $ERROR"
               fi
               rm /tmp/ifconfig_"$NETDEV_0".tmp
       fi
```

#### 5. Reboot the system.

Use the /usr/sbin/shutdown -r command to shut down your system and have it reboot automatically; for example:

# /usr/sbin/shutdown -r now

#### 7.1.1.2 Data Link Provider Interface Primitives

Note that the STREAMS device driver can be a style-1 or a style-2 DLPI provider, as described in the Data Link Provider Interface specification, which is located in /usr/share/doclib/dlpi/dlpi.ps. Note that you must have the OSFPGMR400 subset installed to access the DLPI specification on line.

The driver must support the following DLPI primitives. For detailed information about these primitives and how to use them, see the DLPI specification:

```
DL_PHYS_ADDR_REQ/DL_PHYS_ADDR_ACK

DL_BIND_REQ/DL_BIND_ACK

DL_UNBIND_REQ

DL_UNITDATA_REQ/DL_UNITDATA_IND/DL_UDERROR_IND

DL_OK_ACK/DL_ERROR_ACK
```

## 7.2 Bridging BSD Drivers to STREAMS Protocol Stacks

The dlb STREAMS pseudodevice driver allows you to bridge BSD-style device drivers and STREAMS protocol stacks. The STREAMS pseudodevice driver is the Stream end in a Stream wanting to communicate with BSD-based drivers. The STREAMS pseudodevice driver provided by Digital UNIX has two interfaces, a subset of the DLPI interface that communicates with STREAMS protocol stacks and another interface that accesses the ifnet layer interface of the sockets framework.

Figure 7-2 highlights the dlb STREAMS pseudodriver and shows its place in the network programming environment.

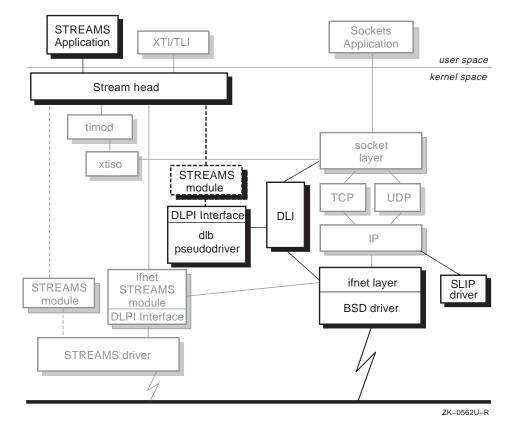


Figure 7-2: DLPI STREAMS Pseudodriver

#### 7.2.1 Supported DLPI Primitives and Media Types

The dlb STREAMS pseudodriver supports the following connectionless mode primitive and media types. For detailed information about these primitives and how to use them, see the Data Link Provider Interface specification which is in /usr/share/doclib/dlpi/dlpi.ps.

```
DL_ATTACH_REQ/DL_DETACH_REQ/DL_OK_ACK

DL_BIND_REQ/DL_BIND_ACK/DL_UNBIND_REQ

DL_ENABMULTI_REQ/DL_DISABLMULTI_REQ

DL_PROMISCON_REQ/DL_PROMISCONOFF_REQ

DL_PHYS_ADDR_REQ/DL_PHYS_ADDR_ACK

DL_SET_PHYS_ADDR_REQ

DL_UNITDATA_REQ/DL_UNITDATA_IND

DL_SUBS_BIND_REQ/DL_SUBS_BIND_ACK

DL_SUBS_UNBIND_REQ/DL_SUBS_UNBIND_ACK

The Ethernet bus (DL_ETHER) is the media type supported by the
```

### 7.2.2 Using the STREAMS Pseudodriver

To use the dlb STREAMS pseudodriver the DLPI option must be configured into your kernel. The DLPI option may have been configured into your system at installation time.

You can check to see if the DLPI option is configured by issuing the following command:

# /usr/sbin/strsetup -c

STREAMS pseudodriver.

If dlb appears in the Name column, the option is configured in you kernel. If not, you must add it using the doconfig command.

For a description of how to reconfigure your kernel with the doconfig command, see Section 7.1.1.1.

For more information on reconfiguring your kernel or the doconfig command see the *System Administration* manual and the doconfig(8) reference page. For information on configuring options during installation, see the *Installation Guide*.

# Sample STREAMS Module A

The spass module is a simple STREAMS module that passes all messages put to it to the putnext() procedure. The spass module delays the call to putnext() for the service procedure to handle. It has flow control code built in, and both the read and write sides share a service procedure.

The following is the code for the spass module:

```
#include <sys/stream.h>
#include <sys/stropts.h>
#include <sys/sysconfig.h>
               spass_close();
static int
static int spass_open();
static int spass_rput();
static int spass_srv();
static int spass_wput();
static struct module_info minfo = {
       0, "spass", 0, INFPSZ, 2048, 128
};
static struct qinit rinit = {
    spass_rput, spass_srv, spass_open, spass_close, NULL, &minfo
static struct qinit winit = {
    spass_wput, spass_srv, NULL, NULL, NULL, &minfo
struct streamtab spassinfo = { &rinit, &winit };
cfg_subsys_attr_t bufcall_attributes[] = {
   {"", 0, 0, 0, 0, 0, 0} /* must be the last element */
spass_configure(op, indata, indata_size, outdata, outdata_size)
        ulong
                        indata_size;
        caddr_t
                        outdata;
        ulong
                        outdata_size;
{
     struct streamadm
                              sa;
     dev_t
                              devno = NODEV;
     if (op != CFG_OP_CONFIGURE)
             return EINVAL;
```

```
sa.sa_version
                            = OSF_STREAMS_10;
     sa.sa_flags
                              = STR_IS_MODULE | STR_SYSV4_OPEN;
                             = 0;
     sa.sa_ttys
     sa.sa_sync_level
                            = SQLVL_QUEUE;
                              = 0;
     sa.sa_sync_info
     strcpy(sa.sa_name,
                             "spass");
     if ( (devno = strmod_add(devno, &spassinfo, &sa)) == NODEV ) {
             return ENODEV;
     return 0;
}
/* Called when module is popped or the Stream is closed */
static int
spass_close (q, credp)
    queue_t * q;
cred_t * credp;
     return 0;
/* Called when module is pushed */
static int
spass_open (q, devp, flag, sflag, credp)
    queue_t * q;
int * devp;
     int
            flag;
     int
            sflag;
     cred_t * credp;
     return 0;
* Called to process a message coming upstream. All messages
* but flow control messages are put on the read side service
* queue for later processing.
static int
spass_rput (q, mp)
    queue_t * q;
    mblk_t * mp;
{
     switch (mp->b_datap->db_type) {
     case M_FLUSH:
             if (*mp->b_rptr & FLUSHR)
                    flushq(q, 0);
             putnext(q, mp);
             break;
     default:
             putq(q, mp);
             break;
     return 0;
}
* Shared by both read and write sides to process messages put
* on the read or write service queues. When called from the
```

```
\mbox{\scriptsize *} write side, sends all messages on the write side queue
* downstream until flow control kicks in or all messages are
^{\star} processed. When called from the read side sends all messages
* on its read side service queue upstreams until flow control
* kicks in or all messages are processed.
static int
spass_srv (q)
     queue_t * q;
     mblk_t *
     while (mp = getq(q)) {
              if (!canput(q->q_next))
                   return putbq(q, mp);
              putnext(q, mp);
     return 0;
}
\mbox{\scriptsize \star} Called to process a message coming downstream. All messages but
\mbox{* flow control} messages are put on the write side service queue for
* later processing.
static int
spass_wput (q, mp)
     queue_t * q;
mblk_t *
                      mp;
{
     switch (mp->b_datap->db_type) {
     case M_FLUSH:
              if (*mp->b_rptr & FLUSHW)
                     flushq(q, 0);
              putnext(q, mp);
              break;
     default:
              putq(q, mp);
              break;
     return 0;
```

# Socket and XTI Programming Examples B

This appendix contains annotated files for a sample server/client<sup>1</sup> credit card authorization program. Clients access a server on the merchant's behalf and request authorization from the server to put a charge on the client's credit card. The server maintains a database of authorized merchants and their passwords, as well as a database of credit card customers, their credit limit, and current balance. It either authorizes or rejects a client request based on the information in its database.

Several variations on the credit card authorization program are presented, including connection-oriented and connectionless modes. The connection-oriented and connectionless modes each contain socket and XTI code for the server and client portions of the program.

Although the program uses network programming in a real world application, it has the following limitations:

- Error handling is not robust
- · Accepts only integer amounts
- Performs no child process clean up
- In the case of the connection-oriented protocol examples in Section B.1, for each request received, the server program forks a child process to handle the request. The database information is "detached" in the child process' private data area. When the child process analyzes the request and reduces the customer's credit balance appropriately, it needs to update this information in the original server's data area (and on some persistent storage as well) so that the next request for the same customer is handled correctly. To avoid unnecessary complexity, this logic is not included in the program.

The information is organized as follows:

- Connection-oriented mode programs
  - Socket

 $<sup>^{1}</sup>$  The term client in this appendix refers to the program initiated by the merchant which interacts with the server program.

- \* Server
- \* Client
- XTI
  - \* Server
  - \* Client
- Connectionless mode programs
  - Socket
    - \* Server
    - \* Client
  - XTI
    - \* Server
    - \* Client
- Common files

You can obtain copies of these example programs from /usr/examples/network\_programming.

## **B.1 Connection-Oriented Programs**

This section contains sockets and XTI variations of the same server and client programs, written for connection-oriented modes communication.

### **B.1.1 Socket Server Program**

Example B-1 implements a server using the socket interface.

#### **Example B-1: Connection-Oriented Socket Server Program**

```
Example B-1: (continued)
                              sockfd;
        int
        int
                              newsockfd;
        struct sockaddr_in
                             serveraddr;
        struct sockaddr_in
                              clientaddr;
                              clientaddrlen = sizeof(clientaddr);
        int
        struct hostent
                              *he;
        int
                              pid;
        signal(SIGCHLD, SIG_IGN);
        if ((sockfd = socket(AF_INET, SOCK_STREAM, 0)) < 0)</pre>
                perror("socket_create");
                exit(1);
        }
        bzero((char *) &serveraddr,
              sizeof(struct sockaddr_in));
                                                                2
        serveraddr.sin_family = AF_INET;
        serveraddr.sin_addr.s_addr = htonl(INADDR_ANY);
        serveraddr.sin_port = htons(SERVER_PORT);
        if ( bind(sockfd,
                  (struct sockaddr *)&serveraddr,
                  sizeof(struct sockaddr_in)) < 0) {</pre>
                perror("socket_bind");
                exit(2);
        }
        listen(sockfd, 8);
                                                                6
        while(1) {
                if ((newsockfd =
                      accept(sockfd,
                                                                7
                             (struct sockaddr *) &clientaddr,
                             &clientaddrlen)) < 0) {
                        if (errno == EINTR) {
                                printf("Bye...\n");
                                exit(0);
                        } else {
                                perror("socket_accept");
                                exit(3);
                        }
                }
                pid = fork();
                switch(pid) {
```

#### Example B-1: (continued)

```
case -1:
                                        /* error */
                                  perror("dosession_fork");
                                  break;
                         default:
                                  close(newsockfd);
                                  break;
                         case 0:
                                          /* child */
                                  close(sockfd);
                                  transactions(newsockfd);
                                  close(newsockfd);
                                  return(0);
        }
}
transactions(int fd)
                                                                   8
        int
                bytes;
                 *reply;
        char
        int
                 dcount;
        char
                datapipe[MAXBUFSIZE+1];
         \mbox{\ensuremath{\star}} Look at the data buffer and parse commands,
         * keep track of the collected data through
         * transaction_status.
         * /
         while (1) {
                 if ((dcount=recv(fd, datapipe, MAXBUFSIZE))
                     < 0) {
                         perror("transactions_receive");
                         break;
                 if (dcount == 0) {
                         return(0);
                 datapipe[dcount] = '\0';
                 if ((reply=parse(datapipe)) != NULL) {
                         send(fd, reply, strlen(reply), 0);
                                                                   10
         }
}
```

- 1 Create a socket with the socket call.
  - AF\_INET specifies the Internet communication domain. Alternatively, if OSI transport were supported, a corresponding constant such as AF\_OSI would be required here. The socket type SOCK\_STREAM is specified for TCP or connection-oriented communication. This parameter indicates that the socket is connection-oriented.
  - Contrast the socket call with the t\_open call in the XTI server example (Section B.1.3).
- 2 The serveraddr is of type sockaddr\_in, which is dictated by the communication domain of the socket (AF\_INET). The socket address for the Internet communication domain contains an Internet address and a 16-bit port number, which uniquely identifies an entity on the network. For the TCP/IP and UDP/IP this is the Internet address of the server and the port number on which it is listening.
  - Note that the information contained in the sockaddr\_in structure is dependent on the address family, which is AF\_INET in this example. If AF\_OSI were used instead of AF\_INET, then sockaddr\_osi would be required for the bind call instead of sockaddr\_in.
- 3 INADDRANY signifies any attached interface adapter on the system. All numbers must be converted to the network format using appropriate macros. See the following reference pages for more information: hton1(3), htons(3), ntoh1(3), and ntohs(3).
- **4** SERVER\_PORT is defined in the common.h header file. It is a short integer, which helps identify the server process from other application processes. Numbers from 0 to 1024 are reserved.
- **5** Bind the server's address to this socket with the bind call. The combination of the address and port number identify it uniquely on the network.
- 6 Specify the number of pending connections the server can queue while it finishes processing the previous accept call. This value governs the success rate of connections while the server processes accept calls. Use a larger number to obtain a better success rate if multiple clients are sending the server connect requests. The operating system imposes a ceiling on this value.
- **7** Accept connections on this socket. For each connection, the server forks a child process to handle the session to completion. The server then resumes listening for new connection requests. This is an example of a concurrent server. You can also have an iterative server, meaning that the server handles the data itself. See Section B.2 for an example of iterative servers.

8 Each incoming message packet is accepted and passed to the parse function, which tracks the information provided, such as the merchant's login ID, password, and customer's credit card number. This process is repeated until the parse function identifies a complete transaction and returns a response packet, to be sent to the client program.

The client program can send information packets in any order (and in one or more packets), so the parse function is designed to remember state information sufficient to deal with this unstructured message stream.

Since the program uses a connection-oriented protocol for data transfer, this function uses send and receive messages, respectively.

- **9** Receive data with the recv call.
- 10 Send data with the send call.

#### **B.1.2 Socket Client Program**

Example B-2 implements a client program that can communicate with the socketserver interface shown in Example B-1.

#### **Example B-2: Connection-Oriented Socket Client Program**

```
* This generates the client program.
 * usage: socketclient [serverhostname]
 * If a host name is not specified, the local
 * host is assumed.
#include "client.h"
main(int argc, char *argv[])
{
                                 sockfd;
        int
        struct sockaddr_in serveraddr;
struct hostent *he;
        int
                                 n;
                                 *serverhost = "localhost";
*serverhostp;
        char
        struct hostent
                               *servel....buffer[1024];
        char
        char
                                  inbuf[1024];
        if (argc>1) {
```

# **Example B-2: (continued)**

```
serverhost = argv[1];
}
init();
if ((sockfd = socket(AF_INET, SOCK_STREAM, 0)) < 0)</pre>
        perror("socket_create");
        exit(1);
}
bzero((char *) &serveraddr,
      sizeof(struct sockaddr_in));
                                                        2
serveraddr.sin_family
                          = AF_INET;
if ((serverhostp = gethostbyname(serverhost)) ==
    (struct hostent *)NULL) {
        fprintf(stderr, "gethostbyname on %s failed\n",
                serverhost);
        exit(1);
bcopy(serverhostp->h_addr,
      (char *)&(serveraddr.sin_addr.s_addr),
                  serverhostp->h_length);
serveraddr.sin_port
                           = htons(SERVER_PORT);
/* Now connect to the server */
if (connect(sockfd, &serveraddr, sizeof(serveraddr))
    < 0) {
        perror ("connect");
        exit(2);
}
while(1) {
        /* Merchant record */
        sprintf(buffer, "%%%%m%s##%%%%p%s##",
                merchantname, password);
        printf("\n\nSwipe card, enter amount: ");
        fflush(stdout);
        if (scanf("%s", inbuf) == EOF) {
                printf("bye...\n");
                exit(0);
        }
        soundbytes();
        sprintf(buffer, "%s%%%a%s##%%%%n%s##",
                buffer, inbuf, swipecard());
        if (send(sockfd, buffer, strlen(buffer), 0)
```

# Example B-2: (continued)

1 Create a socket with the socket call.

AF\_INET is the socket type for the Internet communication domain. Note that this parameter must match the protocol and type selected in the corresponding server program.

Contrast the socket call with the t\_open call in the XTI client example (Section B.1.4).

2 The serveraddr is of type sockaddr\_in, which is dictated by the communication domain of the socket (AF\_INET). The socket address for the Internet communication domain contains an Internet address and a 16-bit port number, which uniquely identifies an entity on the network. For the TCP/IP protocol suite, this is the Internet address of the server and the port number on which it is listening.

Note that the information contained in the sockaddr\_in structure is dependent on the address family (or the protocol).

- **3** Getting information about the server depends on the protocol or the address family. To get the IP address of the server, you can use the gethostbyname routine.
- 4 SERVER\_PORT is defined in the common.h header file. It is imperative that the same port number be used to connect to the socket server program. The server and client select the port number, which functions as a well known address for communication.

- 5 Client issues a connect call to connect to the server. When the connect call is used with a connection-oriented protocol, it allows the client to build a connection with the server before sending data. This is analogous to dialing a phone number.
- 6 Send data with the send call.
- 7 Receive data with the recy call.

# **B.1.3 XTI Server Program**

Example B-3 implements a server using the XTI library for network communication. It is an alternative design for a communication program that makes it transport independent. Compare this program with the socket server program in Section B.1.1. This program has the same limitations described at the beginning of the appendix.

# **Example B-3: Connection-Oriented XTI Server Program**

```
* This file contains the main XTI server code
* for a connection-oriented mode of communication.
 * Usage: xtiserver
#include "server.h"
                         *parse(char *);
char
struct transaction
                         *verifycustomer(char *, int, char *);
main(int argc, char *argv[])
                                 xtifd;
       struct sockaddr_in serveraddr;
struct hostent *he;
int
        int
        int pid;
struct t_bind *bindreqp;
struct t_bind *bindretp;
        signal(SIGCHLD, SIG_IGN);
        if ((xtifd = t_open("/dev/streams/xtiso/tcp", O_RDWR, 1
                            NULL)) < 0) {
                 xerror("xti_open", xtifd);
                 exit(1);
```

# Example B-3: (continued)

bzero((char \*) &serveraddr, sizeof(struct sockaddr\_in)); serveraddr.sin\_family = AF\_INET; serveraddr.sin\_addr.s\_addr = htonl(INADDR\_ANY); 3 serveraddr.sin\_port = htons(SERVER\_PORT); /\* allocate structures for the t\_bind call \*/ if (((bindreqp=(struct t\_bind \*) t\_alloc(xtifd, T\_BIND\_STR, T\_ALL)) == NULL) || ((bindretp=(struct t\_bind \*) t\_alloc(xtifd, T\_BIND\_STR, T\_ALL)) == NULL)) { xerror("xti\_alloc", xtifd); exit(3); } = (char \*)&serveraddr; = sizcof' bindreqp->addr.buf bindreqp->addr.len = sizeof(serveraddr); \* Specify how many pending connections can be \* maintained, until finish accept processing \* / bindreqp->qlen = 8; if (t\_bind(xtifd, bindreqp, (struct t\_bind \*)NULL) < 0) { xerror("xti\_bind", xtifd); exit(4);} \* Now the socket is ready to accept connections. \* For each connection, fork a child process in charge \* of the session, and then resume accepting connections. \* / while(1) { struct t\_call call; if (t\_listen(xtifd, &call) < 0) {</pre> 7 if (errno == EINTR) { printf("Bye...\n"); exit(0); } else {

xerror("t\_listen", xtifd);

# **Example B-3: (continued)**

```
}
                 * Create a new transport endpoint on which
                 * to accept a connection
                 * /
                if ((newxtifd=t_open("/dev/streams/xtiso/tcp", 8
                                      O_RDWR, NULL)) < 0) {
                        xerror("xti_newopen", xtifd);
                        exit(5);
                if (t_bind(newxtifd,
                                                                 9
                           (struct t_bind *)NULL,
                            (struct t_bind *)NULL) < 0) {
                        xerror("xti_newbind", xtifd);
                        exit(6);
                /* accept connection */
                if (t_accept(xtifd, newxtifd, &call) < 0) {</pre>
                                                                 10
                        xerror("xti_accept", xtifd);
                        exit(7);
                }
                pid = fork();
                switch(pid) {
                        case -1:
                                         /* error */
                                 xerror("dosession_fork", xtifd);
                                 break;
                        default:
                                 t_close(newxtifd);
                                 break;
                        case 0:
                                         /* child */
                                 t_close(xtifd);
                                 transactions(newxtifd);
                                 t_close(newxtifd);
                                 return(0);
               }
        }
}
transactions(int fd)
                                                                 11
{
                bytes;
        int
```

exit(4);

### **Example B-3: (continued)**

```
char
        *reply;
        dcount;
int
int
        flags;
        datapipe[MAXBUFSIZE+1];
char
 ^{\star} Look at the data buffer and parse commands, if more data
 * required go get it
 * Since the protocol is SOCK_STREAM oriented, no data
 * boundaries will be preserved.
 * /
while (1) {
        if ((dcount=t_rcv(fd, datapipe, MAXBUFSIZE,
                                                         12
                           &flags)) < 0){
                /* if disconnected bid a goodbye */
                if (t_errno == TLOOK) {
                         int tmp = t_look(fd);
                         if (tmp != T_DISCONNECT) {
                                 t_scope(tmp);
                         } else {
                                 exit(0);
                xerror("transactions_receive", fd);
                break;
        if (dcount == 0) {
                /* consolidate all transactions */
                return(0);
        datapipe[dcount] = ' ';
        if ((reply=parse(datapipe)) != NULL) {
                if (t_snd(fd, reply, strlen(reply), 0) 13
                        xerror("xti_send", fd);
                        break;
                }
 }
```

1 The t\_open call specifies a device special file name; for example /dev/streams/xtiso/tcp. This file name provides the necessary abstraction for the TCP transport protocol over IP. Unlike the socket interface, where you specify the address family (for example, AF\_INET), this information is already represented in the choice of the device special file. The /dev/streams/xtiso/tcp file implies both TCP transport and IP. See the Chapter 5 for information about STREAMS devices.

}

As mentioned in Section B.1.1, if the OSI transport were available you would use a device such as /dev/streams/xtiso/cots.

Contrast the t\_open call with the socket call in Section B.1.1.

- 2 Selection of the address depends on the choice of the transport protocol. Note that in the socket example the address family was the same as used in the socket system call. With XTI, the choice is not obvious and you must know the appropriate mapping from the transport protocol to sockaddr. See Chapter 3 for more information.
- 3 INADDRANY signifies any attached interface adapter on the system. All numbers must be converted to the network format using appropriate macros. See the following reference pages for more information: hton1(3), htons(3), ntoh1(3), ntohs(3).
- **4** SERVER\_PORT is defined in the common.h header file. It has a data type of short integer which helps identify the server process from other application processes. Numbers from 0 to 1024 are reserved.
- **5** Specify the number of pending connections the server can queue while it processes the last request.
- Bind the server's address with the t\_bind call. The combination of the address and port number uniquely identify it on the network. After the server process' address is bound, the server process is registered on the system and can be identified by the lower level kernel functions to which to direct any requests.
- 7 Listen for connection requests with the t\_listen function.
- **8** Create a new transport endpoint with another call to the t\_open function.
  - Bind the server's address with the t\_bind call. The combination of the address and port number identify it uniquely on the network.
- **9** Bind to the new transport endpoint with the t bind function.
- **10** Accept the connection request with the t\_accept function.
- 11 Each incoming message packet is accepted and passed to the parse function, which tracks the information provided (such as the merchant's login ID, password, and customer's credit card number). This process is repeated until the parse function identifies a complete transaction and returns a response packet, to be sent to the client program.

The client program can send information packets in any order (and in one or more packets), so the parse function is designed to remember state information sufficient to deal with this unstructured message stream.

Since the program uses a connection-oriented protocol for data transfer, this function uses t\_snd and t\_rcv to send and receive data, respectively.

- **12** Receive data with the t\_rcv function.
- 13 Send data with the t\_snd function.

# **B.1.4 XTI Client Program**

Example B-4 This sample program implements a client program that can communicate with the xtiserver interface shown in Section B.1.3. Compare this program with the socket client program in Example B-3.

### **Example B-4: Connection-Oriented XTI Client Program**

```
* This file contains the main XTI client code
* for a connection-oriented mode of communication.
* Usage: xticlient [serverhostname]
* If a host name is not specified, the local
* host is assumed.
#include "client.h"
main(int argc, char *argv[])
      inbuf[1024];
sndcall;
       char
       struct t_call sndcall;
struct t_call rcvcall;
                           flags = 0;
       if (argc>1) {
             serverhost = argv[1];
       init();
       if ((xtifd = t_open("/dev/streams/xtiso/tcp", O_RDWR, 1
                       NULL)) < 0) {
              xerror("xti_open", xtifd);
              exit(1);
```

# Example B-4: (continued)

```
bzero((char *) &serveraddr,
                                                         2
      sizeof(struct sockaddr_in));
serveraddr.sin_family = AF_INET;
                                                         3
if ((serverhostp = gethostbyname(serverhost)) ==
    (struct hostent *)NULL) {
        fprintf(stderr, "gethostbyname on %s failed\n",
                serverhost);
        exit(1);
bcopy(serverhostp->h_addr,
      (char *)&(serveraddr.sin_addr.s_addr),
                  serverhostp->h_length);
serveraddr.sin_port
                        = htons(SERVER_PORT);
if (t_bind(xtifd, (struct t_bind *)NULL,
           (struct t_bind *)NULL) < 0) {
        xerror("bind", xtifd);
        exit(2);
}
sndcall.opt.maxlen
                        = 0;
sndcall.udata.maxlen = 0;
sndcall.uuacu.....sndcall.addr.buf = (char ^)&selve___
= sizeof(serveraddr);
                        = (char *)&serveraddr;
rcvcall.opt.maxlen = 0;
rcvcall.udata.maxlen = 0;
rcvcall.addr.buf = (char *)&clientaddr;
rcvcall.addr.maxlen = sizeof(clientaddr);
                                                         7
if (t_connect(xtifd, &sndcall,
              (struct t_call *)NULL) < 0) {
        xerror ("t_connect", xtifd);
        exit(3);
}
while(1) {
        /* Merchant record */
        sprintf(buffer, "%%%%m%s##%%%%p%s##",
                 merchantname, password);
        printf("\n\nSwipe card, enter amount: ");
        fflush(stdout);
        if (scanf("%s", inbuf) == EOF) {
                printf("bye...\n");
                 exit(0);
        soundbytes();
```

# Example B-4: (continued)

```
sprintf(buffer, "%s%%%a%s##%%%%n%s##",
                         buffer, inbuf, swipecard());
                if (t_snd(xtifd, buffer, strlen(buffer), 0)
                     < 0) {
                         xerror("t_snd", xtifd);
                         exit(1);
                 }
                if ((n = t_rcv(xtifd, buffer, 1024, &flags)) 9
                     < 0) {
                         xerror("t_rcv", xtifd);
                         exit(1);
                buffer[n] = ' \setminus 0';
                if ((n=analyze(buffer))== 0) {
                         printf("transaction failure,"
                                " try again\n");
                 } else if (n<0) {</pre>
                         printf("login failed, try again\n");
                         init();
                }
        }
}
```

1 AF\_INET is the socket type for the Internet communication domain. If AF\_OSI were supported, it could be used to create a socket for OSI communications. The socket type SOCK\_STREAM is specified for TCP or connection-oriented communication.

The t\_open call specifies a special device file name instead of the socket address family, socket type, and protocol that the socket call requires.

Contrast the socket call in Section B.1.2 with the t\_open call.

- 2 The serveraddr is of type sockaddr\_in, which is dictated by the communication domain of the socket (AF\_INET). The socket address for the Internet communication domain contains an Internet address and a 16-bit port number, which uniquely identifies an entity on the network. For the TCP/IP protocol suite, which includes UDP, this is the Internet address of the server and the port number on which it is listening.
  - Note that the information contained in the sockaddr\_in structure is dependent on the address family (or the protocol).
- **3** AF\_INET specifies the Internet communication domain. If AF\_OSI were supported, it could be used to create a socket for OSI communications.

- 4 Obtaining information about the server depends on the protocol or the address family. To get the IP address of the server, you can use the gethostbyname routine.
- **5** SERVER\_PORT is defined in the <common.h> header file. It is imperative that the same port number be used to connect to the XTI server program. Numbers from 0 through 1024 are reserved.
- **6** Bind the server address with the t\_bind function to enable the client to start sending and receiving data.
- 7 Initiate a connection with the server using the t\_connect function.
- 8 Send data with the t\_snd function.
- **9** Receive data with the t\_rcv function.

# **B.2 Connectionless Programs**

This section contains sockets and XTI variations of the same server and client programs, written for connectionless modes of communication.

# **B.2.1 Socket Server Program**

Example B-5 implements the server portion of the application in a manner similar to the socket server described in Section B.1.1. Instead of using a connection-oriented paradigm, this program uses a connectionless (datagram/UDP) paradigm for communicating with client programs. This program has the limitations described at the beginning of the appendix.

# **Example B-5: Connectionless Socket Server Program**

### Example B-5: (continued)

```
struct sockaddr_in
                                clientaddr;
                                clientaddrlen = sizeof(clientaddr);
        int
                                *he;
       struct hostent
                                pid;
        int
       signal(SIGCHLD, SIG_IGN);
        /* Create a socket for the communications */
        if ((sockfd = socket(AF_INET, SOCK_DGRAM, 0))
            < 0) {
                perror("socket_create");
                exit(1);
       }
       bzero((char *) &serveraddr,
                                                             2
             sizeof(struct sockaddr_in));
       serveraddr.sin_family = AF_INET;
       serveraddr.sin_addr.s_addr = htonl(INADDR_ANY);
       serveraddr.sin_port = htons(SERVER_PORT);
                                                             5
        if (bind(sockfd,
                 (struct sockaddr *)&serveraddr,
                 sizeof(struct sockaddr_in)) < 0) {</pre>
                perror("socket_bind");
                exit(2);
       }
       transactions(sockfd);
}
transactions(int fd)
                             bytes;
       int
       char
                              *reply;
                              dcount;
        int
                             datapipe[MAXBUFSIZE+1];
        struct sockaddr_in
                             serveraddr;
                              serveraddrlen = sizeof(serveraddr);
       int
       bzero((char *) &serveraddr, sizeof(struct sockaddr_in));
       serveraddr.sin_family = AF_INET;
       serveraddr.sin_addr.s_addr = htonl(INADDR_ANY);
       serveraddr.sin_port
                               = htons(SERVER_PORT);
        * Look at the data buffer and parse commands.
         * Keep track of the collected data through
```

# **Example B-5: (continued)**

```
* transaction_status.
         * /
         while (1) {
                                                                7
                if ((dcount=recvfrom(fd, datapipe,
                                MAXBUFSIZE, 0,
                                 (struct sockaddr *)&serveraddr,
                                 &serveraddrlen)) < 0){
                         perror("transactions_receive");
                         break;
                if (dcount == 0) {
                         return(0);
                datapipe[dcount] = ' \setminus 0';
                if ((reply=parse(datapipe)) != NULL) {
                         if (sendto(fd, reply, strlen(reply), 8
                             (struct sockaddr *)&serveraddr,
                             serveraddrlen) < 0) {
                                 perror("transactions_sendto");
                         }
         }
}
```

1 Create a socket with the socket call.

AF\_INET specifies the Internet communication domain. The socket type SOCK\_DGRAM is specified for UDP or connectionless communication. This parameter indicates that the program is connectionless.

Contrast the socket call with the t\_open call in the XTI server example (Section B.2.3).

- 2 The serveraddr is of type sockaddr\_in, which is dictated by the communication domain of the socket (AF\_INET). The socket address for the Internet communication domain contains an Internet address and a 16-bit port number, which uniquely identifies an entity on the network. For the TCP/IP protocol suite, which includes UDP, this is the Internet address of the server and the port number on which it is listening.
  - The information contained in the sockaddr\_in structure is dependent on the address family, which is AF\_INET in this example. If AF\_OSI were used instead of AF\_INET, then sockaddr\_osi would be required for the bind call instead of sockaddr\_in.
- **3** INADDRANY signifies any attached interface adapter on the system. All numbers must be converted to the network format using appropriate

- macros. See the following reference pages for more information: hton1(3), htons(3), ntoh1(3), and ntohs(3).
- **4** SERVER\_PORT is defined in the <common.h> header file. It has a data type of short integer which helps identify the server process from other application processes.
- **5** Bind the server's address to this socket with the bind call. The combination of the address and port number identify it uniquely on the network.
  - After the server process' address is bound, the server process is registered on the system and can be identified by the lower level kernel functions to which to direct requests.
- 6 Each incoming message packet is accepted and passed to the parse function, which tracks the information provided (such as the merchant's login ID, password, and customer's credit card number). This process is repeated until the parse function identifies a complete transaction and returns a response packet, to be sent to the client program.
  - Since this program uses a connectionless (datagram) protocol, it uses sendto and recyfrom to send and receive data, respectively.
- 7 Receive data with the recyfrom call.
- 8 Send data with the sendto call.

# **B.2.2 Socket Client Program**

Example B-6 implements a socket client that can communicate with the socket server in Example B-5. Section B.2.1. It uses the socket interface in the connectionless, or datagram, mode.

### **Example B-6: Connectionless Socket Client Program**

# **Example B-6: (continued)**

```
struct sockaddr_in
                      serveraddr;
int.
                       serveraddrlen;
struct hostent
                       *he;
                      n;
int
                       *serverhost = "localhost";
char
struct hostent
                      *serverhostp;
                      buffer[1024];
char
                       inbuf[1024];
char
if (argc>1) {
       serverhost = argv[1];
init();
/* Create a socket for the communications */
if ((sockfd = socket(AF_INET, SOCK_DGRAM, 0)) < 0)</pre>
       perror("socket_create");
        exit(1);
}
bzero((char *) &serveraddr,
                                                      2
     sizeof(struct sockaddr_in));
serveraddr.sin_family = AF_INET;
if ((serverhostp = gethostbyname(serverhost)) ==
    (struct hostent *)NULL) {
        fprintf(stderr, "gethostbyname on %s failed\n",
               serverhost);
        exit(1);
bcopy(serverhostp->h_addr,
      (char *)&(serveraddr.sin_addr.s_addr),
      serverhostp->h_length);
serveraddr.sin_port
                         = htons(SERVER_PORT);
/* Now connect to the server
                                                      5
if (connect(sockfd, &serveraddr,
           sizeof(serveraddr)) < 0) {</pre>
       perror ("connect");
       exit(2);
while(1) {
        /* Merchant record */
        sprintf(buffer, "%%%%m%s##%%%p%s##",
```

### Example B-6: (continued)

```
printf("\n\nSwipe card, enter amount: ");
        fflush(stdout);
        if (scanf("%s", inbuf) == EOF) {
                printf("bye...\n");
                exit(0);
        }
        soundbytes();
        sprintf(buffer, "%s%%%%a%s##%%%%n%s##",
                buffer, inbuf, swipecard());
        if (sendto(sockfd, buffer, strlen(buffer),
                    Ο,
                    &serveraddr, sizeof(serveraddr)) < 0) {
                perror("sendto");
                exit(1);
        }
        /* receive info */
        if ((n = recvfrom(sockfd, buffer, 1024, 0,
                           &serveraddr, &serveraddrlen))
            < 0) {
                perror("recvfrom");
                exit(1);
        buffer[n] = ' \setminus 0';
        if ((n=analyze(buffer))== 0) {
                printf("transaction failure, "
                        "try again\n");
        } else if (n<0) {</pre>
                printf("login failed, try again\n");
                init();
        }
}
```

merchantname, password);

1 Create a socket with the socket call.

AF\_INET specifies the Internet communication domain. If AF\_OSI were supported, it could be used to create a socket for OSI communications. The socket type SOCK\_DGRAM is specified for UDP or connectionless communication.

Contrast the socket call with the t\_open call in the XTI client example (Section B.2.4).

2 The serveraddr is of type sockaddr\_in, which is dictated by the communication domain of the socket (AF\_INET). The socket address for the Internet communication domain contains an Internet address and a

}

16-bit port number, which uniquely identifies an entity on the network. For the TCP/IP protocol suite, which includes UDP, this is the Internet address of the server and the port number on which it is listening.

Note that the information contained in the sockaddr\_in structure is dependent on the address family (or the protocol).

- **3** Getting information about the server depends on the protocol or the address family. To get the IP address of the server, you can use the gethostbyname routine.
- **4** SERVER\_PORT is defined in the <common.h> header file. It is a short integer, which helps identify the server process from other application processes.
- 5 Client issues a connect call to connect to the server. When the connect call is used with a connectionless protocol, it allows the client to store the server's address locally. This means that the client does not have to specify the server's address each time it sends a message.
- 6 Send data with the sendto call.
- 7 Receive data with the recyfrom call.

# **B.2.3 XTI Server Program**

Example B-7 implements a server using the XTI library for network communication. It is an alternative design for a communication program that makes it transport independent. Compare this program with the socket server program in Example B-5. This program has the limitations described at the beginning of the appendix.

# **Example B-7: Connectionless XTI Server Program**

### Example B-7: (continued)

```
struct sockaddr_in
                        serveraddr;
struct hostent
                        *he;
                        pid;
int
struct t_bind
                        *bindreqp;
struct t_bind
                        *bindretp;
signal(SIGCHLD, SIG_IGN);
/* Create a transport endpoint for the communications */
if ((xtifd = t_open("/dev/streams/xtiso/udp",
                    O_RDWR, NULL)) < 0) {
        xerror("xti_open", xtifd);
        exit(1);
}
bzero((char *) &serveraddr,
                                                      2
      sizeof(struct sockaddr_in));
serveraddr.sin_family = AF_INET;
                                                      3
serveraddr.sin_addr.s_addr = htonl(INADDR_ANY);
serveraddr.sin_port = htons(SERVER_PORT)
/* allocate structures for the t_bind call */
if (((bindreqp=(struct t_bind *)t_alloc(xtifd,
                                        T_BIND_STR,
                                        T_ALL))
        == NULL)
    ((bindretp=(struct t_bind *)t_alloc(xtifd,
                                        T_BIND_STR,
                                        T_ALL))
        == NULL)) {
        xerror("xti_alloc", xtifd);
        exit(3);
}
                   = (char *)&serveraddr;
bindreqp->addr.buf
bindreqp->addr.len
                       = sizeof(serveraddr);
 * Specify how many pending connections can be
 * maintained, while we finish "accept" processing
 * /
bindreqp->qlen
                        = 8;
                                                      6
if (t_bind(xtifd, bindreqp, (struct t_bind *)NULL)
                                                      7
        xerror("xti_bind", xtifd);
        exit(4);
}
```

# **Example B-7: (continued)**

```
\ensuremath{^{\star}} 
 Now the server is ready to accept connections
         * on this socket. For each connection, fork a child
         * process in charge of the session, and then resume
         * accepting connections.
         * /
        transactions(xtifd);
}
                                                                8
transactions(int fd)
        int
                                bytes;
        char
                                 *reply;
        int
                                 dcount;
        int
                                flags;
                                datapipe[MAXBUFSIZE+1];
        char
        struct t_unitdata unitdata;
struct sockaddr_in clientadd
                                clientaddr;
         \mbox{\ensuremath{^{\star}}} Look at the data buffer and parse commands.
         * If more data required, go get it.
         * /
         while (1) {
                unitdata.udata.maxlen = MAXBUFSIZE;
                = 0;
                unitdata.opt.maxlen
                if ((dcount=t_rcvudata(fd, &unitdata, &flags))9
                         /* if disconnected bid a goodbye */
                         if (t_errno == TLOOK) {
                                 int tmp = t_look(fd);
                                 if (tmp != T_DISCONNECT) {
                                         t_scope(tmp);
                                 } else {
                                         exit(0);
                         xerror("transactions_receive", fd);
                         break;
                if (unitdata.udata.len == 0) {
                         return(0);
```

# Example B-7: (continued)

- 1 The t\_open call specifies a device special file name, which is /dev/streams/xtiso/udp in this example. This file name provides the necessary abstraction for the UDP transport protocol over IP. Unlike the socket interface, where you specify the address family (for example, AF\_INET), this information is already represented in the choice of the device special file. The /dev/streams/xtiso/udp file implies both UDP transport and Internet Protocol. See the Chapter 5 for information about STREAMS devices. Contrast the t\_open call with the socket call in Section B.2.1.
- 2 The serveraddr is of type sockaddr\_in, which is dictated by the communication domain or address family of the socket (AF\_INET). The socket address for the Internet communication domain contains an Internet address and a 16-bit port number, which uniquely identifies an application entity on the network. For TCP/IP and UDP/IP this is the Internet address of the server and the port number on which it is listening. The information contained in the sockaddr\_in structure is dependent on the address family (or the protocol).
- **3** AF INET specifies the Internet communication domain or address family.
- 4 INADDRANY signifies any attached interface adapter on the system. All numbers must be converted to the network format using appropriate macros. See the following reference pages for more information: hton1(3), htons(3), ntoh1(3), and ntohs(3).
- **5** SERVER\_PORT is defined in the <common.h> header file. It is a short integer, which helps identify the server process from other application processes. Numbers from 0 to 124 are reserved.
- **6** Specify the number of pending connections the server can queue while it processes the last request.

- 7 Bind the server's address with the t\_bind call. The combination of the address and port number identify it uniquely on the network. After the server process' address is bound, the server process is registered on the system and can be identified by the lower level kernel functions to which to direct any requests.
- 8 Each incoming message packet is accepted and passed to the parse function, which tracks the information provided, such as the merchant's login ID, password, and customer's credit card number. This process is repeated until the parse function identifies a complete transaction and returns a response packet, to be sent to the client program.

The client program can send information packets in any order (and in one or more packets), so the parse function is designed to remember state information sufficient to deal with this unstructured message stream.

Since this program uses a connectionless (datagram) protocol, it uses t\_sndudata and t\_rcvudata to send and receive data, respectively.

- **9** Receive data with the t rcvudata function.
- **10** Send data with the t\_sndudata function.

# **B.2.4 XTI Client Program**

Example B-8 This sample program implements an XTI client that can communicate with the XTI server in Example B-7. It uses the XTI interface in the connectionless, or datagram, mode.

# **Example B-8: Connectionless XTI Client Program**

```
* This file contains the main XTI client code
 * for a connectionless mode of communication.
 * usage: client [serverhostname]
#include "client.h"
main(int argc, char *argv[])
                               xtifd;
        int
        struct sockaddr_in serveraddr; struct hostent *he;
        int
                                 *serverhost = "localhost";
        char
                             ^serverhost =
*serverhostp;
        struct hostent
                               buffer[MAXBUFSIZE+1];
        char
        char
                                 inbuf[MAXBUFSIZE+1];
```

### Example B-8: (continued)

```
struct t_unitdata
                        unitdata;
                        sndcall;
struct t_call
                        rcvcall;
struct t_call
                        flags = 0;
int
if (argc>1) {
        serverhost = argv[1];
init();
if ((xtifd = t_open("/dev/streams/xtiso/udp",
                    O_RDWR, NULL)) < 0) {
        xerror("xti_open", xtifd);
        exit(1);
}
bzero((char *) &serveraddr,
                                                      2
     sizeof(struct sockaddr_in));
serveraddr.sin_family = AF_INET;
                                                      3
if ((serverhostp = gethostbyname(serverhost)) ==
    (struct hostent *)NULL) {
        fprintf(stderr, "gethostbyname on %s failed\n",
               serverhost);
        exit(1);
bcopy(serverhostp->h_addr,
      (char *)&(serveraddr.sin_addr.s_addr),
      serverhostp->h_length);
 * SERVER_PORT is a short which identifies
 * the server process from other sources.
 * /
serveraddr.sin_port
                           = htons(SERVER_PORT);
if (t_bind(xtifd, (struct t_bind *)NULL,
          (struct t_bind *)NULL) < 0) {</pre>
        xerror("bind", xtifd);
        exit(2);
}
while(1) {
        /* Merchant record */
        sprintf(buffer, "%%%%m%s##%%%%p%s##",
                merchantname, password);
        printf("\n\nSwipe card, enter amount: ");
        fflush(stdout);
```

### Example B-8: (continued)

}

```
printf("bye...\n");
                exit(0);
        soundbytes();
        sprintf(buffer, "%s%%%a%s##%%%n%s##",
                buffer, inbuf, swipecard());
        unitdata.addr.buf
                               = (char *)&serveraddr;
        unitdata.addr.len
                               = sizeof(serveraddr);
        unitdata.udata.buf
                               = buffer;
        unitdata.udata.len
                               = strlen(buffer);
        unitdata.opt.len
                                = 0;
        if (t_sndudata(xtifd, &unitdata) < 0) {</pre>
                xerror("t_snd", xtifd);
                exit(1);
        }
        unitdata.udata.maxlen = MAXBUFSIZE;
        unitdata.addr.maxlen
                                = sizeof(serveraddr);
        /* receive info */
        if ((t_rcvudata(xtifd, &unitdata, &flags))
            < 0) {
                xerror("t_rcv", xtifd);
                exit(1);
        }
        buffer[unitdata.udata.len] = '\0';
        if ((n=analyze(buffer))== 0) {
                printf("transaction failure, "
                       "try again\n");
        } else if (n<0) {</pre>
                printf("login failed, try again\n");
                init();
        }
}
```

if (scanf("%s", inbuf) == EOF) {

1 The t\_open call specifies a device special file name; for example /dev/streams/xtiso/udp. This file name provides the necessary abstraction for the UDP transport protocol over IP. Unlike the socket interface, where you specify the address family (for example, AF\_INET), this information is already represented in the choice of the device special file. The /dev/streams/xtiso/udp file implies both UDP transport and Internet Protocol. See the Chapter 5 for information about STREAMS devices.

Contrast the t\_open call with the socket call in Section B.2.2.

- 2 The serveraddr is of type sockaddr\_in, which is dictated by the communication domain of the socket (AF\_INET). The socket address for the Internet communication domain contains an Internet address and a 16-bit port number, which uniquely identifies an entity on the network. For the TCP/IP protocol suite, which includes UDP, this is the Internet address of the server and the port number on which it is listening.
  - The information contained in the sockaddr\_in structure is dependent on the address family (or the protocol).
- **3** AF\_INET specifies the Internet communication domain. If AF\_OSI were supported it could be used to create a socket for OSI communications.
- **4** Getting information about the server depends on the protocol or the address family. To get the IP address of the server, you can use the gethostbyname(3) routine.
- **5** SERVER\_PORT is defined in the <common.h> header file. It is a short integer, which helps identify the server process from other application processes.
- 6 Bind the server address with the t\_bind function to enable the client to start sending and receiving data.
- **7** Send data with the t\_sndudata function.
- **8** Receive data with the t\_rcvudata function.

# **B.3 Common Code**

The following header and database files are required for all or several of the client and server portions of this application:

- <common.h>
- <server.h>
- serverauth.c
- serverdb.c
- xtierror.c
- <client.h>
- clientauth.c
- clientdb.c

#### B.3.1 The common.h Header File

Example B-9 shows the <common.h> header file. It contains common header files and constants required by all sample programs.

# Example B-9: The common.h Header File

```
#include <sys/types.h>
#include <sys/socket.h>
#include <sys/errno.h>
#include <netinet/in.h>
#include <netdb.h>
#include <string.h>
#include <stdio.h>
#include <signal.h>
#include <stdlib.h>
#include <fcntl.h>
#include <xti.h>
#define SEPARATOR
#define PREAMBLE "%%"
#define PREAMBLELEN 2
#define POSTAMBLE "##"
#define POSTAMBLELEN 2
/* How to contact the server */
#define SERVER_PORT 1234
/* How to contact the client (for datagram only) */
#define CLIENT_PORT 1235
```

#define MAXBUFSIZE 4096

1 List of header files to include.

- 2 These statements define constants that allow more effective parsing of data exchanged between the server and client.
- 3 SERVER\_PORT is a well known port that is arbitrarily assigned by the programmer so that clients can communicate with the server. SERVER\_PORT is used to identify the service to which you want to connect. Port numbers 0 through 1024 are reserved for the system. Programmers can choose a number, as long as it does not conflict with any other applications. While debugging, this number is chosen randomly (and by trial and error). For a well-distributed application, some policy must be used to avoid conflicts with other applications.

#### B.3.2 The server.h Header File

Example B-10 shows the <server.h> header file. It contains the data structures for accessing the server's database, as well as the data structures for analyzing and synthesizing messages to and from clients.

# Example B-10: The server.h Header File

```
#include "common.h"
struct merchant {
       char *name;
        char *passwd;
};
struct customer {
                                 *cardnum;
        char
        char
                                 *name;
       int
                                 limit;
        int balance;
struct transaction *tlist;
        /* presumably other data */
};
struct transaction {
       struct transaction *nextcust; struct transaction *nextglob;
        struct customer
                                 *whose;
                                 *merchantname;
        char
                                 amount;
        int
        char
                                 *verification;
};
extern struct transaction *alltransactions;
extern struct merchant merchant[];
extern int
                        merchantcount;
extern struct customer customer[];
extern int
                       customercount;
#define INVALID (struct transaction *)1
#define MERCHANTAUTHERROR
                                 "%%A##"
                                "%%U##"
#define USERAUTHERROR
#define USERAMOUNTERROR
                                "%%V##"
#define TRANSMITERROR
                                "deadbeef"
/* define transaction_status flags */
#define NAME
                                 0x01
#define PASS
                                 0 \times 02
#define AMOUNT
                                 0 \times 04
#define NUMBER
                                 0x08
                                 0x03
#define AUTHMASK
#define VERIMASK
                                  0x0C
```

# Example B-10: (continued)

# B.3.3 The serverauth.c File

Example B-11 shows the serverauth.c file.

# Example B-11: The serverauth.c File

```
* Authorization information (not related to the
* networking interface)
* /
#include "server.h"
* Currently a simple non-encrypted password method to search db
* /
authorizemerchant(char *merch, char *password)
        struct merchant *mp;
        for(mp = merchant; (mp)->name != (char *)NULL; mp++) {
                if (!strcmp(merch, (mp)->name)) {
                       return (!strcmp(password, (mp)->passwd));
        return(0);
}
struct transaction *
verifycustomer(char *num, int amount, char *merchant)
        char buf[64];
        struct customer
                                *cp;
                               *tp;
        struct transaction
        for(cp = customer; (cp)->cardnum != NULL; cp++) {
                if (!strcmp(num, (cp)->cardnum)) {
                        if (amount <= (cp)->balance) {
                                (cp)->balance -= amount;
                                if ((tp = malloc(sizeof(
                                           struct transaction)))
                                    == NULL) {
                                        printf("Malloc error\n");
                                        return(NULL);
```

# Example B-11: (continued)

```
tp->merchantname = merchant;
                                  tp->amount = amount;
sprintf(buf, "v%012d", time(0));
                                  if ((tp->verification =
                                       malloc(strlen(buf)+1))
                                       == NULL) {
                                           printf("Malloc err\n");
                                           return(NULL);
                                  }
                                  strcpy(tp->verification, buf);
                                  tp->nextcust = cp->tlist;
                                  tp->whose = cp;
                                  cp->tlist = tp;
                                  tp->nextglob = alltransactions;
                                  alltransactions = tp;
                                  return(tp);
                          } else {
                                  return(NULL);
        return(INVALID);
}
int
                          transaction_status;
int
                          authorized = 0;
int
                          amount = 0;
                         number[256];
char
char
                         Merchant[256];
char
                         password[256];
char *
parse(char *cp)
                          *dp, *ep;
        char
        unsigned char
                         type;
                          doauth = 0;
        int
        char
                          *buffer;
        dp = cp;
        if ((buffer=malloc(256)) == NULL) {
                 return(TRANSMITERROR);
        while (*dp) {
                 /* terminate the string at the postamble */
                 if (!(ep=strstr(dp, POSTAMBLE))) {
                          return(TRANSMITERROR);
                 \star ep = ' \setminus 0';
```

# Example B-11: (continued)

```
ep = ep + POSTAMBLELEN;
/* search for preamble */
if (!(dp=strstr(dp, PREAMBLE))) {
        return(TRANSMITERROR);
dp += PREAMBLELEN;
/* Now get the token */
type = *dp++;
switch(type) {
        case 'm':
                strcpy(Merchant, dp);
                transaction_status |= NAME;
                break;
        case 'p':
                strcpy(password, dp);
                transaction_status |= PASS;
        case 'n':
                transaction_status |= NUMBER;
                strcpy(number, dp);
                break;
        case 'a':
                transaction_status |= AMOUNT;
                amount = atoi(dp);
                break;
        default:
                printf("Bad command\n");
                return(TRANSMITERROR);
if ((transaction_status & AUTHMASK) == AUTHMASK) {
        transaction_status &= ~AUTHMASK;
        authorized = authorizemerchant(
                     Merchant,
                      password);
        if (!authorized) {
                printf("Merchant not"
                       " authorized\n");
                return(MERCHANTAUTHERROR);
        }
 * If both amount and number gathered,
 * do verification
* /
if ((authorized) &&
    ((transaction_status&VERIMASK)
     ==VERIMASK)) {
        struct transaction *tp;
```

### Example B-11: (continued)

```
transaction_status &= ~VERIMASK;
                /* send a verification back */
                if ((tp=verifycustomer(number,
                                        amount,
                                        Merchant))
                    == NULL) {
                        return(USERAMOUNTERROR);
                } else if (tp==INVALID) {
                        return(USERAUTHERROR);
                } else {
                         sprintf(buffer,
                            "%%%%%s##%%%%C%s##%%%%m%s##",
                            tp->verification,
                            tp->whose->name,
                            tp->merchantname);
                        return(buffer);
                }
        dp = ep;
return(NULL);
```

1 This function parses the incoming data, which includes the merchant authorization information, customer's credit card number, and the amount the customer is charging. Note that the function can not assume that all of the information is available in one message because the underlying TCP protocol is stream-oriented. This function can be simplified if a datagram type service is used or if a protocol that uses sequenced packets (SEQPACKET) is used. The function is designed to accept pieces of information in any order and in one or more message blocks.

### B.3.4 The serverdb.c File

Example B-12 shows the serverdb.c file.

# Example B-12: The serverdb.c File

### **Example B-12: (continued)**

```
{"magic",
                                    "magic"},
          "gasco",
                                     "gasco"},
          "furnitureco",
                                     "abc" } ,
          "groceryco",
                                     "groceryco" },
          "bakeryco",
                                     "bakeryco" },
          "restaurantco",
                                    "restaurantco"},
                                    NULL}
         {NULL.
};
int merchantcount = sizeof(merchant)/sizeof(struct merchant)-1;
struct customer customer[] = {
           "4322546789701000", "John Smith", "4322546789701001", "Bill Stone",
                                                                800 },
                                                      1000,
                                                       2000,
                                                                200
           "4322546789701002", "Dave Adams",
                                                      1500,
                                                                500
           "4322546789701003", "Ray Jones",
                                                                800
                                                      1200,
           "4322546789701004", "Tony Zachry",
                                                                100 },
                                                      1000.
           "4322546789701005", "Danny Einstein", 5000,
           "4322546789701006", "Steve Simonyi", 10000, "4322546789701007", "Mary Ming", 1100,
                                                                5800},
                                                                700 },
           "4322546789701008", "Joan Walters",
                                                                780 },
                                                     800,
           "4322546789701009", "Gail Newton", "4322546789701010", "Jon Robertson",
                                                      1000,
                                                                900 },
                                                      1000,
                                                                1000},
           "4322546789701011", "Ellen Bloop",
                                                      1300,
                                                                800 },
           "4322546789701012", "Sue Svelter",
                                                      1400,
                                                                347
           "4322546789701013", "Suzette Ring",
                                                                657
                                                      1200,
           "4322546789701014", "Daniel Mattis", 1600,
                                                                239 },
           "4322546789701015", "Robert Esconis", 1800,
                                                                768 },
           "4322546789701016", "Lisa Stiles",
                                                                974 },
                                                      1100,
           "4322546789701017", "Bill Brophy",
                                                                800 },
                                                      1050,
           "4322546789701018", "Linda Smitten", 4000,
                                                                200
           "4322546789701019", "John Norton", 1400, "4322546789701020", "Danielle Smith", 2000,
                                                                900
                                                                640
           "4322546789701021", "Amy Olds",
                                                      1300,
                                                                100
           "4322546789701022", "Steve Smith",
                                                      2000,
                                                                832 },
           "4322546789701023", "Robert Smart",
                                                                879 },
                                                     3000,
           "4322546789701024", "Jon Harris",
                                                      500,
                                                                146 },
           "4322546789701025", "Adam Gershner", 1600,
                                                                111 },
           "4322546789701026", "Mary Papadimis", 2000,
                                                                382 },
           "4322546789701027", "Linda Jones", 1300, 
"4322546789701028", "Lucy Barret", 1400, 
"4322546789701029", "Marie Gilligan", 1000,
                                                                578 },
                                                                865
                                                                904 },
           "4322546789701030", "Kim Coyne",
                                                      3000,
                                                                403
           "4322546789701031", "Mike Storm",
                                                                5183},
                                                      7500,
           "4322546789701032", "Cliff Clayden", 750,
                                                                430 },
           "4322546789701033", "John Turing",
                                                                800 },
                                                      4000,
           "4322546789701034", "Jane Joyce",
                                                                8765},
                                                      10000,
           "4322546789701035", "Jim Roberts",
                                                      4000,
                                                                3247},
           "4322546789701036", "Stevw Stephano", 1750,
                                                                894 },
         NULL, NULL
};
struct transaction *
```

# Example B-12: (continued)

```
alltransactions = NULL;
int customercount = sizeof(customer)/sizeof(struct customer)-1;
```

#### B.3.5 The xtierror.c File

Example B-13 shows the xtierror.c file. It is used to generate a descriptive message in case of an error. Note that for asynchronous errors or events, the t\_look function is used to get more information.

# Example B-13: The xtierror.c File

```
#include <xti.h>
#include <stdio.h>
int
xerror(char *marker, int fd)
        fprintf(stderr, "%s error [%d]\n", marker, t_errno);
        t_error("Transport Error");
        if (t_errno == TLOOK) {
               t_scope(t_look(fd));
}
int
t_scope(int tlook)
        char *tmperr;
        switch(tlook) {
                case T_LISTEN:
                        tmperr = "connection indication";
                        break;
                case T_CONNECT:
                        tmperr = "connect confirmation";
                        break;
                case T_DATA:
                        tmperr = "normal data received";
                        break;
                case T_EXDATA:
                        tmperr = "expedited data";
                        break;
                case T_DISCONNECT:
                        tmperr = "disconnect received";
                        break;
                case T_UDERR:
                        tmperr = "datagram error";
                        break;
                case T_ORDREL:
                        tmperr = "orderly release indication";
```

# Example B-13: (continued)

#### B.3.6 The client.h Header File

Example B-14 shows the client.h header file.

### Example B-14: The client.h File

```
#include "common.h"

extern char merchantname[];
extern char password[];
```

# B.3.7 The clientauth.c File

Example B-15 shows the clientauth.c file. It contains the code that obtains the merchant's authorization, as well as the logic to analyze the message sent from the server. The resulting message is interpreted to see if the authorization was granted or rejected by the server.

#### Example B-15: The clientauth.c File

```
#include "client.h"
init()
{
         printf("\nlogin: "); fflush(stdout);
         scanf("%s", merchantname);

         printf("Password: "); fflush(stdout);
         scanf("%s", password);

         srandom(time(0));
}
```

# **Example B-15: (continued)**

```
/* simulate some network activity via sound */
soundbytes()
{
        int i;
        for(i=0;i<11;i++) {
                printf("");
                fflush(stdout);
                usleep(27000*(random()%10+1));
        }
}
analyze(char *cp)
        char *dp, *ep;
        unsigned char type;
        char customer[128];
        char verification[128];
        customer[0] = verification[0] = '\0';
        dp = cp;
        while ((dp!=NULL) && (*dp)) {
                /* terminate the string at the postamble */
                if (!(ep=strstr(dp, POSTAMBLE))) {
                        return(0);
                *ep = ' \setminus 0';
                ep = ep + POSTAMBLELEN;
                /* search for preamble */
                if (!(dp=strstr(dp, PREAMBLE))) {
                        return(0);
                dp += PREAMBLELEN;
                /* Now get the token */
                type = *dp++;
                switch(type) {
                        case 'm':
                                if (strcmp(merchantname, dp)) {
                                         return(0);
                                 }
                                break;
                        case 'c':
                                 strcpy(customer, dp);
                                 break;
                        case 'U':
                                 printf("Authorization denied\n");
                                 return(1);
```

# **Example B-15: (continued)**

```
case 'V':
                                 printf("Amount exceeded\n");
                                 return(1);
                        case 'A':
                                 return(-1);
                         case 'v':
                                 strcpy(verification, dp);
                                 break;
                        default:
                                 return(0);
                dp = ep;
        if (*customer && *verification) {
                printf("%s, verification ID: %s\n",
                       customer, verification);
                return(1);
        return(0);
}
```

# B.3.8 The clientdb.c File

Example B-16 shows the clientdb.c file. It contains a database of customer credit card numbers used to simulate the card swapping action. In a real world application, a magnetic reader reads the numbers through an appropriate interface. Also, the number cache is not required for a real world application.

# Example B-16: The clientdb.c File

```
/*
  *
  * Database of customer credit card numbers to simulate
  * the card swapping action. In practice the numbers
  * will be read by magnetic readers through an
  * appropriate interface.
  */

#include <time.h>

char merchantname[256];
char password[256];

char *numbercache[] = {
    "4322546789701000",
    "4322546789701001",
    "4322546789701002",
    "4222546789701002",
    "4222546789701002",
    "4222546789701002",
    "4222546789701002",
    "45ake id */
```

# **Example B-16: (continued)**

```
"4322546789701003",
        "4322546789701004",
        "4322546789701005",
        "4322546789701006",
        "4322546789701007",
        "4322546789701008",
        "4322546789701009",
        "4322546789701010",
        "4322546789701011",
        "4322546789701012",
        "4322546789701013",
        "4322546789701014",
        "4322546789701015",
        "4322546789701016",
        "4322546789701017",
        "4322546789701018",
        "4222546789701018",
                                        /* fake id */
        "4322546789701019",
        "4322546789701020",
        "4322546789701021",
        "4322546789701022",
        "4322546789701023",
        "4322546789701024",
        "4322546789701025",
                                        /* fake id */
        "2322546789701025",
        "4322546789701026",
        "4322546789701027",
        "4322546789701028",
        "4322546789701029",
        "4322546789701030",
        "4322546789701031",
        "4322546789701032",
        "4322546789701033",
        "4322546789701034",
        "4322546789701035",
        "4322546789701036",
};
#define CACHEENTRIES (sizeof(numbercache)/sizeof(char *))
char *
swipecard()
{
        return(numbercache[random()%CACHEENTRIES]);
}
```

# TCP Specific Programming Information

This appendix contains information about performance aspects of the Transport Control Protocol (TCP). It discusses how programs can influence TCP throughput by controlling the window size used by TCP via socket options.

# C.1 TCP Throughput and Window Size

TCP throughput depends on the transfer rate, which is the rate at which the network can accept packets, and the round-trip time, which is the delay between the time a TCP segment is sent and the time an acknowledgement arrives for that segment. These factors determine the amount of data that must be buffered (the window) prior to receiving acknowledgment to obtain maximum throughput on a TCP connection.

If the transfer rate or the round-trip time or both is high, the default window size used by TCP may be insufficient to keep the pipe fully loaded. Under these circumstances, TCP throughput can be limited because the sender is required to stall until acknowledgements for prior data are received.

The receive socket buffer size determines the maximum receive window for a TCP connection. The transfer rate from a sender can also be limited by the send socket buffer size. DEC OSF/1 currently uses a default value of 32768 bytes for TCP send and receive buffers.

# C.2 Programming the TCP Socket Buffer Sizes

An application can override the default TCP send and receive socket buffer sizes by using the setsockopt system call specifying the SO\_SNDBUF and SO\_RCVBUF options, prior to establishing the connection. The largest size that can be specified with the SO\_SNDBUF and SO\_RCVBUF options is limited by the kernel variable sb\_max. See Section C.3.1 for information about increasing this value.

For maximum throughput, Digital recommends send and receive socket buffers on both ends of the connection be of equal size.

When writing programs that use the setsockopt system call to change a TCP socket buffer size (SO\_SNDBUF, SO\_RCVBUF), note that the actual socket buffer size used for a TCP connection can be larger than the specified

value. This situation occurs when the specified socket buffer size is not a multiple of the TCP Maximum Segment Size (MSS) to be used for the connection.

TCP determines the actual size, and the specified size is rounded up to the nearest multiple of the negotiated MSS. For local network connections, the MSS is generally determined by the network interface type and its maximum transmission unit (MTU).

# **C.3 TCP Window Scale Option**

DEC OSF/1 implements the TCP window scale option, as defined in RFC 1323: *TCP Extensions for High Performance*. The TCP window scale option, which allows larger windows to be used, was designed to increase throughput of TCP over high bandwidth, long delay networks. This option may also increase throughput of TCP in local FDDI networks.

The window field in the TCP header is 16 bits. Therefore, the largest window that can be used without the window scale option is 2\*\*16 (64KB). When the window scale option is used between cooperating systems, windows up to (2\*\*30)-1 bytes are allowed. The option, transmitted between TCP peers at the time a connection is established, defines a scale factor which is applied to the window size value in each TCP header to obtain the actual window size.

The maximum receive window, and therefore the scale factor offered by TCP during connection establishment, is determined by the maximum receive socket buffer space.

If the receive socket buffer size is greater than 65535 bytes, during connection establishment, TCP will specify the Window Scale option with a scale factor based on the size of the receive socket buffer. Both systems involved in the TCP connection must send the Window Scale option in their SYN segments for window scaling to occur in either direction on the connection. As stated previously, Digital recommends that, for maximum throughput, send and receive buffers on both ends of the connection be of equal size.

## C.3.1 Increasing the Socket Buffer Size Limit

The sb\_max kernel variable limits the amount of socket buffer space that can be allocated for each send and receive buffer. The current default is 128KB but optionally you can increase it.

For local FDDI connections, the current value is sufficient. For long delay, high bandwidth paths, values greater than 128KB may be required.

To change the sb\_max kernel variable, use the dbx -k command as root. The following example shows how to increase the sb\_max variable in the

kernel disk image, as well as the kernel currently in memory, to 150KB:

See dbx(1) for a description of the dbx assign and patch commands.

# Information for Token Ring Driver Developers

This appendix contains the following information for developers of Token Ring drivers for Digital UNIX:

- Enabling source routing
- Using canonical addresses
- Avoiding unaligned access
- Setting fields in the softc structure of the driver

## **D.1 Enabling Source Routing**

Source routing is a bridging mechanism that systems on a Token Ring local area network (LAN) use to send messages to a system on another interconnected Token Ring LAN. Under this mechanism, the system that is the source of a message uses a route discover process to determine the optimum route over Token Ring LANs and bridges to a destination system.

To use the Token Ring source routing module you must add the TRSRCF option to your kernel configuration file. Use the doconfig -c command to add the TRSRCF option, as follows:

- 1. Enter the doconfig -c HOSTNAME command from the superuser prompt (#). HOSTNAME is the name of your system in uppercase letters; for example, for a system called host1 you would enter:
  - # doconfig -c HOST1
- 2. Add TRSRCF to the options section of the kernel configuration file.

Enter y when the system asks whether you want to edit the kernel configuration file. The doconfig command allows you to edit the configuration file with the ed editor. For information about using the ed editor, see ed(1).

The following ed editing session shows how to add the TRSRCF option to the kernel configuration file for host1. The number of the line after

which you append the new line can differ between kernel configuration files:

3. After the new kernel is built, you must move it from the directory where doconfig places it to the root directory ( / ) and reboot your system.

For detailed information on reconfiguring your kernel or the doconfig command see the *System Administration* manual.

The Token Ring source routing functionality is initialized if the trn\_units variable is greater than or equal to 1. The trn\_units variable indicates the number of Token Ring adapters initialized on the system.

The driver should declare trn\_units as follows:

```
extern int trn_units;
```

At the end of its attach routine, the driver should increment the trn\_units variable as follows:

```
trn units++;
```

For information on source routing management see the *Network Administration* manual.

# **D.2 Using Canonical Addresses**

The Token Ring driver requires that the destination address (DA) and source address (SA) in the Media Access Control (MAC) header be in the canonical form while presenting it to the layers above the driver.

The canonical form is also known as the Least Significant Bit (LSB) format. It differs from the noncanonical form, known as the Most Significant Bit

(MSB) format, in that it transmits the LSB first. The noncanonical form transmits the MSB first. The two formats also differ in that the bit order within each octet is reversed.

For example, the following address is in noncanonical form:

```
10:00:d4:f0:22:c4
```

The same address in canonical form is as follows:

```
08-00-2b-0f-44-23
```

If the hardware does not present the driver with a canonical address in the MAC header, you should convert the address to canonical form before passing it up to the higher layers. The haddr\_convert kernel routine is available for converting canonical addresses to noncanonical, and vice versa. It has the following format:

```
haddr_convert( addr) unsigned char *addr
```

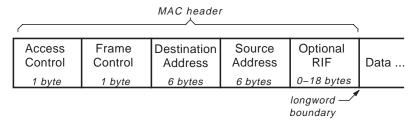
The *addr* variable is a pointer to the 6 bytes of the address that require conversion from either noncanonical to canonical or canonical to noncanonical form. The converted address is returned in the same buffer.

# **D.3 Avoiding Unaligned Access**

The frame that the driver receives consists of the Media Access Control (MAC) header, which includes the Routing Information Field (RIF) and data. Because the length of the RIF can vary between 0 and 18 bytes, the data after the RIF may not be aligned to a longword boundary. To avoid degraded performance, Digital recommends that you pad the RIF field so that data always starts on a longword boundary.

Figure D-1 illustrates the relationship between the components of the MAC header and the data in a typical frame.

Figure D-1: Typical Frame



ZK-0894U-R

# D.4 Setting Fields in the softc Structure of the Driver

The softc structure contains driver-specific information.

You must set the following field of the softc structure in the attach routine of the driver:

sc->isac.ac\_arphrd=ARPHRD\_802;

Here, sc is a pointer to the softc structure, and ARPHRD\_802 is the value of the hardware type used in an Address Resolution Protocol (ARP) packet sent from this interface. A value of 6 for ARPHRD\_802 indicates an IEEE 802 network.

# The Data Link Interface

The data link interface (DLI) is a programming interface that allows programs on a Digital UNIX system directly to use the data link facility to communicate with data link programs running on a remote system.

See Section E.5 for client and server DLI programming examples.

#### E.1 **Prerequisites for DLI Programming**

DLI programming requires both a thorough knowledge of the C programming language and experience writing system programs. If you intend to use the Ethernet substructure, you should be familiar with the Ethernet protocol. If you intend to use the 802 substructure, you should be familiar with the 802.2, 802.3, and FDDI protocols.

You should be also be familiar with the following concepts before attempting to write programs to the DLI interface:

Datagram sockets

Your application uses sockets to send and receive Ethernet, 802.3 and FDDI frames. DLI uses datagram sockets only.

For more information about using sockets, see Chapter 4.

Logical Link Control (LLC)

LLC is a sublayer of DLI that provides a set of services determined by a value in the 802.2 frame format.

Physical and multicast addressing

You can send and receive messages over the network using physical or multicast addresses. You can use physical addresses to send messages to a single destination system. Multicast addresses are not associated with any specific system; instead, a packet sent to a multicast address is received by all systems with the multicast address enabled.

For more information about multicast addressing, see Section 4.6.

Standard frame formats

The Ethernet frame format is a proprietary standard that belongs to Digital Equipment Corporation, Intel Corporation, and Xerox Corporation. The IEEE 802.3 frame format is a standard for multivendor networking. The FDDI and IEEE 802.3 frame formats are very similar. Both contain the LLC (or 802.2) frame within them. See Section E.3.1 for more information.

Note that running DLI applications on Digital UNIX requires superuser or root privileges.

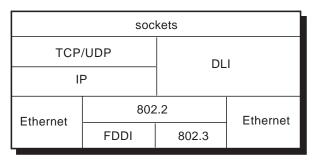
## E.2 DLI Overview

DLI programs transfer data over networks using the standard Ethernet frame format, the Open Systems Interconnect (OSI) 802.3 frame format, or the FDDI frame format. Your Digital UNIX system can run Internet, DECnet, and DLI programs concurrently.

Digital UNIX supports both Ethernet and 802.2 data link services. DLI and IP both run over Ethernet and 802.2. FDDI and 802.3 use the 802.2 Logical Link Control (LLC) as their data link sublayer. TCP and UDP run over IP, providing data delivery and message routing services to the programs that use them. Because DLI provides direct access to the data link layer it does not provide the higher-level services that TCP and UDP do.

Figure E-1 illustrates in greater detail the relationships between DLI and IP, DLI and Ethernet, and DLI and 802.2.

Figure E-1: DLI and the Digital UNIX Network Programming Environment



ZK-0812U-R

Sockets are the user application interface and facilitate access to TCP, UDP, and DLI. See Chapter 4 for information about opening sockets in the DLI communication domain (AF\_DLI).

## E.2.1 DLI Services

DLI provides the following services at the data link layer:

- Datagram service
- Logical Link Control (LLC) layer
  - ISO 802.2 Class I, Type I service
- · Multicast address mode
- Medium Access Control (MAC) layer
  - Ethernet frames
  - 802.3 frames
  - FDDI frames

## **E.2.2** Hardware Support

DLI requires no knowledge of the underlying hardware. It uses Ethernet or FDDI device drivers, which each use the probe routine to determine what devices a particular system has configured. For a complete list of the network devices that Digital UNIX supports, see the Digital UNIX Operating System Version 3.0 Software Product Description 41.61.xx.

To determine which network devices are configured on your system, use the /usr/sbin/netstat -i command, as follows:

#### % /usr/sbin/netstat -i

Name	Mtu	Network	Address	Ipkts	Ierrs	Opkts	0errs	Coll
ln0	1500	<link/>		746	0	234	0	18
ln0	1500	orange-net	host1	746	0	234	0	18
s10*	296	<link/>		0	0	0	0	0
sl1*	296	<link/>		0	0	0	0	0
100	1536	<link/>		74	0	74	0	0
100	1536	100p	localhost	74	0	74	0	0

The output displayed on your screen contains information about the interfaces or devices that your system has configured. In this example, an Ethernet hardware device (ln) is configured, as are two Serial Line Interface Protocol devices (sl0 and sl1). The asterisk (\*) following the sl0 and sl1 indicates that the support for the interfaces has not been turned on yet.

## E.2.3 Using DLI to Access the Local Area Network

A data link on a single local area network (LAN) controller supports multiple concurrent users. Each station represents an available port on the network channel.

Because multiple users simultaneously access the network channel, your program must use addressing mechanisms that ensure delivery of messages to

the correct recipient. Any message you transmit on the network must include an Ethernet or FDDI address that identifies the destination system. The message must also include an additional identifier that directs the message to the correct user on the destination system; this identifier varies according to the frame format you choose to use. DLI builds frames according to the Ethernet, IEEE 802.3, or FDDI standards.

## E.2.4 Including Higher-Level Services

DLI provides only datagram services. Because DLI is a direct interface to the data link layer, it does not offer higher-level services normally provided by Internet and DECnet. Therefore, your application should provide the following kinds of services:

- Packet routing and guaranteed delivery
- Flow control
- Error recovery
- Data segmentation

## E.3 The DLI Socket Address Data Structure

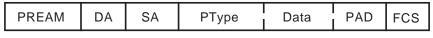
This section describes the Ethernet, 802.3, and FDDI standard frame formats, and the function of the DLI socket address data structure (sockaddr\_dl). It explains how you use sockaddr\_dl to specify the domain address, the network device, and the Ethernet, 802.3, or FDDI substructure.

#### E.3.1 Standard Frame Formats

The following diagrams illustrate the differences and similarities between the Ethernet, 802.3, and FDDI frames.

Figure E-2 illustrates the Ethernet frame format.

Figure E-2: The Ethernet Frame Format



ZK-0687U-R

Figure E-3 illustrates the 802.3 frame format. Note that the 802.3 frame format contains the 802.2 structure, which is illustrated in Figure E-5.

Figure E-3: The 802.3 Frame Format

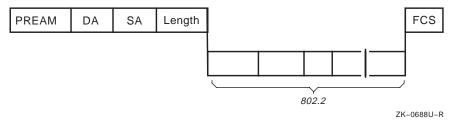


Figure E-4 illustrates the FDDI frame format. The FDDI frame format also contains within it the 802.2 structure illustrated in Figure E-5.

Figure E-4: The FDDI Frame Format

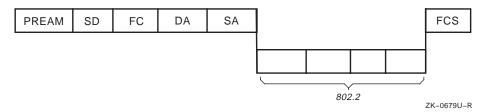
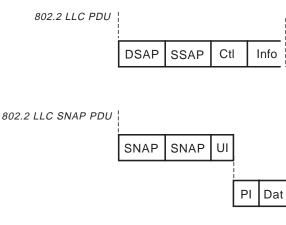


Figure E-5 illustrates the 802.2 LLC PDU and the 802.2 LLC SNAP PDU. One of these two structures is contained within the 802.3 and FDDI frame formats.

Figure E-5: The 802.2 Structures



ZK-0822U-R

Typically, 802 applications use the 802.2 LLC PDU format; however, an application developer may choose to use the 802.2 LLC SNAP PDU format for the following reasons:

- Using the SNAP\_SAP is a convenient way to map Ethernet protocol types on to 802.2 protocols. This is useful for applications that operate over both Ethernet and 802.2, or are migrating from Ethernet to 802.2.
- The I/O control flags (DLI\_NORMAL, DLI\_EXCLUSIVE and DLI\_DEFAULT) are valid only for Ethernet and 802.2 SNAP frames. These flags are meaningless when the non-SNAP 802.2 LLC PDU is used.
- Using the SNAP\_SAP allows a greater number of applications to run over 802.2 because the SNAP SAP has a five byte Protocol ID associated with it. The normal 802.2 LLC PDU, on the other hand, is multiplexed on the 7 most significant bits of the DSAP.

## E.3.2 How the sockaddr\_dl Structure Works

DLI provides a socket address data structure through which you can configure the set of services required for communication at the data link layer. The data structure <code>sockaddr\_dl</code> is used to convey information to DLI when an application binds to the network, or when it transmits a packet to the network. DLI also uses it to convey information to the application when it receives a packet from the network. This includes network device information, the packet format to be used, and addressing information.

The following example shows the DLI socket address structure, which is defined in the header file <dli/dli var.h>:

```
#define DLI_ETHERNET
                           Λ
#define DLI_802
struct sockaddr_dl {
    u_char dli_len;
                                       /* length of sockaddr */
    u_char dli_family; /* address family (AF_DLI) */
struct dli_devid dli_device; /* id of comm device to use */
    u_char dli_substructype;
                                       /* id to interpret following */
                                       /* structure */
                                                  /* Ethernet */
        struct sockaddr_edl dli_eaddr;
         struct sockaddr_eql qli_eaqur, , __emel..., /* oSI 802 support */
/* this needs to have
        caddr_t dli_aligner1;
                                                  /* this needs to have */
                                                   /* longword alignment */
    } choose_addr;
};
```

Any single application can send and receive both Ethernet and 802 substructures. The Ethernet substructure enables applications to communicate across an Ethernet. The 802 substructure enables applications to use 802.2, 802.3, and FDDI protocols to communicate with each other.

You can use system calls to specify values within the socket address structure by using either the Ethernet or 802 substructures.

The fields within the substructures are updated as a function of the system call. For example, the bind system call is used to specify the domain, network device, and most of the substructure. When using the sendto system call to transmit data, the domain, network device, and part of the substructure must be specified. When using the recvfrom system call to receive data, DLI fills in the entire sockaddr structure.

The dli\_econn and dli\_802\_3\_conn user-written subroutines open a socket and bind the associated domain, network device name, protocol type, and other substructure information to the socket. See Section E.5 for examples of the dli\_econn and dli\_802\_3\_conn user-written subroutines.

The following sections describe the functions that the Ethernet and 802.2 substructures provide within the DLI sockaddr\_dl data structure.

## **E.3.3** The Ethernet Substructure

The following example shows the DLI Ethernet socket address substructure:

```
#define DLI EADDRSIZE
struct sockaddr_edl {
   u_char dli_ioctlflg;
                                        /* i/o control flags */
                            /* Ethernet options */
/* Ethernet protocol type */
    u_char dli_options;
    u_short dli_protype;
    u_char dli_target[DLI_EADDRSIZE]; /* Ethernet address of */
                                        /* destination system */
    u_char dli_dest[DLI_EADDRSIZE];
                                       /* Ethernet address used to */
                                        /* address the local system; */
};
                                        /* DLI places the destination */
                                        /* address of an incoming */
                                         /* packet here to be used in */
                                         /* the recyfrom call. This */
                                         /* address can be the sys- */
                                         /* tem's address or a multi */
                                         /* cast address. */
```

The Ethernet substructure specifies the following:

- An I/O control flag for the protocol type (dli\_ioctlflg)
- Whether Ethernet is padded (dli\_options)

The PAD is a 2-byte length field in little-endian after the MAC/LLC header. The following line, from <dli/dli\_var.h>, is the bit that must be set in the dli\_options field to turn on padding:

```
#define DLI_ETHERPAD 0x01 /* Protocol is padded */
```

- The DLI protocol type (dli\_prototype)
- The Ethernet address of the destination system (dli\_target)
- The Ethernet address used to address the local system (dli\_dest)

This information is used to create the Ethernet frame format.

#### E.3.3.1 How Ethernet Frames Work

All Ethernet frames contain a 16-bit identification number called an Ethernet protocol type (PType). When a message arrives at the controller, the protocol type is used to identify which port receives the frame. DLI applications that communicate across the Ethernet must always enable the same Ethernet protocol type. In addition to using protocol types to select a user for an incoming packet, you can configure DLI to select a user as a function of both the protocol type and the physical address of the remote system. This allows several applications in the same system to use the same type, which can make input/output simpler for the application.

## **E.3.3.2** Defining Ethernet Substructure Values

The user specifies the values for the following fields in the Ethernet socket address substructure. The other fields are filled in either by system calls or DLI:

- Destination address (dli\_target[DLI\_EADDRSIZE])
  You can use the dli\_target field to specify the destination address.
- Protocol type (dli\_protype)
   You can use the dli\_prototype field to specify the protocol to be used for data transmission.
- I/O Control Flag (dli\_ioctlflg)

The following sections define the values for user-definable members in the Ethernet substructure.

## **Destination Node Physical Address**

The destination system physical address (DA in Figure E-2) is a 48-bit unique value assigned by the manufacturer to a station on the Ethernet. For example, 08-00-2b-XX-XX-XX is the form a valid Ethernet address takes, with the Xs being replaced by hexadecimal digits. DA is the address of the remote system with respect to the local system.

If you do not specify the DA value with the bind call, you must specify it when sending data by using the sendto call. In addition, you should use the recvfrom call to determine the source of a data message. You can use either the physical address or a multicast address to send messages in the sendto system call.

#### **Protocol Type**

The protocol type (PType in Figure E-2) is a 16-bit value in the Ethernet frame following the source address. The Ethernet driver passes the protocol type to DLI for use in determining the recipient of the data in the frame. With the exception of reserved values, you can use any Ethernet protocol type if it is assigned to you by the manufacturer and not used elsewhere in your system.

The following hexadecimal values are reserved for use by the system:

- 0X 0200 PUP Protocol
- 0X 0800 Internet Protocol
- 0X 0806 Address Resolution Protocol
- 0X 6004 Local Area Transport

- 0X 6003 Phase IV DECnet
- 0X 6002 MOP CCR Protocol
- 0X 6001 MOP Downline Load Protocol
- 0X 9000 MOP Loopback Protocol
- 0X 1000 to 0X 100f Internet Trailer Protocol (used by VAX only)

## I/O Control Flag

The I/O control flag, defined in the header file <dli/dli\_var.h>, is a value that DLI uses to determine how your program reserves a protocol type. It is used by DLI to determine whether to select a user as a function of the protocol type alone or as a function of the combination of the protocol type and the target audience. The following list defines the possible I/O control flags and describes the conditions for their use:

#### NORMAL

Allows your program to exchange messages with one destination system, using only the specified protocol type. When using the NORMAL flag, you must specify the destination system physical address in the bind call and you can use any of the data transfer calls to send and receive data. DLI forwards to the user all messages containing the specified protocol type from the specified target.

#### • EXCLUSIVE

Gives your program exclusive use of the specified protocol type and allows the program to exchange data with any other system using this protocol type. In other words, the program receives all messages with the specified protocol type. When you use the EXCLUSIVE flag, do not specify the target address with the bind call. You must use the sendto and recvfrom calls to exchange data with other systems, and you must specify the target address with the sendto call. In the address structure (returned with recvfrom), DLI fills in the target address with the source address in the Ethernet frame. It also fills in the destination address with the destination address in the Ethernet frame.

#### • DEFAULT

Allows your program to receive messages that contain the specified protocol type and that are meant for no other program on the system. If no other program is bound exclusively to the protocol type or the protocol type/address pair in the message, the socket bound to the protocol type gets the message by default. This mode of operation is recommended for use in programs that listen for messages but do not necessarily send them. When you use the DEFAULT flag, do not specify the target address with the bind call. Use the recyfrom call to receive

data from other systems. If you are using the DEFAULT flag, DLI fills in the target address with the source address in the Ethernet frame. It also fills in the destination address with the destination address in the Ethernet frame.

#### E.3.4 The 802.2 Substructure

The 802.2 substructure enables applications to communicate with each other using the 802.2, 802.3, and FDDI protocols. It uses two basic modes of operation: Class I, Type 1 service, and the services supplied by your application using the 802.2 protocol.

The following example shows the DLI 802.3 socket address substructure:

The 802.2 substructure subsumes both the 802.3 and FDDI frame formats. You can specify values for the following fields:

- Destination system physical address (DA in Figure E-3 and Figure E-4)
- Service class
- Destination service access points (DSAP in Figure E-5)
  - Individual
  - Group
- Source service access point (SSAP in Figure E-5)

The protocol identifier and I/O control field may be required, depending on the type of SSAP you enable.

Control field

#### E.3.4.1 Defining 802 Substructure Values

The following sections define the possible values for all members in the 802 substructure.

#### **Destination Node Physical Address**

The destination system physical address (DA) is a 48-bit unique value assigned by the manufacturer to a station on an Ethernet or FDDI network. For example, 08-00-2b-XX-XX-XX is a valid Ethernet or FDDI address, with the Xs being replaced by hexadecimal digits. This is the address of the remote system with which the application attempts to exchange packets. It

must be specified in the bind call, except when the I/O control field is either EXCLUSIVE or DEFAULT and the service access point (SAP) is a SNAP\_SAP type. The SAP must be specified in the sendto call.

#### **Service Class**

The service class is a value in the 802.2 substructure that determines the capabilities and features provided by the Logical Link Control (LLC) sublayer of the data link layer. The possible service classes are:

• TYPE1

This value causes DLI to interpret all header information and provide Class I, Type 1 service.

#### Note

When Type 1 service is used, the DLI software handles the XID and TEST packets. This is transparent to the application.

DLI uses the source and destination service access points to determine who should receive the message; it interprets the control field on behalf of the user. Whether DLI passes the data field to the user depends on the value of the control field.

• USER

This value provides few services. The user must, therefore, implement most of the 802.2 protocol. In other words, the application must handle the XID and TEST packets. DLI uses the source and destination service access points, but it passes the control field with the data to the user. The user must interpret the control field. This mode must be selected if the application needs to implement Class II, Type 2 service.

## **Destination Service Access Point**

The destination service access point (DSAP) is a field in the 802.2 frame that identifies the application for which the message is intended.

You can use individual or group DSAPs to identify one user or a group of users. You can use group DSAPs only when the service class is set to USER. The possible values for this field are:

• Individual DSAPs

NULL\_SAP — A DSAP consisting of all zeros. You can send TEST and XID commands and responses, but no data, to a NULL\_SAP. (TEST and XID are explained later in this section.) The data link layer uses the

NULL\_SAP to talk to another data link layer, primarily for testing.

User-defined DSAP — Identifies one user for whom the message is intended. The user-defined individual DSAP must be an even number greater than or equal to 2 and less than or equal to 254.

SNAP SAP — The 802.3 Subnetwork Access Protocol.

• Group DSAP (user defined)

Identifies more than one user for whom the message is intended. You can send data to a maximum of 127 group DSAPs on one socket. The user defined group DSAP must be an odd number greater than or equal to 3 and less than or equal to 255. Note that the 255 number is the global SAP and must be enabled like any other group SAP. You can use group SAPs only when the service class is set to USER.

#### **Source Service Access Point**

The source service access point (SSAP) is a field in the 802.2 frame that identifies the address of the application that sent the message. You can enable only one SSAP on a socket. The SSAP must be an even number greater than or equal to 2 and less than or equal to 254.

#### Note

When using the SNAP\_SAP, both the DSAP and SSAP must be set to SNAP\_SAP. In addition, you must specify the protocol identifier and control field. The protocol identifier is five bytes. The control field is one byte. Enabling the SNAP\_SAP is allowed only when the service class is TYPE1.

Note also that IEEE 802.2 standard reserves for its own definition all SAP addresses with the second least significant bit set to 1. It is suggested that you use these SAP values for their intended purposes, as defined in the IEEE 802.2 standard.

## **Control Field**

The control field specifies the packet type. The following values are defined for Class I, Type 1 service, and can also be used in the user-supplied mode to provide Class II, Type 2 service.

#### Note

An application using this user mode is responsible for providing the correct services. For other operations supported by CLASS II service, see the *IEEE Standards for Local Area Networks: Logical Link Control*, published by the Institute of Electrical and Electronics Engineers, Inc.

#### • Exchange Identification

The value XID identifies the exchange identification command or response. An 8-bit format identifier and a 16-bit parameter follow the XID control field. The 16-bit parameter identifies the supported LLC services and the receive window size. The LLC is the top sublayer in the data link layer of the IEEE/Std 802 Local Area Network Protocol. The following values of XID are defined in the DLI header file <dli/dli\_var.h>:

#### - XID PCMD

Exchange identification command with the poll bit set. The exchange identification command conveys the types of LLC services supported and the receive window size to the destination LLC. This command causes the destination LLC to reply with the XID response Protocol Data Unit (PDU) at the earliest opportunity. The poll bit is set to 1, soliciting a response PDU.

#### - XID NPCMD

Exchange identification command with no poll bit set. This command is identical to the previous command, except that you clear the poll bit. No response is expected.

#### - XID PRSP

Exchange identification response with the poll bit set. The Data Link layer uses the exchange identification response to reply to an XID command at the earliest opportunity. The XID response PDU identifies the responding LLC and includes an information field like that defined for the XID command PDU, regardless of what information is present in the information field of the received XID command PDU. The final bit is set to 1, indicating that this response is sent by the LLC as a reply to a soliciting command PDU.

#### XID\_NPRSP

Exchange identification response with no poll bit set. This response is identical to the previous one, except that the final bit is cleared.

#### LLC Protocol Data Unit Test

The value TEST identifies the LLC PDU command or response test. The TEST control field can be followed by a data field. The following values of TEST are defined in the DLI header file <dli/dli\_var.h>:

- TEST\_PCMD

TEST command with the poll bit set. The TEST command tests the LLC-to-LLC transmission path by causing the destination LLC to respond with the TEST response PDU at the earliest opportunity. An information field is optional with this control field value. If used, the receiving LLC returns the information rather than passing it to the user. The poll bit is set to 1, soliciting a response PDU.

- TEST NPCMD

TEST command with no poll bit set. This command is identical to the previous command, except that the poll bit is cleared.

- TEST PRSP

TEST response with the poll bit set. The TEST response PDU is a reply to the TEST command PDU. An information field, if present in the TEST command PDU, is returned in the corresponding TEST response PDU. The final bit is set to 1, indicating that this response is sent by the LLC as a reply to a soliciting command PDU.

- TEST NPRSP

TEST response with no poll bit set. This response is identical to the previous one, except that the final bit is cleared.

## Unnumbered Information Command

The unnumbered information command with no poll set (UI\_NPCMD) sends information to one or more LLCs. The UI\_NPCMD command does not have an LLC response PDU. This is usually passed up to the application. Class I, Type 1 applications generally send and receive data using this command.

# **E.4 Writing DLI Programs**

This section explains how to use Digital UNIX system calls to write DLI programs and describes procedures for specifying values within the Ethernet and 802 substructures.

Section E.5 contains DLI programming examples of the procedures described in this section.

For additional information about how to use sockets and system calls to write application programs, see Chapter 4.

## E.4.1 Supplying Data Link Services

Because DLI provides only a datagram service, a DLI application should provide the services that the higher levels of network software normally provide:

- Flow control DLI programs running on different systems must synchronize data transfer or they will lose data.
- Error recovery DLI reports errors, but your application must recover from them.
- Data segmentation Your application must segment data during transmission. (See Section E.4.7 for information about the buffer size for Ethernet, 802.3 and FDDI packets.)

## E.4.2 Using Digital UNIX System Calls

Your DLI program uses the socket interface with input arguments, structures, and substructures specific to DLI. For example, when issuing the socket system call, your program uses the address format AF\_DLI and the protocol DLPROTO\_DLI.

The beginning of any DLI program must include the header file <dli/dli\_var.h>. Then it should follow the calling sequence shown in Table E-1.

Table E-1: Calling Sequence for DLI Programs

Function	System Call
Create a socket.	socket
Bind the socket to a device by specifying the address family, the frame format type, and the device over which the program will send the data using the sockaddr_dl structure.	bind
Set socket options. This call is optional.	setsockopt
Transfer data.	send recv read write

## Table E-1: (continued)

Function System Call

sendto recvfrom

Deactivate the socket descriptor.

close

See Chapter 4 and the reference page for each system call for more information.

The following sections describe DLI functions, input arguments, and structures.

## E.4.3 Creating a Socket

Your DLI application must create a socket by using the socket system call with the following input arguments:

Address family: AF\_DLI
Socket type: SOCK\_DGRAM
Protocol: DLPROTO\_DLI

The value AF\_DLI specifies the DLI address family. SOCK\_DGRAM creates a datagram socket, which is the only type of socket that DLI allows. DLI does not supply the services necessary for connecting to other programs and for using other socket types. The value DLPROTO\_DLI specifies the DLI protocol module.

The following example shows how the socket call is used to open a socket to DLI:

```
int so;
...
...
if ( (so = socket(AF_DLI,SOCK_DGRAM,DLPROTO_DLI))<0)
    {
        perror("cannot open DLI socket");
        return (-1);
}</pre>
```

## **E.4.4 Setting Socket Options**

Use the setsockopt call to set the following socket options within the sockaddr\_dl structure:

Option	Description
DLI_ENAGSAP	Enables a group service access point (GSAP)
DLI_DISGSAP	Disables a group service access point (GSAP)
DLI_SET802CTL	Sets the 802 control field
DLI_MULTICAST	Enables the reception of all messages addressed to a multicast address

The following code examples show how to use the setsockopt call to set the socket options.

The following example shows how the setsockopt call is used to enable the GSAP option:

```
/* enable GSAPs supplied by user */
j = 3;
i = 0;
while (j < argc ) {
    sscanf(argv[j++], "%x", &k);
    out_opt[i++] = k;
}
optlen = i;
if
(setsockopt(sock,DLPROTO_DLI,DLI_ENAGSAP,&out_opt[0],optlen) < 0){
    perror("dli_setsockopt: Can't enable gsap");
    exit(1);
}</pre>
```

The following example shows how the setsockopt call is used to disable the GSAP option:

```
/* disable all but the last 4 or all GSAPs, */
/* whichever is smallest */
if ( optlen > 4 )
    optlen -= 4;
if
(setsockopt(sock,DLPROTO_DLI,DLI_DISGSAP,&out_opt[0],optlen) < 0){
    perror("dli_setsockopt: Can't disable gsap");
}</pre>
```

The following example shows how the setsockopt call is used to set the

#### 802 control field:

The following example shows how the setsockopt call is used to enable two multicast addresses:

See Section E.5 for more detailed code examples.

## E.4.5 Binding the Socket

After you create the socket, your application must bind the socket to a network device. At this point, you specify the type of format for the message. You assign a name to the socket, where the variable <code>name</code> is a pointer to a structure of the type <code>sockaddr\_dl</code>. Then, you must fill in the <code>sockaddr\_dl</code> data structure and include the appropriate substructure (Ethernet or 802).

To bind the socket, use the following system call:

```
int bind,(
    int socket,
    struct sockaddr_dl *name,
    int namelen );
```

For more information about the bind system call, see the bind(2) reference page.

## E.4.6 Filling in the sockaddr\_dl Structure

Fill in the sockaddr\_dl structure with the following information:

- Address family
- I/O device ID
- Substructure type

## **E.4.6.1** Specifying the Address Family

To specify the address family, use the value AF\_DLI in the socket call.

## E.4.6.2 Specifying the I/O Device ID

The I/O device is the controller over which your program sends and receives data to and from the target system. The I/O device ID consists of the device name, dli\_devname, and the device number, dli\_devnumber. Definitions for each variable follow:

• dli devname

The netstat -i command lists the devices that are available on your system.

• dli\_devnumber

The device number is set up in the system configuration file.

## E.4.6.3 Specifying the Substructure Type

The substructure specifies the type of frame format that the program will use. Definitions for each variable follow:

• dli\_eaddr
Ethernet frame format ( DLI ETHERNET)

• dli\_802addr 802.3 frame format ( DLI\_802)

A program can send and receive Ethernet, 802.3, and FDDI frames, as long as it has a socket associated with each type. For example, your DLI program might communicate with one system using the Ethernet frames and another system using 802.3 or FDDI frames. Your choice of frame formats depends on the frame types used by the target program; however, only one type of frame per socket is allowed.

Your program specifies the packet header for sending your message by filling in the substructure of your choice. Example E-1 shows how to fill the sockaddr\_dl structure for the Ethernet protocol. Example E-2 shows how to fill the sockaddr\_dl structure for the 802 protocol:

## Example E-1: Filling the sockaddr\_dl structure for Ethernet

#### Example E-2: Filling the sockaddr\_dl structure for 802.2

```
* Fill out sockaddr_dl structure for the bind call.
 * Note that we need to determine whether the
 * control field is 8 bits (unnumbered format) or
 * 16 bits (informational/supervisory format). We do this
 * by checking the low order 2 bits, which are both 1 only
 * for unnumbered control fields.
 * /
bzero(&out_bind, sizeof(out_bind));
out_bind.dli_family = AF_DLI;
out_bind.dli_substructype = DLI_802;
bcopy(devname, out_bind.dli_device.dli_devname, i);
out_bind.dli_device.dli_devnumber = devunit;
out_bind.choose_addr.dli_802addr.ioctl = ioctl;
out_bind.choose_addr.dli_802addr.svc = svc;
if(ctl & 3)
   out_bind.choose_addr.dli_802addr.eh_802.ctl.U_fmt=(u_char)ctl;
   out_bind.choose_addr.dli_802addr.eh_802.ctl.I_S_fmt = ctl;
out_bind.choose_addr.dli_802addr.eh_802.ssap = sap;
out_bind.choose_addr.dli_802addr.eh_802.dsap = dsap;
if (ptype)
   bcopy(ptype,out_bind.choose_addr.dli_802addr.eh_802.osi_pi,5);
if (taddr)
    bcopy(taddr, out_bind.choose_addr.dli_802addr.eh_802.dst,
           DLI_EADDRSIZE);
```

## Example E-2: (continued)

```
if ( bind(sock, &out_bind, sizeof(out_bind)) < 0 )
{
    perror("dli_802, can't bind DLI socket");
    return(-1);
}
return(sock);</pre>
```

## E.4.7 Calculating the Buffer Size

The buffer size must be no larger than the controllers on the communicating systems can handle, or you will lose data. The maximum buffer size for Ethernet packets is 1500 bytes.

The maximum buffer size for 802.3 packets is calculated as follows:

The number of bytes in the control field and in the Source SAP are specified in the bind call.

The maximum buffer size for FDDI packets 4352 bytes.

## E.4.8 Transferring Data

A DLI program can use the write, send, or sendto calls to send data and the read, recv, or recvfrom calls to receive data. The X's in Table E-2 indicate the conditions under which you can use the system calls as a function of the I/O control flag set up during the bind call.

#### Note

You must set the target address in the bind call when using the Normal control flag. You do not need to set the target address in the bind call when using the Exclusive or Default control flags. However, if you do not set the target address then you must use the sendto and recvfrom system calls.

Table E-2: Data Transfer System Calls Used With DLI

System Calls	Normal	Exclusive	Default
	Control	Control	Control
write send	X X		

Table E-2: (continued)

System Calls	Normal Control	Exclusive Control	Default Control
sendto	X	X	X
read	X		
recv	X		
recvfrom	X	X	X

When you set the control flag to NORMAL, set the target address in the bind call. Then use any of the following calls to transfer data: write, send, sendto, read, recy, recyfrom.

When you set the control flag to EXCLUSIVE, make the value of the target address in the bind call zero. Then, set the target address in the sendto call. Use only the sendto and recvfrom calls to transfer data.

When you set the control flag to DEFAULT, make the value of the target address in the bind call zero. Then use the sendto call to send data and set the target address in that call. Use the recvfrom call to determine the source address of any data.

## **E.4.9** Deactivating the Socket

When you have finished sending or receiving data, deactivate the socket by issuing the close system call.

# **E.5 DLI Programming Examples**

This section includes the following DLI programming examples:

- A sample DLI client program using Ethernet format packets
- A sample DLI server program using Ethernet format packets
- A sample DLI client program using 802.3 format packets
- A sample DLI server program using 802.3 format packets
- A sample DLI program using getsockopt and setsockopt system calls

These programming examples are also available on line in the /usr/examples/dli directory.

## E.5.1 Sample DLI Client Program Using Ethernet Format Packets

```
#include <stdio.h>
#include <errno.h>
#include <string.h>
#include <memory.h>
#include <stdlib.h>
#include <unistd.h>
#include <sys/types.h>
#include <sys/socket.h>
#include <sys/ioctl.h>
#include <net/if.h>
#include <net/route.h>
#include <dli/dli_var.h>
       dli_example:dli_eth
  Description: This program sends out a message to a node where a
                companion program, dli_ethd, echoes the message back.
                The ethernet packet format is used. The ethernet
                address of the node where the companion program is
                running, the protocol type, and the message are
                supplied by the user. The companion program should
                be started before executing this program.
 * Inputs:
               device, target address, protocol type, short message.
 * Outputs:
               Exit status.
* To compile: cc -o dli_eth dli_eth.c
* Example:
                dli_eth ln0 08-00-2b-02-e2-ff 6006 "Echo this"
 * Comments:
                This example demonstrates the use of the "NORMAL" \ensuremath{\text{I/O}}
                control flag. The use of the "NORMAL" flag means that
                we can communicate only with a single specific node
                whose address is specified during the bind. Because
                of this, we can use the normal write \&\ \text{read} system
                calls on the socket, because the source/destination of
                all data that is read/written on the socket is fixed.
* Digital Equipment Corporation supplies this software example on
* an "as-is" basis for general customer use. Note that Digital
* does not offer any support for it, nor is it covered under any
* of Digital's support contracts.
* /
main(
   int argc,
   char **argv)
   struct sockaddr_dl sdl;
   size_t sdllen;
   int ch, fd, rsize, itarget[6], ptype, ioctlflg = DLI_NORMAL, errflg = 0;
   u_char inbuf[4800], u_char *src;
```

```
memset(&sdl, 0, sizeof(sdl));
while ((ch = getopt(argc, argv, "xp:")) != EOF) {
 case 'x': ioctlflg = DLI_EXCLUSIVE; break;
 case 'p': {
     if (sscanf(optarg, "%x", &ptype, &ch) != 1) {
  fprintf(stderr, "%s: invalid protocol type "s
             argv[0], optarg);
        errflg++;
       break;
      }
 default:
              errflg++; break;
if (errflg || argc - optind < 5) {
    fprintf(stderr, "%s %s %s\n",</pre>
        "usage:",
        argv[0],
        "device lan-address short-message");
    exit(1);
}
 * Get device name and unit number.
if (sscanf(argv[optind], "%[a-z]%hd%c", sdl.dli_device.dli_devname,
    &sdl.dli_device.dli_devnumber, &ch) != 2) {
  fprintf(stderr, "%s: invalid device name
       argv[0], argv[optind]);
 exit(1);
}
 * Get the address to which we will be sending
if (sscanf(argv[++optind], "%x%*[:-]%x%*[:-]%x%*[:-]
                              %x%*[:-]%x%*[:-]%x%c",
     &itarget[0], &itarget[1], &itarget[2],
     &itarget[3], &itarget[4], &itarget[5], &ch) != 6) {
  fprintf(stderr, "%s: invalid lan address
       argv[0], argv[optind]);
 exit(1);
       If the LAN Address is a multicast, then we can't
  use DLI_NORMAL. Use DLI_DEFAULT instead.
* /
if ((itarget[0] & 1) && ioctflg == DLI_NORMAL)
 ioctlflg = DLI_DEFAULT;
* fill out sockaddr structure for bind/sento/recvfrom
sdl.dli_family = AF_DLI;
if (ptype < GLOBAL_SAP) {
 sdl.dli_substructype = DLI_802;
 sdl.choose_addr.dli_802addr.ioctl = ioctlflg;
 sdl.choose_addr.dli_802addr.svc = TYPE1;
 sdl.choose_addr.dli_802addr.eh_802.dsap = ptype;
```

```
sdl.choose_addr.dli_802addr.eh_802.ssap = ptype;
 sdl.choose_addr.dli_802addr.eh_802.ctl.U_fmt = UI_NPCMD;
 src = sdl.choose_addr.dli_802addr.eh_802.dst;
 sdl.dli_substructype = DLI_ETHERNET;
 sdl.choose_addr.dli_eaddr.dli_ioctlflg = ioctlflg;
 sdl.choose_addr.dli_eaddr.dli_protype = ptype;
 src = sdl.choose_addr.dli_eaddr.dli_target;
\mbox{*} If we are using DLI_NORMAL, we must supply
if (ioctlflg == DLI_NORMAL) {
 src[0] = itarget[0]; src[1] = itarget[1]; src[2] = itarget[2];
 src[3] = itarget[3]; src[4] = itarget[4]; src[5] = itarget[5];
* Open a socket to DLI and then bind to our protocol/address.
if ((fd = socket(AF_DLI, SOCK_DGRAM, DLPROTO_DLI)) < 0) {</pre>
 fprintf(stderr, "%s: DLI open failed: %s\n",
       argv[0], strerror(errno));
 exit(1);
}
if (bind(fd, (struct sockaddr *) &sdl, sizeof(sdl)) < 0) {</pre>
 fprintf(stderr, "%s: DLI bind failed: %s\n",
       argv[0], strerror(errno));
 exit(2);
if (ioctlflg != DLI_NORMAL) {
 src[0] = itarget[0]; src[1] = itarget[1]; src[2] = itarget[2];
 src[3] = itarget[3]; src[4] = itarget[4]; src[5] = itarget[5];
/* send response to originator. */
sdllen = sizeof(sdl);
if (sendto(fd, argv[4], strlen(argv[4]), 0,
        (struct sockaddr *) &sdl, sdllen) < 0) {
 fprintf(stderr, "%s: DLI transmission failed: %s\n",
       argv[0], strerror(errno));
 exit(1);
if ((rsize = recvfrom(fd, inbuf, sizeof(inbuf), 0,
              (struct sockaddr *) &sdl, &sdllen)) < 0 ) {
 fprintf(stderr, "%s: DLI reception failed: %s\n",
       argv[0], strerror(errno));
 exit(1);
/* check header */
if (sdllen != sizeof(struct sockaddr_dl)) {
 fprintf(stderr, "%s, incorrect header supplied\n", argv[0]);
 exit(1);
}
if (from.dli_substructype == DLI_802)
```

```
src = from.dli_choose_addr.dli_802addr.eh_802.dst;
else
 src = from.dli_choose_addr.dli_eaddr.dli_target;
/* anv data? */
fprintf(stderr, "%s: %sdata received from ", argv[0],
       rsize ? "" : "NO ");
fprintf(stderr, "%02x-%02x-%02x-%02x-%02x",
     src[0], src[1], src[2], src[3], src[4], src[5]);
if (from.dli_substructype == DLI_802)
 fprintf(stderr, "SAP %02x\n\n",
       sdl.choose_addr.dli_802addr.eh_802.ssap & ~1);
 fprintf(stderr, " on protocol type 04x\n\n",
       sdl.choose_addr.dli_eaddr.dli_protype);
/* print results */
printf("%s\n", inbuf);
close(fd);
return 0;
```

## E.5.2 Sample DLI Server Program Using Ethernet Format Packets

```
#ifndef lint
static char *rcsid = "@(#)$RCSfile: netprog.ap-dli,v $ \
        $Revision: 1.1.8.7 $ (DEC) $Date: 1996/02/01 21:37:40 $";
#endif
#include <stdio.h>
#include <ctype.h>
#include <errno.h>
#include <strings.h>
#include <sys/types.h>
#include <sys/socket.h>
#include <net/if.h>
#include <netinet/in.h>
#include <netinet/if_ether.h>
#include <dli/dli var.h>
#include <sys/ioctl.h>
extern int errno;
       dli_example:dli_ethd
* Description: This daemon program transmits any message it
 * receives to the originating system, i.e., it echoes the
^{\star} \, message back. The device and protocol type are supplied
  by the user. The program uses ethernet format packets.
* Inputs:
               device, protocol type.
* Outputs:
               Exit status.
* To compile: cc -o dli_ethd dli_ethd.c
```

```
* Example:
                dli_ethd de0 6006
* Comments:
                This example demonstrates the use of the "DEFAULT"
* I/O control flag, and the recvfrom & sendto system calls.
* By specifying "DEFAULT" when binding the DLI socket to
 ^{\star} the device we inform the system that this program will
 * receive any ethernet format packet with the given
 * protocol type which is not meant for any other program
   on the system. Since packets may arrive from
 * different systems we use the recvfrom call to read the
 * packets. This call gives us access to the packet
* header information so that we can determine where the
 * packet came from. When we write on the socket we must
* use the sendto system call to explicitly give the
* destination of the packet.
*/
* Digital Equipment Corporation supplies this software
* example on an "as-is" basis for general customer use. Note
\mbox{\ensuremath{^{\star}}} that Digital does not offer any support for it, nor is it
* covered under any of Digital's support contracts.
main(argc, argv, envp)
int argc;
char **argv, **envp;
    u_char inbuf[1500], outbuf[1500];
    u_char devname[16];
    u_char target_eaddr[6];
    char *cp;
    int rsize;
    unsigned int devunit;
    int i, sock, fromlen;
    unsigned int ptype;
    struct sockaddr_dl from;
    if ( argc < 3 )
        fprintf(stderr,
                "usage: %s device hex-protocol-type\n", argv[0]);
        exit(1);
    /* get device name and unit number. */
    bzero(devname, sizeof(devname));
    i = 0;
    cp = argv[1];
    while ( isalpha(*cp) )
       devname[i++] = *cp++;
    sscanf(cp, "%d", &devunit);
    /* get protocol type */
    sscanf(argv[2], "%x", &ptype);
    /* open dli socket */
    if
```

```
DLI_DEFAULT))<0)</pre>
       perror("dli_ethd, dli_econn failed");
       exit(1);
   while ( 1 ) {
       /* wait for message */
       from.dli_family = AF_DLI;
       fromlen = sizeof(struct sockaddr_dl);
       if ((rsize = recvfrom(sock, inbuf, sizeof(inbuf),
                           NULL, &from, &fromlen)) < 0 ) \{
           sprintf(inbuf, "%s: DLI reception failed", argv[0]);
           perror(inbuf);
           exit(2);
       }
       /* check header */
       if ( fromlen != sizeof(struct sockaddr_dl) ) {
           fprintf(stderr,"%s, incorrect header supplied\n",argv[0]);
           continue;
       }
       /* any data? */
       if (! rsize)
           fprintf(stderr, "%s, NO data received from ", argv[0]);
       else
       fprintf(stderr, "%x%s",
                   from.choose_addr.dli_eaddr.dli_target[i],
                   ((i<5)?"-":" "));
       fprintf(stderr, "on protocol type %x\n",
               from.choose_addr.dli_eaddr.dli_protype);
       /* send response to originator. */
       if ( sendto(sock, inbuf, rsize, NULL, &from, fromlen) < 0 ) {
           sprintf(outbuf, "%s: DLI transmission failed", argv[0]);
           perror(outbuf);
           exit(2);
       }
   }
}
               dli_econn
* Description:
       This subroutine opens a dli socket, then binds an associated
       device name and protocol type to the socket.
 * Inputs:
                      = ptr to device name
       devname
       devunit
                      = device unit number
                      = protocol type
       ptype
       taddr
                      = target address
       ioctl
                      = io control flag
```

((sock = dli\_econn(devname, devunit, ptype, NULL, \

```
* Outputs:
                       = socket handle if success, otherwise -1
 */
dli_econn(devname, devunit, ptype, taddr, ioctl)
char *devname;
unsigned devunit;
unsigned ptype;
u_char *taddr;
u_char ioctl;
    int i, sock;
    struct sockaddr_dl out_bind;
    if ( (i = strlen(devname)) >
          sizeof(out_bind.dli_device.dli_devname) )
         fprintf(stderr, "dli_ethd: bad device name");
         return(-1);
    if ((sock = socket(AF_DLI, SOCK_DGRAM, DLPROTO_DLI)) < 0)</pre>
         perror("dli_ethd, can't open DLI socket");
         return(-1);
    * fill out bind structure
    * /
    bzero(&out_bind, sizeof(out_bind));
    out_bind.dli_family = AF_DLI;
    out_bind.dli_substructype = DLI_ETHERNET;
    bcopy(devname, out_bind.dli_device.dli_devname, i);
    out_bind.dli_device.dli_devnumber = devunit;
    out_bind.choose_addr.dli_eaddr.dli_ioctlflg = ioctl;
    out_bind.choose_addr.dli_eaddr.dli_protype = ptype;
    if ( taddr )
         bcopy(taddr, out_bind.choose_addr.dli_eaddr.dli_target,
               DLI_EADDRSIZE);
    if ( bind(sock, &out_bind, sizeof(out_bind)) < 0 )</pre>
         perror("dli_ethd, can't bind DLI socket");
         return(-1);
    return(sock);
```

# E.5.3 Sample DLI Client Program Using 802.3 Format Packets

```
#ifndef lint
static char
             *sccsid = "@(#)dli_802.c 1.1 (DEC OSF/1) 5/29/92";
#endif lint
#include <stdio.h>
#include <ctype.h>
#include <errno.h>
#include <strings.h>
#include <sys/types.h>
#include <sys/socket.h>
#include <net/if.h>
#include <netinet/in.h>
#include <netinet/if_ether.h>
#include <dli/dli_var.h>
#include <sys/ioctl.h>
extern int errno;
#define PROTOCOL_ID
                       \{0x00, 0x00, 0x00, 0x00, 0x5\}
u_char protocolid[] = PROTOCOL_ID;
       dli_example:dli_802
* Description: This program sends out a message to a system
     where a companion program, dli_802d, echoes the message
     back. The 802.3 packet format is used. The ethernet
     address of the system where the companion program is
     running, the sap, and the message are supplied by the
     user. The companion program should be started before
     executing this program.
* Inputs:
               device, target address, sap, short message.
* Outputs:
               Exit status.
#ifndef lint
static char *rcsid = "@(#)$RCSfile: netprog.ap-dli,v $ \
           $Revision: 1.1.8.7 $ (DEC) $Date: 1996/02/01 21:37:40 $";
#endif
#include <stdio.h>
#include <ctype.h>
#include <errno.h>
#include <strings.h>
#include <sys/types.h>
#include <sys/socket.h>
#include <net/if.h>
#include <netinet/in.h>
#include <netinet/if_ether.h>
#include <dli/dli_var.h>
#include <sys/ioctl.h>
extern int errno;
```

```
#define PROTOCOL_ID
                        \{0x00, 0x00, 0x00, 0x00, 0x5\}
u_char protocolid[] = PROTOCOL_ID;
       dli_example:dli_802
 \mbox{\scriptsize \star} Description: This program sends out a message to a system
      where a companion program, dli_802d, echoes the message
      back. The 802.3 packet format is used. The ethernet
      address of the system where the companion program is
      running, the sap, and the message are supplied by the
      user. The companion program should be started before
      executing this program.
 * Inputs:
                device, target address, sap, short message.
 * Outputs:
               Exit status.
 * To compile: cc -o dli_802 dli_802.c
* Example:
               dli_802 ge0 08-00-2b-02-e2-ff ac "Echo this"
 * Comments:
               This example demonstrates the use of 802 "TYPE1"
      service. With TYPE1 service, the processing of
      XID and TEST messages is handled transparently by
     DLI, i.e., this program doesn't have to be concerned
      with handling them. If the SNAP SAP (0xAA) is
      selected, a 5 byte protocol id is also required.
     This example automatically uses a protocol id of
      of PROTOCOL_ID when the SNAP SAP is used. Also,
     note the use of DLI_NORMAL for the i/o control flag.
      DLI makes use of this only when that SNAP_SAP/Protocol
     ID pair is used. DLI will filter all incoming messages
     by comparing the Ethernet source address and Protocol
      ID against the target address and Protocol ID set up
     in the bind call. Only if a match occurs will DLI
     pass the message up to the application.
* Digital Equipment Corporation supplies this software
* example on an "as-is" basis for general customer use. Note
 * that Digital does not offer any support for it, nor is it
* covered under any of Digital's support contracts.
main(argc, argv, envp)
int argc;
char **argv, **envp;
   u_char inbuf[1500], outbuf[1500];
   u_char target_eaddr[6];
   u_char devname[16];
   int rsize, devunit;
   char *cp;
   int i, sock, fromlen;
   struct sockaddr_dl from;
```

```
unsigned int obsiz, byteval;
u_int sap;
u_char *pi = 0;
if ( argc < 5 )
    fprintf(stderr, "%s %s %s\n",
            "usage:",
            argv[0],
            "device ethernet-address hex-sap short-message");
   exit(1);
}
/* get device name and unit number. */
bzero(devname, sizeof(devname));
i = 0;
cp = argv[1];
while ( isalpha(*cp) )
   devname[i++] = *cp++;
sscanf(cp, "%d", &devunit);
/* get phys addr of remote system */
bzero(target_eaddr, sizeof(target_eaddr));
i = 0;
cp = argv[2];
while ( *cp ) {
   if ( *cp == '-' ) {
       cp++;
       continue;
    else {
      sscanf(cp, "%2x", &byteval );
     target_eaddr[i++] = byteval;
       cp += 2;
   }
/* get sap */
sscanf(argv[3], "%x", &sap);
/* get message */
bzero(outbuf, sizeof(outbuf));
if ( (obsiz = strlen(argv[4])) > 1500 ) {
    fprintf(stderr, "%s: message is too long\n", argv[0]);
    exit(2);
strcpy(outbuf, argv[4]);
/* open dli socket. notice that if (and only if) the */
/* snap sap was selected then a protocol id must also */
/* be provided. */
if ( sap == SNAP_SAP )
   pi = protocolid;
if ( (sock = dli_802_3_conn(devname, devunit, pi, target_eaddr,
              DLI_NORMAL, TYPE1, sap, sap, UI_NPCMD)) < 0 ) {</pre>
   perror("dli_802, dli_econn failed");
    exit(3);
```

```
}
    /* send message to target. minimum message size is 46 bytes. */
   if ( write(sock, outbuf, (obsiz < 46 ? 46 : obsiz)) < 0 ) \{
        sprintf(outbuf, "%s: DLI transmission failed", argv[0]);
       perror(outbuf);
        exit(4);
    /* wait for response from correct address */
   while (1) {
       bzero(&from, sizeof(from));
        from.dli_family = AF_DLI;
        fromlen = sizeof(struct sockaddr_dl);
        if ((rsize = recvfrom(sock, inbuf, sizeof(inbuf),
                            NULL, &from, &fromlen)) < 0 ) {
              sprintf(inbuf, "%s: DLI reception failed", argv[0]);
              perror(inbuf);
              exit(5);
        if ( fromlen != sizeof(struct sockaddr_dl) ) {
              fprintf(stderr,"%s, invalid address size\n",argv[0]);
        if ( bcmp(from.choose_addr.dli_802addr.eh_802.dst,
                 target_eaddr, sizeof(target_eaddr)) == 0 )
   }
   if (! rsize ) {
        fprintf(stderr, "%s, no data returned\n", argv[0]);
        exit(7);
    /* print message */
   printf("%s\n", inbuf);
   close(sock);
}
               d 1 i _8 0 2 _ 3 _ c o n n
 * Description:
       This subroutine opens a dli 802.3 socket, then binds an
       associated device name and protocol type to the socket.
 * Inputs:
       devname
                       = ptr to device name
       devunit
                       = device unit number
       ptype
                       = protocol type
       taddr
                      = target address
       ioctl
                       = io control flag
       svc
                       = service class
       sap
                       = source sap
       dsap
                       = destination sap
        ctl
                       = control field
```

```
* Outputs:
                       = socket handle if success, otherwise -1
 * /
dli_802_3_conn (devname, devunit, ptype, taddr, ioctl, svc, sap, dsap, ctl)
char *devname;
u_short devunit;
u_char *ptype;
u_char *taddr;
u_char ioctl;
u_char svc;
u_char sap;
u_char dsap;
u_short ctl;
    int i, sock;
    struct sockaddr_dl out_bind;
    if ( (i = strlen(devname)) >
         sizeof(out_bind.dli_device.dli_devname) )
         fprintf(stderr, "dli_802: bad device name");
         return(-1);
    if ((sock = socket(AF_DLI, SOCK_DGRAM, DLPROTO_DLI)) < 0)</pre>
         perror("dli_802, can't open DLI socket");
         return(-1);
     * fill out bind structure. note that we need to determine
     * whether the ctl field is 8 bits (unnumbered format) or
     * 16 bits (informational/supervisory format). We do this
     * by checking the low order 2 bits, which are both 1 only
     * for unnumbered control fields.
    bzero(&out_bind, sizeof(out_bind));
    out_bind.dli_family = AF_DLI;
    out_bind.dli_substructype = DLI_802;
    bcopy(devname, out_bind.dli_device.dli_devname, i);
    out_bind.dli_device.dli_devnumber = devunit;
    out_bind.choose_addr.dli_802addr.ioctl = ioctl;
    out_bind.choose_addr.dli_802addr.svc = svc;
    if(ctl & 3)
        out_bind.choose_addr.dli_802addr.eh_802.ctl.U_fmt=\
              (u_char)ctl;
    else
        out_bind.choose_addr.dli_802addr.eh_802.ctl.I_S_fmt = \
              ctl;
    out_bind.choose_addr.dli_802addr.eh_802.ssap = sap;
    out_bind.choose_addr.dli_802addr.eh_802.dsap = dsap;
    if ( ptype )
        bcopy(ptype,out_bind.choose_addr.dli_802addr.eh_802.osi_pi,\
```

# E.5.4 Sample DLI Server Program Using 802.3 Format Packets

```
#ifndef lint
static char *rcsid = "@(#)$RCSfile: netprog.ap-dli,v $ \
           $Revision: 1.1.8.7 $ (DEC) $Date: 1996/02/01 21:37:40 $";
#endif
#include <stdio.h>
#include <ctype.h>
#include <errno.h>
#include <strings.h>
#include <sys/types.h>
#include <sys/socket.h>
#include <net/if.h>
#include <netinet/in.h>
#include <netinet/if_ether.h>
#include <dli/dli_var.h>
#include <sys/ioctl.h>
extern int errno;
#define PROTOCOL_ID
                       \{0x00, 0x00, 0x00, 0x00, 0x5\}
u_char protocolid[] = PROTOCOL_ID;
       dli_example:dli_802d
* Description: This daemon program transmits any message it
    receives to the originating system, i.e., it echoes the
    message back. The device and sap are supplied by the
    user. The program uses 802.3 format packets.
 * Inputs:
             device, sap.
* Outputs:
               Exit status.
* To compile: cc -o dli_802d dli_802d.c
* Example:
               dli_802d de0 ac
* Comments: This example demonstrates the recvfrom & sendto
    system calls. Since packets may arrive from different
```

```
This call gives us access to the packet header information
     so that we can determine where the packet came from.
     When we write on the socket we must use the sendto
     system call to explicitly give the destination of
     the packet. The use of the "DEFAULT" I/O control flag
     only applies (i.e. only has an affect) when the SNAP SAP
     is used. When the SNAP SAP is used, any arriving packets
     which have the specified protocol id and which are not
    destined for some other program will be given to this
 * Digital Equipment Corporation supplies this software
 * example on an "as-is" basis for general customer use.
 * Note that Digital does not offer any support for it, nor
 * is it covered under any of Digital's support contracts.
main(argc, argv, envp)
int argc;
char **argv, **envp;
    u_char inbuf[1500], outbuf[1500];
    u_char devname[16];
    u_char target_eaddr[6];
    char *cp;
    int rsize, devunit;
    int i, sock, fromlen;
    u_char tmpsap, sap;
    struct sockaddr_dl from;
    u_char *pi = 0;
    if ( argc < 3 )
        fprintf(stderr, "usage: %s device hex-sap\n", argv[0]);
        exit(1);
    /* get device name and unit number. */
    bzero(devname, sizeof(devname));
    i = 0;
    cp = argv[1];
    while ( isalpha(*cp) )
        devname[i++] = *cp++;
    sscanf(cp, "%d", &devunit);
    /* get sap */
    sscanf(argv[2], "%x", &sap);
    /* open dli socket. note that if (and only if) the snap sap */
    /* was selected then a protocol id must also be specified. */
    if ( sap == SNAP_SAP )
        pi = protocolid;
    if ((sock = dli_802_3_conn(devname, devunit, pi, target_eaddr,
                    DLI_DEFAULT, TYPE1, sap, sap, UI_NPCMD)) < 0) {</pre>
```

systems we use the recvfrom call to read the packets.

```
perror("dli_802d, dli_conn failed");
    exit(1);
/* listen and respond */
while ( 1 ) {
   /* wait for message */
    from.dli_family = AF_DLI;
    fromlen = sizeof(struct sockaddr_dl);
    if ((rsize = recvfrom(sock, inbuf, sizeof(inbuf), NULL,
                         &from, &fromlen)) < 0 ) {
        sprintf(inbuf, "%s: DLI reception failed", argv[0]);
        perror(inbuf);
        exit(2);
    }
    /* check header */
    if ( fromlen != sizeof(struct sockaddr_dl) ) {
        fprintf(stderr, "%s, incorrect header supplied\n", \
            argv[0]);
        continue;
    }
    ^{\star} Note that DLI swaps the source & destination saps and
    * lan addresses in the sockaddr_dl structure returned
     * by the recvfrom call. That is, it places the DSAP in
     * eh_802.ssap and the SSAP in eh_802.dsap; it also places
     * the destination lan address in eh_802.src and the source
     * lan address in eh_802.dst. This allows for minimal to
     * no manipulation of the address structure for subsequent
     * sendto or dli connection calls.
    /* any data? */
    if (! rsize)
        fprintf(stderr, "%s: NO data received from ", \
            argv[0]);
    else
       fprintf(stderr, "%s: data received from ", argv[0]);
    for (i = 0; i < 6; i++)
        fprintf(stderr, "%x%s",
                from.choose_addr.dli_802addr.eh_802.dst[i],
                ((i<5)?"-":" "));
    fprintf(stderr, "\n on dsap %x ",
            from.choose_addr.dli_802addr.eh_802.ssap);
    if (from.choose_addr.dli_802addr.eh_802.dsap == \
            SNAP_SAP )
        fprintf(stderr,
           "(SNAP SAP), protocol id = x-x-x-x-xn
           from.choose_addr.dli_802addr.eh_802.osi_pi[0],
           from.choose_addr.dli_802addr.eh_802.osi_pi[1],
           from.choose_addr.dli_802addr.eh_802.osi_pi[2],
           from.choose_addr.dli_802addr.eh_802.osi_pi[3],
           from.choose_addr.dli_802addr.eh_802.osi_pi[4]);
    fprintf(stderr, " from ssap %x ",
            from.choose_addr.dli_802addr.eh_802.dsap);
    fprintf(stderr, "\n\n");
    /* send response to originator. */
```

```
SNAP_SAP )
            bcopy(protocolid,
                  from.choose_addr.dli_802addr.eh_802.osi_pi, 5);
        if ( sendto(sock, inbuf, rsize, NULL, &from, fromlen) \setminus
                 < 0 ) {
            sprintf(outbuf, "%s: DLI transmission failed", \
                 argv[0]);
            perror(outbuf);
            exit(2);
    }
}
                d 1 i _8 0 2 _ 3 _ c o n n
  Description:
        This subroutine opens a dli 802.3 socket, then binds an
        associated device name and protocol type to the socket.
 * Inputs:
                       = ptr to device name
        devname
        devunit
                        = device unit number
        ptype
                        = protocol type
        taddr
                        = target address
        ioctl
                        = io control flag
        svc
                        = service class
        sap
                        = source sap
        dsap
                        = destination sap
        ctl
                        = control field
 * Outputs:
                        = socket handle if success, otherwise -1
       returns
dli_802_3_conn (devname,devunit,ptype,taddr,ioctl,svc,sap,\
        dsap,ctl)
char *devname;
u_short devunit;
u_char *ptype;
u_char *taddr;
u_char ioctl;
u_char svc;
u_char sap;
u_char dsap;
u_short ctl;
    int i, sock;
    struct sockaddr_dl out_bind;
    if ( (i = strlen(devname)) >
         sizeof(out_bind.dli_device.dli_devname) )
         fprintf(stderr, "dli_802d: bad device name");
```

if ( from.choose\_addr.dli\_802addr.eh\_802.dsap == \

```
return(-1);
}
if ((sock = socket(AF_DLI, SOCK_DGRAM, DLPROTO_DLI)) < 0)</pre>
     perror("dli_802d, can't open DLI socket");
    return(-1);
\mbox{*} fill out bind structure. note that we need to determine
* whether the ctl field is 8 bits (unnumbered format) or
 \star 16 bits (informational/supervisory format). We do this
 * by checking the low order 2 bits, which are both 1 only
 * for unnumbered control fields.
bzero(&out_bind, sizeof(out_bind));
out_bind.dli_family = AF_DLI;
out_bind.dli_substructype = DLI_802;
bzero(&out_bind, sizeof(out_bind));
out_bind.dli_family = AF_DLI;
out_bind.dli_substructype = DLI_802;
bcopy(devname, out_bind.dli_device.dli_devname, i);
out_bind.dli_device.dli_devnumber = devunit;
out_bind.choose_addr.dli_802addr.ioctl = ioctl;
out_bind.choose_addr.dli_802addr.svc = svc;
if(ctl & 3)
    out_bind.choose_addr.dli_802addr.eh_802.ctl.U_fmt=\
        (u_char)ctl;
else
    out_bind.choose_addr.dli_802addr.eh_802.ctl.I_S_fmt = \
        ctl;
out_bind.choose_addr.dli_802addr.eh_802.ssap = sap;
out_bind.choose_addr.dli_802addr.eh_802.dsap = dsap;
if (ptype)
    bcopy(ptype,out_bind.choose_addr.dli_802addr.eh_802.osi_pi,\
        5);
if (taddr)
     bcopy(taddr, out_bind.choose_addr.dli_802addr.eh_802.dst,
           DLI_EADDRSIZE);
if ( bind(sock, &out_bind, sizeof(out_bind)) < 0 )</pre>
    perror("dli_802d, can't bind DLI socket");
     return(-1);
return(sock);
```

# E.5.5 Sample DLI Program Using getsockopt and setsockopt

```
#ifndef lint
             *sccsid = "@(#)dli_setsockopt.c 1.5 3/27/90";
static char
#endif lint
#include <stdio.h>
#include <ctype.h>
#include <errno.h>
#include <strings.h>
#include <sys/types.h>
#include <sys/socket.h>
#include <net/if.h>
#include <netinet/in.h>
#include <netinet/if_ether.h>
#include <dli/dli_var.h>
#include <sys/ioctl.h>
extern int errno;
int debug = 0;
#define PROTOCOL_ID
                         \{0x00, 0x00, 0x00, 0x00, 0x5\}
#define CUSTOMER0
                         {0xab, 0x00, 0x04, 0x00, 0x00, 0x00}
#define CUSTOMER1
                         \{0xab, 0x00, 0x04, 0x00, 0x00, 0x01\}
u_char mcast0[] = CUSTOMER0;
u_char mcast1[] = CUSTOMER1;
u_char protocolid[] = PROTOCOL_ID;
       dli example: dli setsockopt
 * Description: This program demonstrates the use of the DLI
      get- and setsockopt calls. It opens a socket, enables
      2 multicast addresses, changes the 802 control
      field, enables a number of group saps supplied by
      the user, and reads the group saps that are enabled.
 * Inputs:
               device, sap, group-saps.
* Outputs:
               Exit status.
 * To compile: cc -o dli_setsockopt dli_setsockopt.c
* Example:
               dli setsockopt ge0 ac 5 9 d
               When a packet arrives with a group dsap,
      all dli programs that have that group sap enabled will
      receive copies of that packet. Group saps are
      those with the low order bit set. Group sap 1
      is currently not allowed for customer use. Group
      saps with the second bit set (eg 3,7,etc) are
      reserved by IEEE.
*/
* Digital Equipment Corporation supplies this software example
\mbox{\scriptsize *} on an "as-is" basis for general customer use. Note that
* Digital does not offer any support for it, nor is it covered
```

```
* under any of Digital's support contracts.
main(argc, argv, envp)
int argc;
char **argv, **envp;
   u_char inbuf[1500], outbuf[1500];
    u_char devname[16];
    u_char target_eaddr[6];
    char *cp;
   int rsize, devunit;
   int i, j, k, sock, fromlen;
    u_short obsiz;
   u_char tmpsap, sap;
    struct sockaddr_dl from;
    u_char *pi = 0;
    u_char out_opt[1000], in_opt[1000];
    int optlen, ioptlen = sizeof(in_opt);
    if (argc < 4)
        fprintf(stderr, "usage: %s device hex-sap hex-groupsaps\n",
        exit(1);
    /* get device name and unit number. */
    bzero(devname, sizeof(devname));
    i = 0;
    cp = argv[1];
    while ( isalpha(*cp) )
    devname[i++] = *cp++;
    sscanf(cp, "%d", &devunit);
    /* get protocol type */
    sscanf(argv[2], "%x", &sap);
    /* open dli socket */
    if ( sap == SNAP_SAP ) {
        fprintf(stderr,
               "%s: can't use SNAP_SAP in USER mode\n", argv[0]);
        exit(1);
    if ( (sock = dli_802_3_conn(devname, devunit, pi,\
           target_eaddr,
                     DLI_DEFAULT, USER, sap, sap, UI_NPCMD)) \
           < 0 ) {
        perror("dli_setsockopt: dli_conn failed");
        exit(1);
    /* enable two multicast addresses */
    bcopy(mcast0, out_opt, sizeof(mcast0));
    bcopy(mcast1, out_opt+sizeof(mcast0), sizeof(mcast1));
    if ( setsockopt(sock, DLPROTO_DLI, DLI_MULTICAST, \
           &out_opt[0],
```

```
(sizeof(mcast0) + sizeof(mcast1))) < 0)
    perror("dli_setsockopt: can't enable multicast");
/* set 802 control field */
out_opt[0] = TEST_PCMD;
optlen = 1;
i f
(setsockopt(sock,DLPROTO_DLI,DLI_SET802CTL,&out_opt[0],\)
       optlen)<0){
    perror("dli_setsockopt: Can't set 802 control");
    exit(1);
/* enable GSAPs supplied by user */
j = 3;
i = 0;
while (j < argc ) \{
   sscanf(argv[j++], "%x", &k);
    out_opt[i++] = k;
optlen = i;
if
(setsockopt(sock,DLPROTO_DLI,DLI_ENAGSAP,&out_opt[0],\
       optlen) < 0)
    perror("dli_setsockopt: Can't enable gsap");
    exit(1);
/* verify all gsaps are enabled */
bzero(in_opt, (ioptlen = sizeof(in_opt)));
(getsockopt(sock,DLPROTO_DLI,DLI_GETGSAP,in_opt,\
      &ioptlen) < 0){
    perror("dli_setsockopt: DLI getsockopt 2 failed");
    exit(1);
printf("number of enabled GSAPs = %d, GSAPS:", ioptlen);
for(i = 0; i < ioptlen; i++) {</pre>
    if (!(i%10))
       printf("\n");
    printf("%2x ",in_opt[i]);
printf("\n");
/* disable all but the last 4 or all GSAPs, */
/* whichever is smallest */
if (optlen > 4)
    optlen -= 4;
(setsockopt(sock,DLPROTO_DLI,DLI_DISGSAP,&out_opt[0],\)
       optlen) < 0){
    perror("dli_setsockopt: Can't disable gsap");
/* verify some gsaps still enabled */
bzero(in_opt, (ioptlen = sizeof(in_opt)));
(getsockopt(sock,DLPROTO_DLI,DLI_GETGSAP,in_opt,\
       \&ioptlen) < 0){
```

```
perror("dli_setsockopt: getsockopt 3 failed");
        exit(1);
    printf("number of enabled GSAPs = %d, GSAPS:", ioptlen);
    for(i = 0; i < ioptlen; i++) {</pre>
        if (!(i%10))
    printf("\n");
        printf("%2x ",in_opt[i]);
    printf("\n");
}
                d l i _8 0 2 _ 3 _ c o n n
 * Description:
        This subroutine opens a dli 802.3 socket and then binds
        an associated device name and protocol type to it.
 * Inputs:
                 = ptr to device name
        devname
        devunit = device unit number
        ptype
                  = protocol type
                  = target address
        taddr
        ioctl
                  = io control flag
        svc
                  = service class
        sap
                  = source sap
        dsap
                  = destination sap
        ctl
                   = control field
 * Outputs:
        returns
                   = socket handle if success, otherwise -1
dli_802_3_conn (devname,devunit,ptype,taddr,ioctl,svc,sap,\
       dsap,ctl)
char *devname;
u_short devunit;
u_char *ptype;
u_char *taddr;
u_char ioctl;
u_char svc;
u_char sap;
u_char dsap;
u_short ctl;
    int i, sock;
    struct sockaddr_dl out_bind;
    if ( (i = strlen(devname)) >
         sizeof(out_bind.dli_device.dli_devname) )
        fprintf(stderr, "dli_setsockopt: bad device name");
        return(-1);
    }
```

```
if ((sock = socket(AF_DLI, SOCK_DGRAM, DLPROTO_DLI)) < 0)</pre>
        perror("dli_setsockopt: can't open DLI socket");
        return(-1);
     \mbox{*} fill out bind structure
    * /
    bzero(&out_bind, sizeof(out_bind));
    out_bind.dli_family = AF_DLI;
    out_bind.dli_substructype = DLI_802;
    bcopy(devname, out_bind.dli_device.dli_devname, i);
    out_bind.dli_device.dli_devnumber = devunit;
    out_bind.choose_addr.dli_802addr.ioctl = ioctl;
    out_bind.choose_addr.dli_802addr.svc = svc;
    if(ctl & 3)
        out_bind.choose_addr.dli_802addr.eh_802.ctl.U_fmt=\
            (u_char)ctl;
    else
        out_bind.choose_addr.dli_802addr.eh_802.ctl.I_S_fmt = \
            ctl;
    out_bind.choose_addr.dli_802addr.eh_802.ssap = sap;
    out_bind.choose_addr.dli_802addr.eh_802.dsap = dsap;
    if ( ptype )
        bcopy(ptype,out_bind.choose_addr.dli_802addr.eh_802.osi_pi,\
           5);
    if ( taddr )
        bcopy(taddr, out_bind.choose_addr.dli_802addr.eh_802.dst,
              DLI_EADDRSIZE);
    if ( bind(sock, &out_bind, sizeof(out_bind)) < 0 )</pre>
        perror("dli_setsockopt: can't bind DLI socket");
        return(-1);
    return(sock);
}
```

# **Glossary**

### active user

In an XTI transport connection, the transport user that initiated the connection. See also **client process** and **passive user**.

### Address Resolution Protocol (ARP)

The Internet (TCP/IP) Protocol that can dynamically resolve an Internet address to a physical hardware address. ARP can be used only across a single physical network and in networks that support the hardware broadcast feature.

### asynchronous event

See event.

### asynchronous execution

- 1. Execution of processes or threads in which each process or thread does not await the completion of the others before starting.
- 2. In XTI, a mode of execution that notifies the transport user of an event without forcing it to wait.

# Berkeley Software Distribution

UNIX software release of the Computer Systems Research Group (CSRG) of the University of California at Berkeley.

# blocking mode

See synchronous execution.

### BSD socket interface

A transport-layer interface provided for applications to perform interprocess communication between two unrelated processes on a single system or on multiply connected systems. This interprocess communications facility allows programs to use sockets for communications between other programs, protocols, and devices.

### client process

In the client/server model of communication, a process that requests services from a server process. See also **active user**.

### communication domain

An abstraction used by the interprocess communication facility of a system to define the properties of a network. Properties include a set of communication protocols, rules for manipulating and interpreting names, and the ability to transmit access rights.

### connection-oriented mode

A mode of service supported by a transport endpoint for transmitting data over an established connection.

#### connectionless mode

A mode of service supported by a transport endpoint that requires no established connection for transmitting data. Data is delivered in self-contained units, called **datagrams**.

### datagram

A unit of data that is transmitted across a network by the connectionless service of a transport provider. In addition to user data, a datagram includes the information needed for its delivery. It is self-contained, in that it has no relationship to any datagrams previously or successively transmitted.

# datagram socket

Socket that provides datagrams consisting of individual messages for transmission in connectionless mode.

### error

In XTI, an indicator that is returned by a function when it encounters a system or library error in the process of executing. The object is to allow applications to take an action based on the returned error code.

#### **eSNMP**

The Extensible Simple Network Protocol (eSNMP) enables you to create subagents to be managed by an SNMP management station. See Chapter 6.

#### Ethernet

A 10-megabit baseband local area network (LAN) using carrier sense multiple access with collision detection (CSMA/CD). The network

allows multiple stations to access the medium at will without prior coordination, and avoids contention by using carrier sense and deference, and detection and transmission.

#### **ETSDU**

# See Expedited Transport Service Data Unit and out-of-band data.

#### event

An occurrence, or happening, that is significant to a transport user. Events are asynchronous, in that they do not happen as a result of an action taken by the user.

# event management

A mechanism by which transport providers notify transport users of the occurrence of significant events.

### expedited data

Data that is considered urgent. The semantics of this data are defined by the transport provider. See also **out-of-band data**.

### Expedited Transport Service Data Unit

In XTI, an expedited message in which the identity of the data unit is preserved from one end of a transport connection to the other.

# file descriptor

A small unsigned integer that a UNIX system uses to identify a file. A file descriptor is created by a process through issuing an open system call for the file name. A file descriptor ceases to exist when it is no longer held by any process.

### host group

A group of zero or more hosts that, for the purposes of IP multicasting, are identified by a single class D IP destination address. Class D IP addresses have 1110 as their high-order four bits. See **IP Multicasting** for more information.

# **ICMP**

# See Internet Control Message Protocol.

### #include file.h

A C language precompiler directive specifying interpolation of a named file into the file being compiled. The interpolated file is a standard

header file (indicated by placing its name in angle brackets) or any other file containing C language code (indicated by placing its name in double quotation marks).

The absolute path name of header files whose names are placed in angle brackets (<>) is /usr/include/file.h.

# International Standards Organization (ISO)

An international body composed of the national standards organizations of 89 countries. ISO issues standards on a vast number of goods and services, including networking software.

### Internet Control Message Protocol (ICMP)

A host-to-host protocol from the Internet Protocol (IP) suite that provides error and informational messages on the operations of the IP.

# Internet Protocol (IP)

The Internet Protocol that provides a connectionless service for the delivery of datagrams across a network.

### ISO

See International Standards Organization.

# **IP Multicasting**

IP Multicasting is a method for transmitting IP datagrams to a group of hosts identified by a single IP destination address, or host group. Host groups are identified by class D IP addresses. See **host group** for more information.

# Management Information Base See MIB.

### MIB

The Management Information Base (MIB) definitions are a set of data elements that relate to network management. See Chapter 6.

### name server

A daemon running on a system that client processes contact to obtain the addresses of hosts or other objects in a network. This daemon translates a machine's network name to its network IP address.

#### name service

The service provided to client processes for identifying peer systems for communications purposes.

# nonblocking mode

See asynchronous execution.

#### normal data

Regular data that is sent or received in band by a transport user. See also **out-of-band data**.

# Object Identifier

See OID.

### OID

Object Identifiers (OID) are data elements in MIB definitions that can be referred to by name or by a corresponding sequence of numbers. See Chapter 6.

# Open Systems Interconnection (OSI)

The interconnection of open systems in accordance with ISO standards.

# orderly release

In XTI, an optional feature that allows a transport user to gracefully terminate a transport connection with no loss of data.

### OSI

See Open Systems Interconnection.

# out-of-band data

Data that is transmitted out of the flow of normal data because it is considered urgent. The receiving process is notified of the presence of this data so that it can be retrieved.

### passive user

In an XTI transport connection, the peer transport user that responded to the connection request. See also **active user** and **client process**.

# pipe

An I/O stream that has a descriptor and can be used in unidirectional communications between related processes. See also **socketpair**.

### raw socket

A socket that provides privileged users access to internal network protocols and interfaces. These socket types can be used to take advantage of protocol features not available through more normal interfaces or to communicate with hardware interfaces.

### Serial Line Internet Protocol (SLIP)

A transmission line protocol that encapsulates and transfers IP datagrams over asynchronous serial lines.

# server process

In the client/server model of communication, a process that provides services to client processes. See also **passive user**.

#### SLIP

See Serial Line Internet Protocol.

#### socket

In interprocess communications, an endpoint of communication. Also, the system call that creates a socket and the associated data structure.

#### socketpair

A pair of sockets that can be created in the UNIX domain for 2-way communication. Like pipes, socketpairs require communicating processes to be related. See also **pipe**.

#### state

In XTI, the current condition of a process that reflects the function in progress. XTI uses eight different states to manage communications over a transport endpoint.

### stream socket

A socket that provides 2-way byte streams across a transport connection. Also includes a mechanism for handling out-of-band data.

# **STREAMS**

A kernel mechanism specified by AT&T that supports the implementation of device drivers and networking protocol stacks. See also **STREAMS framework**.

### STREAMS framework

Components of the AT&T STREAMS mechanism which define the

interface standards for character I/O within the kernel and between the kernel and user levels. It consists of functions, utility routines, kernel facilities, and data structures.

# synchronous execution

A mode of execution that forces transport primitives to wait for specific events before returning control to the transport user.

#### **TCP**

See Transmission Control Protocol.

# TCP/IP

The two fundamental protocols of the Internet Protocol suite, and an acronym that is frequently used to refer to the Internet Protocol suite. TCP provides for the reliable transfer of data, while IP transmits the data through the network in the form of datagrams. See also **Transmission Control Protocol** and **Internet Protocol**.

# TLI

# See Transport Layer Interface.

### Transmission Control Protocol (TCP)

The Internet transport-layer protocol that provides a reliable, full-duplex, connection-oriented service for applications. TCP uses the IP protocol to transmit information through the network.

# transport endpoint

A communication path over which a transport user can exchange data with a transport provider. See also **Transport Layer Interface**.

### Transport Layer Interface (TLI)

An interface to the transport layer of the OSI model, designed on the ISO Transport service definition.

### transport provider

A transport protocol that offers transport layer services in a system.

# Transport Service Data Unit (TSDU)

In OSI terminology, the item of information, or message, that the transport user passes to the transport provider.

# transport services

The support given by the transport layer in a system to the application layer for the transfer of data between user processes. The two types of services provided are connection-oriented and connectionless. See also **Transport Layer Interface**.

# transport user

A program needing the services of a transport protocol to send data to or receive data from another program or point in a network. See also **Transport Layer Interface**.

### **TSDU**

See Transport Service Data Unit.

#### UDP

See User Datagram Protocol.

# User Datagram Protocol (UDP)

The Internet Protocol that allows application programs on remote machines to send datagrams to one another. UDP uses IP to deliver the datagrams.

# X/Open Transport Interface

Protocol-independent, transport-layer interface for applications. XTI consists of a series of C language functions based on TLI, which in turn was based on the transport service definition for the OSI model.

### XTI

See X/Open Transport Interface.

# Index

Special Characters	accept2 event, 3–15
802.3 frame format	access rights
description of, E-11	and the recvmsg system call, 4-33
example of, E–4	and the sendmsg system call, 4-33
processing, E–11	acknowledged connectionless mode of
802.3 substructure	communication
filling the, E–21	in DLPI, 2–4
802.3 substructure values	acknowledged connectionless mode service
control field, E-13	in DLPI, 2–6
destination service access point, E–12	active user
destination system physical address, E–11	defined, 1, 3–3
exchange identification, E–14	typical state transitions, 3-21f
LLC Protocol Data Unit Test, E–15	address family
Service class, E–12	specifying for DLI, E-17
source service access point, E–13	address generation
Unnumbered Information Command, E–15	comparison of TLI and XTI, 3-44
XID, E-14	addressing in DLPI, 2–7
MD, E 11	identifying components, 2-7
A	PPA, 2–7
•	advanced sockets topics, 4-41 to 4-55
abortive release in XTI, 3–10t	AF_INET domain, 4–4
accept socket call	AF_UNIX
contrast to XTI t_accept function, 3-46	See UNIX domain
accept system call, 4–9t, 4–26	AF_UNIX domain, 4-4
accept1 event, 3–15t	See also UNIX domain
accept2 event, 3–15t	alignment
accept3 event, 3–15t	and the Routing Information Field, D-3

all hosts group	bind system call, 4–9t, 4–22, E–7
defined, 4–47	syntax, E–19
application programming interface	binding names to addresses, 4-42
sockets, 1-1, 4-6 to 4-40	in the UNIX domain, 4-44
STREAMS, 1-1, 5-5 to 5-25	INADDR_ANY wildcard address, 4-42
XTI, 1–1, 3–4 to 3–41	binding names to sockets, 4–22
application programs	blocking mode
porting to XTI, 3-41	See synchronous execution
rewriting for XTI, 3-45	bridging
sockets	BSD drivers to STREAMS protocol stacks,
and the netdb.h header file, 4-10	7–11
applications	STREAMS drivers to sockets protocol
distributed	stacks, 7–2
and the client/server paradigm, 4-8	broadcasting and determining network
asynchronous events in XTI, 3–10	configuration, 4–51
and consuming functions, 3-11t	BSD
T_CONNECT, 3-10t	sockets, 4–37
T_DATA, 3–10t	BSD drivers
T_DISCONNECT, 3-10t	bridging to STREAMS protocol stacks, 7-11
T_EXDATA, 3–10t	BSD socket interface
T_GODATA, 3–10t	binding names to sockets, 4-22
T_GOEXDATA, 3–10t	4.3BSD msghdr data structure, 4-39
T_LISTEN, 3–10t	4.4BSD msghdr data structure, 4-39
T_ORDREL, 3–10t	datagram sockets, 4-5
T_UDERR, 3–10t	establishing connections to sockets in, 4-23
asynchronous execution in XTI	performing blocking and nonblocking
defined, 3–5	operations in, 4–21
audience	raw sockets, 4–5
of manual, xvii	stream sockets, 4-5
	transferring data in, 4-29
В	using socket options in, 4-27
big-endian	buffer size
defined, 4–13	calculating, E–22
bind event, 3–15	
bind socket call	

contrast to XTI t\_bind function, 3-46

С	configuration processing, 5–21
canonical addresses	connect system call, 4-9t, 4-23
and Token Ring drivers, D-2	and TCP, 4-24
client process	and UDP, 4–24
defined, 4–8	connect1 event, 3–15t
establishing connections, 4–23	connect2 event, 3–15t
client/server interaction, 4–8	connection establishment phase
clone device, 5–30	state transitions allowed, 3-18t to 3-20t
close function, 5–7	connection indication
close processing, 5–21	in XTI, 3–10t
close socket call	connection mode
contrast to XTI t_snddis function, 3-47	of communication in DLPI, 2-3
close system call, 4–36	connection mode service
closed event, 3–15t	in DLPI, 2–5
closing sockets, 4–36	connection-oriented applications
CLTS	initializing an endpoint, 3–24
See connectionless service in XTI	writing, 3–24 to 3–36
coexistence	connection-oriented communication
defined for Digital UNIX, 7–1	and TCP, 4–7
of STREAMS and sockets, 7–1 to 7–12	sockets, 4–7
communication bridge	<b>connection-oriented programs</b> , B-2 to B-17
defined, 7–1	connection-oriented service in XTI
dlb STREAMS pseudodriver, 1–8f, 1–7, 7–1	defined, 3–4
ifnet STREAMS module, 1–7f, 1–7, 7–1	connection-oriented transport service
communication domains	state transitions allowed in XTI, 3-18t to
sockets, 4–3	3–20t
Internet domain, 4–4	typical sequence of functions, 3-21
UNIX domain, 4–4	connectionless applications
communication properties of sockets, 4–3	writing, 3–37 to 3–39
comparison	connectionless communication
of XTI and sockets, 3–45	and UDP, 4–8
of XTI and TLI, 3-43	sockets, 4–8
compatibility	connectionless mode of communication
of XTI and TLI, 3–43 to 3–44	in DLPI, 2–3
concurrent programs	connectionless mode service
running, E–2	in DLPI, 2–6

connectionless programs, B-17 to B-30	data structures
connectionless service in XTI	4.3BSD msghdr, 4–39
defined, 3–4	4.4BSD msghdr, 4–39
state transitions allowed, 3-17t	dblk_t, 5–19
typical state transitions, 3-22	hostent, 4–11
connections	mblk_t, 5–19
passing to another endpoint, 3-16	message, 5–18
consuming functions	module, 5–17
for asynchronous XTI events, 3-11t	module_info, 5–18
control field	qinit, 5–17
function of, E-13	streamtab, 5–18
COTS	msghdr, 4–17, 4–18
See connection-oriented transport service	netent, 4–11
	protoent, 4–12
D	servent, 4–12
daemon	sockaddr, 4–16
inetd, 4–54	sockaddr_in, 4-17
data flow	sockaddr_un, 4-17
XTI and a sockets-based transport provider,	data transfer
1–6	using DLI program, E-22
XTI and a STREAMS-based transport	data transfer phase
provider, 1–6	of connectionless service, 3-37
Data Link Interface	state transitions allowed for connectionless
See DLI	transport services, 3-17t
data link interfaces, 1–3, 2–1 to 2–12	data transfer state
DLPI, 2–1	in XTI, 3–13t
Data Link Provider Interface	data units
See DLPI	receiving, 3–38
data link service provider	receiving error information, 3-39
See DLS provider	datagram socket, 4–5, E–1, E–17
data link service providers in DLPI, 2–8	dblk_t data structure, 5–19
data link service user	destination service access point
See DLS user	See DSAP
data segmentation	destination system
providing, E-4, E-16	specifying information, E-8

destination system physical address	DL_TEST_IND primitive, 2–8t
defined, E-9, E-11, E-12	DL_TEST_REQ primitive, 2-8t
specifying, E-9	DL_TEST_RES primitive, 2–8t
device drivers	DL_UDERROR_IND primitive, 2-8t
and Stream ends, 5-4	DL_UNBIND_REQ primitive, 2-8t, 7-12
STREAMS processing routines for, 5-20	DL_UNIDATA_IND primitive, 2-8t
device special files, 5–29	DL_UNIDATA_REQ primitive, 2-8t
Digital UNIX system calls	DL_UNITDATA_IND primitive, 7–12
and DLI, E–16	DL_UNITDATA_REQ primitive, 7–12
distributed applications	DL_XID_CON primitive, 2-8t
and the client/server paradigm, 4-8	DL_XID_IND primitive, 2-8t
DL_ATTACH_REQ primitive, 2–8t, 7–12	DL_XID_REQ primitive, 2-8t
DL_BIND_ACK primitive, 2–8t, 7–12	DL_XID_RES primitive, 2–8t
DL_BIND_REQ primitive, 2–8t, 7–12	dlb STREAMS pseudodriver, 7–11f, 1–9, 2–2,
DL_DETACH_REQ primitive, 7–12	7–1
DL_DETTACH_REQ primitive, 2–8t	DLI
DL_DISABLMULTI_REQ primitive, 7–12	and accessing the local area network, E-3
DL_DISABMULTI_REQ primitive, 2–8t	and transmitting IEEE 802.3 frames, 2-2
DL_ENABMULTI_REQ primitive, 2–8t, 7–12	concepts, E-1 to E-2
DL_ERROR_ACK primitive, 2–8t	definition of, E-1
DL_ETHER media, 7–12	services, E–3
DL_INFO_ACK primitive, 2–8t	using Digital UNIX system calls, E-16
DL_INFO_REQ primitive, 2–8t	using the socket system call, E-17
DL_OK_ACK primitive, 2–8t, 7–12	DLI address family
DL_PHYS_ADDR_ACK primitive, 2–8t, 7–12	specifying, E–17
DL_PHYS_ADDR_REQ primitive, 7–12	DLI client program
DL_PROMISCON_REQ primitive, 7–12	using 802.3 format packets
DL_PROMISCONOFF_REQ primitive, 7–12	example, E–31
DL_SET_PHYS_ADDR_REQ primitive, 7-12	using Ethernet format packets
DL_SUBS_BIND_ACK primitive, 2–8t, 7–12	example, E–24
DL_SUBS_BIND_REQ primitive, 2–8t, 7–12	DLI program
DL_SUBS_UNBIND_ACK primitive, 7–12	including higher-level services, E-4
DL_SUBS_UNBIND_REQ primitive, 2–8t,	transferring data, E-22
7–12	using getsockopt and setsockopt
DL_TEST_CON primitive, 2–8t	example, E–41
	writing, E–15

DLI protocol module	<b>DLPI</b> (cont.)
specifying, E-17	primitives the STREAMS driver must
DLI server program	support, 7–10
using 802.3 format packets	supported media
example, E-36	DL_ETHER, 7–12
using Ethernet packets	supported primitives, 7–12
example, E–27	table of, 2–8
DLI services	types of service
examples of, E–3	acknowledged connectionless mode, 2-6,
dli_802_3_conn subroutine	2–4
example, E-41	connection mode, 2-5
using, E–7	connectionless mode, 2-6
dli_econn subroutine	local management, 2-5
example, E–27	DLPI addressing
using, E–7	identifying components, 2-7
DLPI	<b>DLPI interface</b> , 2–1f
accessing specification online, 2-1n	DLPI option, 7–4
addressing, 2–7	adding to kernel configuration file
PPA, 2–7	at installation, 7–4
and DLS provider, 2–2	with the doconfig command, 7-4
and DLS user, 2-2	DLPI primitives
connection mode of communication in, 2-3	supported in Digital UNIX, 2-8
connection mode service in, 2-5	<b>DLPI service interface</b> , 2–3f
connectionless mode of communication in,	DLS provider
2–3	defined, 2–2
connectionless mode service in, 2-6	DLS user
defined, 2–2	defined, 2–2
DLS providers, 2–8	domain
style 1, 2–8	specifying the, E-7
style 2, 2–8	drivers
local management service in, 2-5	bridging BSD to STREAMS protocol stacks,
modes of communication	7–11
acknowledged connectionless, 2-4, 2-3,	Token Ring, D-1
2–4	DSAP
connection, 2–3	defined, E-12
connectionless, 2-3	

E	eSNMP (cont.)
EAFNOSUPPORT socket error, 4–40t	method routines, 6–33
EBADF socket error, 4–40t	object tables, 6–8
ECONNREFUSED socket error, 4–40t	overview, 6–2
EFAULT socket error, 4-40t	SNMP versions, 6–4
EHOSTDOWN socket error, 4-40t	starting, 6–15
EHOSTUNREACH socket error, 4-40t	function return values, 6–16
EINVAL socket error, 4–40t	order of operation, 6–16
EMFILE socket error, 4–40t	stopping, 6–15
endhostent library call, 4–13t	function return values, 6–16
endnetent library call, 4–13t	order of operation, 6–16
endprotoent library call, 4–13t	subtree, 6–6
endservent library call, 4–13t	subtree_tbl.c file, 6–10
ENETDOWN socket error, 4-40t	subtree_tbl.h file, 6–8
ENETUNREACH socket error, 4-40t	value representation, 6–35
ENOMEM socket error, 4-40t	eSNMP application programming interface
ENOTSOCK socket error, 4-40t	Digital UNIX support for, 6–1
EOPNOTSUPP socket error, 4-40t	Ethernet
EPROTONOSUPPORT socket error, 4-40t	accessing, E–3
EPROTOTYPE socket error, 4–40t	address, E–3
error logging	multiple users, E-3
in STREAMS, 5–31	transmitting messages on, E–3
error recovery	Ethernet frame structure
providing, E-4, E-16	example of, E-4, E-8
errors	function of, E–8
comparison of XTI and sockets, 3-47t	specifying destination system information,
contrast between XTI and TLI, 3-44	E-8
in XTI, 3-64 to 3-65	Ethernet substructure
sockets	filling the, E–20
table of, 4–40t	frame structure, E–8
<b>eSNMP</b> , 1–7	sending and receiving, E–7
application interface, 6-4	ETIMEDOUT socket error, 4–40t
architecture, 6–3	event
components, 6-2	management
implementing a subagent, 6-12	and TLI compatibility, 3–43
introduction, 1–7	

event logging	FDDI (cont.)
in STREAMS, 5–31	source service access point, E-13
events	fdetach library call, 5–14
defined, 3–5	file descriptor
in XTI, 3–10	and protocol independence, 3-42
incoming, 3-16t, 3-10	flow control
outgoing, 3–15t, 3–10	contrast between XTI and TLI, 3-44
tracking in XTI, 3-14	in XTI, 3–10t
tracking of outgoing events	providing, E-4, E-16
in XTI, 3–14	frame format
used by connectionless transport services,	802, E-1
3–37	802.3, E-4, E-11
EWOULDBLOCK socket error, 4–40t	Ethernet, E-1, E-4, E-8
exchange identification	FDDI, E–4
defined, E-14	processing, E-11
function of, E-14	standard, E–1
execution in XTI	frames
modes of, 3–4	building, E–3
expedited data	framework
and connectionless transport services, 3-37	sockets
extensible SNMP	components, 4–2
See eSNMP	STREAMS, 5–1 to 5–32
	components, 5–2
F	end, 5–4
F_GETOWN parameter, 4–58	head, 5–3
F_SETOWN parameter, 4–58	modules, 5–4
fattach library call, 5–14	messages, 5–5
fcntl system call	functions
F_GETOWN parameter, 4–58	allowed sequence of in XTI, 3-21
F_SETOWN parameter, 4–58	and protocol independence, 3-41
fcntl.h file, 3–6t	comparison of XTI and sockets, 3-45t
fd variable	STREAMS, 5–6 to 5–17
and outgoing events, 3-14	
FDDI	
accessing, E–3	
frame format, E–4	

G	Н	
generation of addresses	header files	
comparison of TLI and XTI, 3-44	conventions for specifying, 4-5n	
gethostbyaddr library call, 4–13t	fcntl.h, 3-6t	
gethostbyaddr routine, 4–10	netinet/in.h, 4-15t	
gethostbyname library call, 4–13t	sockets, 4-15 to 4-16	
gethostbyname routine, 4–10	STREAMS, 5–5	
gethostent library call, 4–13t	sys/socket.h, 4-15t	
getmsg function, 5–11	sys/types.h, 4–15t	
getnetbyaddr library call, 4–13t	sys/un.h, 4–15t	
getnetbyaddr routine, 4–11	tiuser.h, 3-6t, 3-43	
getnetbyname library call, 4–13t, 4–13	XTI and TLI, 3-6	
getnetbyname routine, 4–11	xti.h, 3-6t, 3-44	
getnetent library call, 4–13t	high-level services	
getnetent routine, 4–11	providing, E-4, E-16	
getpeername system call, 4–9t	host groups	
getpmsg function, 5–11	defined, 4–47	
getprotobyname library call, 4–13t	hostent data structure, 4–11	
getprotobyname routine, 4–12	htonl library call, 4–13t	
getprotobynumber library call, 4–13t	htons library call, 4–13t	
getprotobynumber routine, 4–12		
getprotoent library call, 4–13t	I	
getprotoent routine, 4–12	I/O control flags functions of, E-10 idle state in XTI, 3-13t ifnet STREAMS module, 7-2f, 1-7, 7-1 using, 7-4 required setup, 7-4 INADDR ANY wildcard address	
getservbyname library call, 4–13t		
getservbyname routine, 4–12		
getservbyport library call, 4–13t		
getservbyport routine, 4–12		
getservent library call, 4–13t		
getservent routine, 4–12		
getsockname system call, 4–9t		
getsockopt system call, 4–28	binding names to addresses, 4–42	
guaranteed delivery	incoming connection pending state	
providing, E-4	in XTI, 3–13t	
	incoming events	
	for tracking by programs, 3–16	

incoming events (cont.)	kernel level functions
in XTI, 3-16t	STREAMS, 5–17 to 5–25
incoming orderly release state	kernel subsystems
in XTI, 3-13t	STREAMS drivers
inet_addr library call, 4–13t	configuring, 5–25
inet_lnaof library call, 4-13t	STREAMS modules
inet_makeaddr library call, 4–13t	configuring, 5–25
inet_netof library call, 4-13t	
inet_network library call, 4-13t	L
inetd daemon, 4–54	libraries
initialization phase	XTI and TLI, 3–6
state transitions allowed, 3–17t input/output multiplexing, 4–55	library calls
	in XTI, 3–8
Internet communication domain	sockets, 4–10 to 4–15
characteristics, 4-4t	STREAMS
interrupt driven socket I/O, 4–58	fattach, 5–14
ioctl function, 5–9	fdetach, 5–14
IP multicasting, 4–46	isastream, 5–13
all hosts group, 4-47	XTI, 3–7
host groups, 4-47	libtli.a library, 3–6
receiving datagrams, 4-49	libxti.a library, 3–6
sending datagrams, 4-47	linking
IP_ADD_MEMBERSHIP, 4–49	with XTI and TLI libraries, 3–6
IP_DROP_MEMBERSHIP, 4–50	listen event, 3–16t
IP_MULTICAST_IF, 4–48	listen system call, 4–9t, 4–25
IP_MULTICAST_LOOP, 4–49	LLC
IP_MULTICAST_TTL, 4–48	sublayer of DLI, E–12
isastream library call, 5–13	LLC Protocol Data Unit Test
	defined, E–15
K	function of, E–15
kernel configuration file	local management service
DLPI option, 7–4	in DLPI, 2–5
STRIFNET option, 7–4	logical data boundaries
kernel implementation	and protocol independence, 3–42
of sockets, 4–6	Logical Link Control
	See LLC

M	modules (cont.)
mapping	STREAMS processing routines for (cont.)
hostnames to addresses, 4–10	configuration processing, 5-21
network names to network numbers, 4–11	open processing, 5–20
protocol names to protocol numbers, 4–12	read side put processing, 5-22
service names to proteon numbers, 4–12	read side service processing, 5-22
master device, 4–59	write side put processing, 5-22
mblk_t data structure, 5–19	write side service processing, 5–22
message block	msghdr data structure, 4–17, 4–18
components, 5–18	and the recvmsg system call, 4-33
data buffer, 5–18	and the sendmsg system call, 4-33
dblk_t control structure, 5–18	Digital UNIX support, 4-18
mblk_t control structure, 5–18	multicast addresses, E-1
message data structures, 5–18	using, E–9
message types	multicasting, 4-46
normal, 5–5	multiple processes
priority, 5–5	synchronization in XTI, 3-20
method routines	multiple users
eSNMP, 6–33	on Ethernet, E-3
mkfifo function, 5–9	multiplexing, 4–55
modes of communication	
sockets, 4–7	N
connection-oriented, 4–7	naming sockets, 4–6
connectionless, 4–8	netdb.h header file, 4–10
modes of execution	netent data structure, 4–11
sockets	netinet/in.h header file, 4–15t
blocking mode, 4–21	network
nonblocking mode, 4-21	accessing, E-3
module data structures, 5–17	network byte order translation, 4–13
module_info, 5-18	network configuration
qinit, 5–17	broadcasting and determining, 4-51
streamtab, 5–18	network device
module_info data structure, 5–18	specifying the, E-7
modules	network library routines, 4–10, 4–11, 4–12
STREAMS processing routines for, 5-20	network programming environment
close processing, 5–21	data link interfaces, 1-3

network programming framework	outgoing orderly release state
sockets, 1-4f	in XTI, 3–13t
STREAMS, 1–4f	
nonblocking mode	Р
See asynchronous execution	packet routing
ntohl library call, 4–13t	providing, E-4
ntohs library call, 4–13t	pass_conn event, 3–16
	passive user
0	defined, 3–3
O_NDELAY value	typical state transitions, 3-21f
support in TLI, 3-43	physical addresses, E–1
object tables	using, E–9
eSNMP, 6–8	physical point of attachment
ocnt variable, 3–17	See PPA
and incoming events, 3-16	pipe function, 5–10
and outgoing events, 3-14	poll function, 5–12
open function, 5–6	in XTI applications, 3-13
open processing, 5–20	porting
opened event, 3–15	and protocol independence, 3-41
option management	guidelines for writing XTI applications, 3-41
and TCP, 3–64	porting applications to XTI, 3-41 to 3-47
optmgmt event, 3–15t	PPA
orderly release	and addressing in DLPI, 2-7
and protocol independence, 3-42	defined, 2–7
defined, 3–9	prerequisites
event indicating, 3–10t	for DLI programming, E-1
out-of-band data	privileges
handling in the socket framework, 4-44	superuser, E–1
receiving, 4–45	processes
sending, 4–45	sharing a single endpoint among multiple,
outgoing connection pending state	3–20
in XTI, 3–13t	synchronization of multiple processes in
outgoing events	XTI, 3–20
for tracking by programs, 3-14	protocol independence
in XTI, 3–15t	for XTI applications, 3–41

protocol type	receiving IP multicast datagrams, 4-49
defined, E–9	recommendations
protocol-specific options	for use of connection-oriented transport and
and protocol independence, 3-42	CLTS, 3–4
protocols	for use of execution modes, 3-5
selecting with the socket system call, 4-42	recompiling TLI programs, 3–43
protoent data structure, 4–12	recv system call, 4–9t, 4–32
pseudoterminals, 4–59	recvfrom system call, 4–9t, 4–33
defined, 4–59	recvmsg system call, 4-9t, 4-35
master device, 4-59	and the msghdr data structure, 4-33
slave device, 4–59	resfd variable
putmsg function, 5–11	and outgoing events, 3-14
putpmsg function, 5–11	round-trip time
	defined, C-1
Q	<b>Routing Information Field</b> , D-3
qinit data structure, 5–17	S
R	sa_family, 4–17
raw sockets, 4–5	select socket call
rcv event, 3–16t	contrast to XTI t_look function, 3-46
reveonnect event, 3–16t	select system call, 4–24
rcvdis1 event, 3–16t	send system call, 4–9t, 4–31
rcvdis3 event, 3–16t	sending IP multicast datagrams, 4–47
reveals event, 3–16t	sendmsg system call, 4–9t, 4–34
rcvudata event, 3–16t	and the msghdr data structure, 4-33
revuderr event, 3–16t	sendto system call, 4–9t, 4–31, E–7
read function, 5–8	sequencing functions
read side put processing, 5–22	in XTI, 3-21
read side service processing, 5–22	servent data structure, 4–12
read system call, 4–29	server process
read-only access	accepting connections, 4-25
support in TLI, 3–43	connection-oriented, 4-25
receiving	connectionless, 4-27
data units, 3–38	defined, 4–8
errors about data units, 3–39	server/client interaction, 4–8

service class	SO_REUSEPORT, 4–50
defined, E–12	SOCK_DGRAM socket, 4-5
values, E-12	SOCK_RAW socket, 4–5
service in XTI	SOCK_STREAM socket, 4–5
modes of, 3–4	sockaddr data structure, 4–16
services	sockaddr structures
providing high-level, E-16	comparing 4.3BSD and 4.4BSD, 4-38f
sethostent library call, 4–13t	sockaddr_dl data structure, E-6
setnetent library call, 4–13t	explanation of, E-4
setprotoent library call, 4–13t	filling in, E-7
setservent library call, 4–13t	sockaddr_dl structure
setsockopt system call, 4–9t, 4–27	and the 802.2 substructure, E-11
IP_ADD_MEMBERSHIP option, 4-49	and the ethernet substructure, E-8
IP_DROP_MEMBERSHIP option, 4-50	sockaddr_in data structure, 4–17
IP_MULTICAST_IF option, 4-48	sockaddr_un data structure, 4–17
IP_MULTICAST_LOOP option, 4-49	socket functions
IP_MULTICAST_TTL option, 4-48	comparison with XTI functions, 3-45t
SO_REUSEPORT option, 4-50	socket interface
setup	and TCP/IP, 4–1
to use the ifnet STREAMS module, 7-4	Digital UNIX support for, 4-1
shared libraries	socketpair system call, 4–9t, 4–20
and TLI, 3–6	sockets
and XTI, 3-6	accept system call, 4-9t
support in XTI, 3–6	advanced topics, 4-41 to 4-55
shutdown system call, 4–9t, 4–36	and handling out-of-band data, 4-44
shutting down sockets, 4-36	application programming interface, 4-6 to
slave device, 4–59	4–40
SNAP_SAP	bind system call, 4-9t
using, E–13	binding names to, 4-22
snd event, 3–15t	BSD, 4–37
snddis1 event, 3–15t	4.4BSD features
snddis2 event, 3–15t	receipt of protocol data, 4-38
sndrel event, 3–15t	variable-length network addresses, 4-38
sndudata event, 3–15t	4.3BSD msghdr data structure, 4-39
SNMP application programming interface	characteristics, 4–3
Digital UNIX support for, 6–1	closing, 4–36

sockets (cont.)	sockets (cont.)
coexistence with STREAMS, 7-1 to 7-12	interrupt driven I/O, 4-58
common errors, 4-40	kernel implementation, 4-6
communication bridge to STREAMS	library calls, 4-10 to 4-15
framework, 7–1	table of, 4–13
communication domains, 4-3	listen system call, 4–9t
Internet domain, 4–4	mapping host names to addresses, 4-10
UNIX domain, 4-4	mapping network names to network
communication properties, 4-3	numbers, 4–11
comparison with XTI, 3-45	mapping protocol names to protocol
connect system call, 4-9t	numbers, 4–12
connection-oriented client program, B-6	mapping service names to port numbers,
connection-oriented programs, B-2 to B-17	4–12
connection-oriented server processes, 4-25	modes of communication, 4-7
connection-oriented server program, B-2	connection-oriented, 4-7
connectionless client program, B-20	connectionless, 4–8
connectionless programs, B-17 to B-30	modes of execution, 4-21
connectionless server processes, 4-27	msghdr data structure, 4-17
connectionless server program, B-17	naming, 4-6
creating, 4-19	programming TCP socket buffer sizes, C-1
data structures, 4-16 to 4-18	reclaiming resources when closing, 4-36
defined, 4–3	recv system call, 4-9t
errors	recvfrom system call, 4-9t
comparison with XTI, 3-47t	recvmsg system call, 4-9t
establishing connections	rewriting applications for XTI, 3-45
clients, 4–23	sample programs
servers, 4–25	client.h file, B-39
fcntl system call	clientauth.c file, B-39
F_GETOWN parameter, 4-58	clientdb.c file, B-41
F_SETOWN parameter, 4-58	common.h file, B-31
flushing data when closing, 4-37	server.h file, B-32
getpeername system call, 4-9t	serverauth.h file, B-33
getsockname system call, 4-9t	serverdb.h file, B-36
getting socket options, 4-27	xtierror.c file, B-38
header files, 4-15 to 4-16	selecting protocols, 4-42
input/output multiplexing, 4-55	send system call, 4-9t

sockets (cont.)	source routing
sendmsg system call, 4-9t	enabling, D-1
sendto system call, 4-9t	source service access point
setsockopt system call, 4-9t	See SSAP
setting process groups for signals, 4-58	SSAP
setting process IDs for signals, 4-58	defined, E-13
setting socket options, 4-27	standard frame formats
shutdown system call, 4-9t	802, E-1
shutting down sockets, 4-36	Ethernet, E–1
sockaddr data structure, 4-16	state transitions
sockaddr_in data structure, 4-17	allowed for data transfer
sockaddr_un data structure, 4-17	connectionless transport services, 3-17t
socket system call, 4-9t	allowed for initialization phase, 3-17
socketpair system call, 4-9t	states
states	comparison of XTI and sockets, 3-47
comparison between sockets and XTI,	in XTI, 3–10, 3–13
3–47	managing in XTI, 3-23
system calls, 4–9 to 4–10	strclean command, 5–31
TCP specific programming information, C-1	Stream
TCP specific programming information, C–1 to C–3	Stream defined, 5–2
to C-3	defined, 5–2
to C–3 transferring data, 4–29	defined, 5–2 ends
to C-3 transferring data, 4-29 types, 4-5	defined, 5–2 ends and device drivers, 5–4
to C-3 transferring data, 4-29 types, 4-5 SOCK_DGRAM, 4-5	defined, 5–2 ends and device drivers, 5–4 head, 5–3f
to C-3 transferring data, 4-29 types, 4-5 SOCK_DGRAM, 4-5 SOCK_RAW, 4-5	defined, 5–2 ends and device drivers, 5–4 head, 5–3f module, 5–4
to C-3 transferring data, 4-29 types, 4-5 SOCK_DGRAM, 4-5 SOCK_RAW, 4-5 SOCK_STREAM, 4-5	defined, 5–2 ends and device drivers, 5–4 head, 5–3f module, 5–4 stream sockets, 4–5
to C-3 transferring data, 4-29 types, 4-5 SOCK_DGRAM, 4-5 SOCK_RAW, 4-5 SOCK_STREAM, 4-5 sockets and STREAMS frameworks	defined, 5–2 ends and device drivers, 5–4 head, 5–3f module, 5–4 stream sockets, 4–5 STREAMS
to C-3 transferring data, 4-29 types, 4-5 SOCK_DGRAM, 4-5 SOCK_RAW, 4-5 SOCK_STREAM, 4-5 sockets and STREAMS frameworks communication between, 1-7	defined, 5–2 ends and device drivers, 5–4 head, 5–3f module, 5–4 stream sockets, 4–5 STREAMS and timeout, 5–25
to C-3 transferring data, 4-29 types, 4-5 SOCK_DGRAM, 4-5 SOCK_RAW, 4-5 SOCK_STREAM, 4-5 sockets and STREAMS frameworks communication between, 1-7 sockets framework, 1-4f, 4-1f, 1-3	defined, 5–2 ends and device drivers, 5–4 head, 5–3f module, 5–4 stream sockets, 4–5 STREAMS and timeout, 5–25 application programming interface, 5–5 to
to C-3 transferring data, 4-29 types, 4-5 SOCK_DGRAM, 4-5 SOCK_RAW, 4-5 SOCK_STREAM, 4-5 sockets and STREAMS frameworks communication between, 1-7 sockets framework, 1-4f, 4-1f, 1-3 components, 4-2	defined, 5–2 ends and device drivers, 5–4 head, 5–3f module, 5–4 stream sockets, 4–5 STREAMS and timeout, 5–25 application programming interface, 5–5 to 5–25
to C-3 transferring data, 4-29 types, 4-5 SOCK_DGRAM, 4-5 SOCK_RAW, 4-5 SOCK_STREAM, 4-5 sockets and STREAMS frameworks communication between, 1-7 sockets framework, 1-4f, 4-1f, 1-3 components, 4-2 relationship to XTI, 1-5f	defined, 5–2 ends and device drivers, 5–4 head, 5–3f module, 5–4 stream sockets, 4–5 STREAMS and timeout, 5–25 application programming interface, 5–5 to 5–25 clone device, 5–30
transferring data, 4–29 types, 4–5 SOCK_DGRAM, 4–5 SOCK_RAW, 4–5 SOCK_STREAM, 4–5 sockets and STREAMS frameworks communication between, 1–7 sockets framework, 1–4f, 4–1f, 1–3 components, 4–2 relationship to XTI, 1–5f sockets header files, 4–15t	defined, 5–2 ends and device drivers, 5–4 head, 5–3f module, 5–4 stream sockets, 4–5 STREAMS and timeout, 5–25 application programming interface, 5–5 to 5–25 clone device, 5–30 close function, 5–7
to C-3 transferring data, 4-29 types, 4-5 SOCK_DGRAM, 4-5 SOCK_RAW, 4-5 SOCK_STREAM, 4-5 sockets and STREAMS frameworks communication between, 1-7 sockets framework, 1-4f, 4-1f, 1-3 components, 4-2 relationship to XTI, 1-5f sockets protocol stacks	defined, 5–2 ends and device drivers, 5–4 head, 5–3f module, 5–4 stream sockets, 4–5 STREAMS and timeout, 5–25 application programming interface, 5–5 to 5–25 clone device, 5–30 close function, 5–7 coexistence with sockets, 7–1 to 7–12
to C-3 transferring data, 4-29 types, 4-5 SOCK_DGRAM, 4-5 SOCK_RAW, 4-5 SOCK_STREAM, 4-5 sockets and STREAMS frameworks communication between, 1-7 sockets framework, 1-4f, 4-1f, 1-3 components, 4-2 relationship to XTI, 1-5f sockets header files, 4-15t sockets protocol stacks bridging to STREAMS drivers, 7-2	defined, 5–2 ends and device drivers, 5–4 head, 5–3f module, 5–4 stream sockets, 4–5 STREAMS and timeout, 5–25 application programming interface, 5–5 to 5–25 clone device, 5–30 close function, 5–7 coexistence with sockets, 7–1 to 7–12 communication bridge to sockets framework

STREAMS (cont.)	STREAMS (cont.)
components (cont.)	synchronization mechanism, 5-23
head, 5–3	using the ifnet STREAMS module, 7-4
modules, 5–4	STREAMS concepts, 5–23
configuring drivers, 5–25	STREAMS drivers
configuring modules, 5–25	bridging STREAMS drivers to sockets
device special files, 5-29	protocol stacks, 7–2
error logging, 5–31	bridging to sockets protocol stacks, 7-2
event logging, 5-31	<b>STREAMS framework</b> , 1–4f, 1–3, 5–1 to 5–32
strclean command, 5-31	relationship to XTI, 1-5f
framework, 5–1 to 5–32	STREAMS header files
functions, 5-6 to 5-17	strlog.h, 5–5
header files, 5–5	stropts.h, 5–5
ioctl function, 5–9	sys/stream.h, 5–5
kernel-level functions, 5-17 to 5-25	STREAMS protocol stacks
library calls, 5–13, 5–14	bridging to BSD drivers, 7-11
message types	STREAMS-based drivers
normal, 5–5	accessing from sockets-based protocol stacks
priority, 5–5	1–8
messages, 5–5	streamtab data structure, 5–18
mkfifo function, 5–9	STRIFNET option, 7–4
open function, 5-6	adding to kernel configuration file
pipe function, 5–10	at installation, 7–4
processing routines	with the doconfig command, 7-4
close processing, 5–21	strlog.h header file, 5–5
configuration processing, 5-21	stropts.h header file, 5–5
for drivers and modules, 5-20	struc sockaddr_in, 4–17
open processing, 5-20	struct sockaddr, 4–16
read side put processing, 5-22	struct sockaddr_un, 4–17
read side service processing, 5-22	structure alignment, D-3
write side put processing, 5-22	subagent
write side service processing, 5-22	implementing, 6-12
putmsg function, 5-11	substructures
putpmsg function, 5-11	802.2, E-11
required setup to use the ifnet STREAMS	Ethernet frame structure, E-8
module, 7–4	filling in, E–7

substructures (cont.)	T_COTS constant, 3–9
sending and receiving, E-7	T_COTS_ORD constant, 3–9
subtree	T_DATA asynchronous event, 3–10t
eSNMP, 6–6	T_DATAXFER state, 3–13t
subtree_tbl.c file	T_DISCONNECT asynchronous event, 3-10t
eSNMP, 6–10	t_errno variable, 3–64
subtree_tbl.h file	T_ERROR event
eSNMP, 6–8	support in TLI, 3-43
synchronization	<b>t_error function</b> , 3–40t, 3–41
of multiple processes in XTI, 3-20	T_EXDATA asynchronous event, 3–10t
synchronous execution in XTI	t_free function, 3–40t, 3–41
defined, 3–5	t_getinfo function, 3–40t, 3–40
syncronization mechanism	t_getstate function, 3–40t, 3–40
in STREAMS, 5–23	T_GODATA asynchronous event, 3–10t
sys/socket.h header file, 4–15t	T_GOEXDATA asynchronous event, 3–10t
sys/stream.h header file, 5–5	T_IDLE state, 3–13t
sys/types.h header file, 4–15t	T_INCON state, 3–13t
sys/un.h header file, 4–15t	<b>T_INREL</b> , 3–13t
system calls	T_LISTEN asynchronous event, 3-10t
calling sequence, E-16	t_listen function, 3–28
sockets, 4–9 to 4–10	<b>t_look function</b> , 3–40t, 3–41
specifying values with, E-7	contrast to select socket call, 3-46
summary of, E-16	T_MORE flag
used to transfer data, E-22	and protocol independence, 3-42
	t_open function, 3–24
T	t_optmgmt function, 3–63
t_accept function, 3–29	T_ORDREL asynchronous event, 3-10t
contrast to accept socket call, 3–46	T_OUTCON state, 3–13t
<b>t_alloc function</b> , 3–40t, 3–41	T_OUTREL state, 3–13t
t_bind function, 3–26	t_rcv function, 3–32
contrast to bind socket call, 3–46	t_revdis function, 3–34
t close function, 3–36	and protocol independence, 3-42
T_CLTS constant, 3–9	t_rcvrel function, 3–35
T_CONNECT asynchronous event, 3–10t	and protocol independence, 3-42
t_connect function, 3–28	t_revudata function, 3–38

t_rcvuderr function, 3–39	Token Ring drivers
and protocol independence, 3-42	and canonical addresses, D-2
t_snd function, 3–31	enabling source routing, D-1
t_snddis function, 3–34	transfer rate
contrast to close socket call, 3-47	defined, C-1
t_sndrel function, 3–35	transferring
and protocol independence, 3-42	state to another endpoint, 3-14
t_sndudata function, 3–37	transferring data
<b>t_sync function</b> , 3–40t, 3–40	with sockets, 4-29
T_UDERR asynchronous event, 3-10t	transitions
t_unbind function, 3–36	between XTI states, 3-17
T_UNBIND state, 3–13t	<b>Transmission Control Protocol</b>
T_UNINIT state, 3–13t	See TCP
purpose of, 3–23	transport endpoint
TCP	defined, 3–2
and round-trip time, C-1	Transport Layer Interface
and the connect system call, 4-24	See TLI
and transfer rate, C-1	transport provider
connection-oriented communication and, 4-7	and state management, 3-23
programming information, C-1 to C-3	defined, 3–2
protocol, 4–7	Transport Service Data Unit
throughput, C-1	See TSDU
window scale option	transport user
configuring the kernel, C-2	defined, 3–2
window size, C-1	trn_units variable
timeout, 5–25	and enabling source routing, D-1
tiuser.h file, 3–6t, 3–43, 3–6	TSDU
TLI	and protocol independence, 3-42
and XTI, 3-1	types of service
compatibility with XTI, 3-43	in DLPI, 2-4
contrast with XTI, 3-43	types of sockets, 4–5
header files, 3-6t	SOCK_DGRAM, 4–5
library and header files, 3-6	SOCK_RAW, 4–5
reference pages, 3-7	SOCK_STREAM, 4–5
TLOOK error message	

XTI events causing, 3-12

U	XID (cont.)
UDP	function of, E-14
and the connect system call, 4–24	XPG3 compliance
protocol, 4–8	and Digital UNIX's XTI, 3-1
unbind event, 3–15t	<b>XTI</b> , 3–2f
unbound state	and TLI, 3–1
in XTI, 3–13t	and XPG3 compliance, 3-1
uninitialized state	application programming interface, 3-4 to
in XTI, 3–13t	3–41
UNIX communication domain, 4–4	code migration XPG3 to XPG4, 3-48
characteristics, 4–4t	comparison with sockets, 3-45
UNIX domain, 4–44	comparison with TLI, 3-43
unnumbered information command	configuring xtiso, 3-65 to 3-69
defined, E–15	during installation, 3-65
function of, E–15	manually, 3–65
User Datagram Protocol	connection indication, 3-10t
See UDP	connection-oriented client program, B-14
	connection-oriented programs, B-2 to B-17
V	connection-oriented server program, B-9
	connectionless client program, B-27
value representation	connectionless programs, B-17 to B-30
eSNMP, 6–35	connectionless server program, B-23
<b>14</b> 7	constants identifying service modes
W	T_CLTS, 3–9
write function, 5–8	T_COTS, 3–9
write side put processing, 5–22	T_COTS_ORD, 3–9
write side service processing, 5–22	contrast with TLI, 3-43
write system call, 4–30	data flow with a sockets-based transport
write-only access	provider, 1–6
support in TLI, 3–43	data flow with a STREAMS-based transport
	provider, 1–6
X	defined, 1-5, 3-1
X/Open Transport Interface	differences between XPG3 and XPG4, 3-47
	to 3–50
See XTI	errors
XID	comparison with sockets, 3-47t
defined, E–14	<u>-</u>

XTI (cont.)	XTI
errors (cont.)	outgoing connection pending state, 3-13t
t_errno variable, 3–64	outgoing events, 3–15t
event tracking, 3–14	outgoing orderly release state, 3-13t
events, 3–10	overview, 3–2
for tracking by programs, 3-14	passing connections to other endpoints, 3-16
incoming, 3–10	phase independent functions
outgoing, 3–10	table of, 3–40t
used by connectionless transport services,	porting applications to, 3-41 to 3-47
3–37	relationship to STREAMS and sockets
events causing T_LOOK error, 3-12	frameworks, 1–5f
functions, 3–7	relationships between users, providers, and
handling errors, 3-64	endpoints, 3–3
header files, 3-6t	rewriting socket applications for, 3-45
incoming events, 3-16t	sample programs
interoperability of XPG3 and XPG4, 3-50	client.h file, B-39
library and header files, 3-6	clientauth.c file, B-39
library calls	clientdb.c file, B-41
table of, 3–8	common.h file, B-31
map of functions, events, and states, 3-17	server.h file, B-32
modes of execution	serverauth.h file, B-33
asynchronous, 3-5, 3-4	serverdb.h file, B-36
synchronous, 3–5	xtierror.c file, B-38
modes of service, 3-4	sequencing functions, 3-21
connection-oriented, 3-4	state management by transport providers,
connectionless, 3-4	3–23
option management, 3-64	states, 3–10, 3–13
xti	comparison between XTI and sockets,
options, 3–50	3–47
format, 3–52	synchronization of multiple processes, 3-20
info argument, 3–62	transport endpoint, 3-3f
management of a transport endpoint, 3-60	using XPG3 programs, 3-49
negotiating, 3–53	writing connection-oriented applications
portability, 3–63	accepting a connection, 3-29
T_UNSPEC, 3–62	binding an address to an endpoint, 3–25, 3–24

```
XTI (cont.)
   writing connection-oriented applications
          (cont.)
      deinitializing endpoints, 3-36
      establishing a connection, 3-27 to 3-30
      initializing an endpoint, 3-24 to 3-64
      initiating a connection, 3-28
      listening for connection indications, 3-27
      negotiating protocol options, 3-63
      opening an endpoint, 3-24
      receiving data, 3-32
      releasing connections, 3-33 to 3-35
      sending data, 3-30
      to deinitialize an endpoints in, 3-36
      to use phase-independent functions, 3-40
      transferring data, 3-30 to 3-33
      using the abortive release of connections,
             3-33
      using the orderly release of connections,
             3-34
   writing connectionless applications
      deinitializing endpoints, 3-39
      initializing endpoints, 3-37
      transferring data, 3-37 to 3-39
XTI asynchronous events
   and consuming functions, 3-11t
   table of, 3-10
XTI states
   table of, 3-13
xti.h file, 3-6
xti.h header file
   and t_errno variable, 3-65, 3-6
xtiso
   configuring, 3-65 to 3-69
```

# **How to Order Additional Documentation**

### **Technical Support**

If you need help deciding which documentation best meets your needs, call 800-DIGITAL (800-344-4825) before placing your electronic, telephone, or direct mail order.

#### **Electronic Orders**

To place an order at the Electronic Store, dial 800-234-1998 using a 1200- or 2400-bps modem from anywhere in the USA, Canada, or Puerto Rico. If you need assistance using the Electronic Store, call 800-DIGITAL (800-344-4825).

## **Telephone and Direct Mail Orders**

<b>Your Location</b>	Call	Contact
Continental USA, Alaska, or Hawaii	800-DIGITAL	Digital Equipment Corporation P.O. Box CS2008 Nashua, New Hampshire 03061
Puerto Rico	809-754-7575	Local Digital subsidiary
Canada	800-267-6215	Digital Equipment of Canada Attn: DECdirect Operations KAO2/2 P.O. Box 13000 100 Herzberg Road Kanata, Ontario, Canada K2K 2A6
International		Local Digital subsidiary or approved distributor
Internal <sup>a</sup>		SSB Order Processing – NQO/V19 or U. S. Software Supply Business Digital Equipment Corporation 10 Cotton Road Nashua, NH 03063-1260

<sup>&</sup>lt;sup>a</sup> For internal orders, you must submit an Internal Software Order Form (EN-01740-07).

#### **Reader's Comments**

**Digital UNIX** Network Programmer's Guide AA-PS2WD-TE

	tions on this manual. Your input will help us to
write documentation that meets your needs. following methods:	Please send your suggestions using one of the

- This postage-paid form
- Internet electronic mail: readers\_comment@zk3.dec.com
- Fax: (603) 881-0120, Attn: UEG Publications, ZKO3-3/Y32

If you are not using this form, please be sure you include the name of the document, the page number, and the product name and version.

Please rate this manual: Accuracy (software works as manual says) Completeness (enough information) Clarity (easy to understand) Organization (structure of subject matter) Figures (useful) Examples (useful) Index (ability to find topic) Usability (ability to access information quickly)		Excellent	Good	Fair	Poor
Please lis	st errors you have found in this ma	anual:			
Page	Description				
	*				
Addition	al comments or suggestions to imp	rovo this mor	anal.		
	-				
What ve	rsion of the software described by	this manual a	re you usi	ng?	
Name/Tit	tle		Dept.		
	<i>.</i>		_		
iviaiiiig A	Address				
	Email		Phor	16	

\_\_ Do Not Cut or Tear – Fold Here and Tape .\_\_\_\_





NO POSTAGE NECESSARY IF MAILED IN THE UNITED STATES

# BUSINESS REPLY MAIL

FIRST-CLASS MAIL PERMIT NO. 33 MAYNARD MA

POSTAGE WILL BE PAID BY ADDRESSEE

DIGITAL EQUIPMENT CORPORATION UEG PUBLICATIONS MANAGER ZKO3-3/Y32 110 SPIT BROOK RD NASHUA NH 03062-9987

Illiandhilliandhaalidhabababababababbbabbbb

Do Not Cut or Tear - Fold Here

Cut on Dotted Line