

The Sun logo, Sun Microsystems, Sun Workstation, NFS, and TOPS are registered trademarks of Sun Microsystems, Inc.

Sun, Sun-2, Sun-3, Sun-4, Sun386i, SPARCstation, SPARCserver, NeWS, NSE, OpenWindows, SPARC, SunInstall, SunLink, SunNet, SunOS, SunPro, and Sun-View are trademarks of Sun Microsystems, Inc.

UNIX is a registered trademark of AT&T; OPEN LOOK is a trademark of AT&T.

All other products or services mentioned in this document are identified by the trademarks or service marks of their respective companies or organizations, and Sun Microsystems, Inc. disclaims any responsibility for specifying which marks are owned by which companies or organizations.

VMEbus is a trademark of Motorola, Incorporated.

VAX is a trademark of Digital Equipment Corporation.

IBM-PC and IBM 370 are trademarks of International Business Machines Corporation.

Cray is a trademark of Cray Research.

All other products or services mentioned in this document are identified by the trademarks or service marks of their respective companies or organizations

The Network Information Service (NIS) was formerly known as Sun Yellow Pages. The functionality of the two remains the same, only the name has changed. The name Yellow Pages™ is a registered trademark in the United Kingdom of British Telecommunications plc and may not be used without permission.

Copyright © 1990 Sun Microsystems, Inc. – Printed in U.S.A.

All rights reserved. No part of this work covered by copyright hereon may be reproduced in any form or by any means – graphic, electronic, or mechanical – including photocopying, recording, taping, or storage in an information retrieval system, without the prior written permission of the copyright owner.

Restricted rights legend: use, duplication, or disclosure by the U.S. government is subject to restrictions set forth in subparagraph (c)(1)(ii) of the Rights in Technical Data and Computer Software clause at DFARS 52.227-7013 and in similar clauses in the FAR and NASA FAR Supplement.

The Sun Graphical User Interface was developed by Sun Microsystems, Inc. for its users and licensees. Sun acknowledges the pioneering efforts of Xerox in researching and developing the concept of visual or graphical user interfaces for the computer industry. Sun holds a non-exclusive license from Xerox to the Xerox Graphical User Interface, which license also covers Sun's licensees.

This product is protected by one or more of the following U.S. patents: 4,777,485 4,688,190 4,527,232 4,745,407 4,679,014 4,435,792 4,719,569 4,550,368 in addition to foreign patents and applications pending.

Contents

Chapter 1	Introduction	1
1.1.	Manual Overview	1
1.2.	STREAMS Context	1
1.3.	A Basic View of a stream	3
	System Calls	4
1.4.	Benefits of STREAMS	5
	Creating Service Interfaces	5
	Manipulating Modules	6
	Protocol Portability	6
	Protocol Substitution	7
	Protocol Migration	7
	Module Reusability	7
1.5.	An Advanced View of a stream	8
	Stream Head	9
	Modules	9
	stream End	10
1.6.	Building a stream	11
	Expanded streams	11
	Pushable Modules	12
1.7.	Basic User Level Functions	13
	STREAMS System Calls	13
	An Asynchronous Protocol stream Example	14
	Initializing the stream	14
	Message Types	15

Sending and Receiving Messages	16
Using Messages in the Example	16
Other User Functions	19
1.8. Kernel Level Functions	19
Messages	19
Message Allocation	21
Put and Service Procedures	22
Put Procedures	22
Service Procedures	22
Kernel Processing	23
Read Side Processing	24
Driver Processing	24
CHARPROC	24
CANONPROC	25
Write Side Processing	25
Analysis	26
1.9. Other Facilities	26
Message Queue Priority	26
Flow Control	27
Multiplexing	29
Monitoring	31
Error and Trace Logging	32
1.10. Driver Design Comparisons	33
Environment	34
Drivers	34
Modules	34
1.11. Glossary	35
Chapter 2 STREAMS Applications Programming	37
2.1. Introduction	37
STREAMS Overview	37
Development Facilities	39
2.2. Basic Operations	40

A Simple Stream	40
Inserting Modules	42
Module and Driver Control	43
2.3. Advanced Operations	46
Advanced Input/Output Facilities	46
Input/Output Polling	46
Asynchronous Input/Output	49
Clone Open	50
2.4. Multiplexed Streams	50
Multiplexor Configurations	50
Building a Multiplexor	52
Dismantling a Multiplexor	57
Routing Data Through a Multiplexor	58
2.5. Message Handling	59
Service Interface Messages	59
Service Interfaces	59
The Message Interface	60
Datagram Service Interface Example	62
Accessing the Datagram Provider	64
Closing the Service	67
Sending a Datagram	67
Receiving a Datagram	68
Chapter 3 STREAMS Module and Driver Programming	71
3.1. Introduction	71
Development Facilities	72
3.2. STREAMS Mechanism	72
Stream Construction	73
Opening a Stream	75
Adding and Removing Modules	76
Closing	76
3.3. Modules	77
Module Declarations	77

Module Procedures	79
Module and Driver Environment	80
3.4. Messages	81
Message Format	81
Message Generation and Reception	83
Filter Module Declarations	83
bappend () Subroutine	84
Message Allocation	85
Put Procedure	86
3.5. Message Queues and Service Procedures	87
The queue_t Structure	87
Service Procedures	88
Message Queues and Message Priority	89
Flow Control	89
Example	90
Procedures	91
3.6. Drivers	94
Overview of Drivers	94
Driver Flow Control	95
Driver Programming	96
Driver Declarations	96
Driver Open	98
Driver Processing Procedures	99
Driver Flush Handling	100
Driver Interrupt	100
Driver and Module Ioctls	101
Driver Close	103
3.7. Example Driver	104
Cloning	104
Loop-Around Driver	104
Write Put Procedure	107
Service Procedures	110
Close	111

3.8. Multiplexing	111
Multiplexing Configurations	112
Connecting Lower Streams	113
Disconnecting Lower Streams	114
Multiplexor Construction Example	115
Multiplexing Driver	117
Upper Write Put Procedure	120
Lower <code>QUEUE</code> Write Service Procedure	123
Lower Read Put Procedure	125
3.9. Service Interface	127
Definition	127
Message Usage	127
Example	128
Declarations	128
Service Interface Procedure	130
3.10. Advanced Topics	132
Recovering From No Buffers	132
Advanced Flow Control	135
Signals	135
Control of Stream Head Processing	136
Read Options	136
Write Offset	137
Chapter 4 SunOS STREAMS Topics	139
4.1. Configuring STREAMS Drivers	139
Module Configuration	140
Tunable Parameters	141
System Error Messages	141
4.2. STREAMS in SunOS	142
STREAM Modules	142
SunOS STREAMS Extension	143
STREAMS Portability	143
User Line Disciplines	144

Appendix A Supplementary STREAMS Material	147
A.1. Kernel Structures	147
streamtab	147
QUEUE Structures	147
A.2. Message Structures	149
iocblk	149
linkblk	150
A.3. Message Types	150
Ordinary Messages	150
High Priority Messages	156
A.4. Utilities	158
Buffer Allocation Priority	159
adjmsg () — Trim Bytes in a Message	159
allocb () — Allocate a Message Block	160
backq () — Get Pointer to Queue Behind a Given Queue	160
bufcall () — Recover from Failure of allocb ()	160
canput () — Test for Room in a Queue	160
copyb () — Copy a Message Block	161
copymsg () — Copy a Message	161
datamsq () — Test Whether Message is a Data Message	161
dupb () — Duplicate a Message Block Descriptor	161
dupmsg () — Duplicate a Message	162
enableok () — Re-allow Queue to be Scheduled	162
esballoc () — Extended STREAMS Buffer Allocation	162
esbbcall () — bufcall for extended buffers	162
flushq () — Flush a Queue	163
freeb () — Free a Message Block	163
freemsg () — Free All Message Blocks in a Message	163
getq () — Get a Message from a Queue	163
insq () — Put a Message at a Specific Place in a Queue	164
linkb () — Concatenate Two Messages into One	164
msgdsiz () — Get Number of Data Bytes in a Message	164
noenable () — Prevent a Queue from Being Scheduled	164

OTHERQ () — Get Pointer to the Mate Queue	164
pullupmsg () — Concatenate Bytes in a Message	165
putbq () — Return a Message to the Beginning of a Queue	165
putctl () — Put a Control Message	165
putctl1 () — Put One-byte Parameter Control Message	165
putnext () — Put a Message to the Next Queue	166
putq () — Put a Message on a Queue	166
qenable () — Enable a Queue	166
qreply () — Send Reverse-Direction Message	167
qsize () — Find the Number of Messages on a Queue	167
RD () — Get Pointer to the Read Queue	167
rmvb () — Remove a Message Block from a Message	167
rmvq () — Remove a Message from a Queue	167
splstr () — Set Processor Level	168
testb () — Check for an Available Buffer	168
unbufcall () — Cancel an outstanding <code>bufcall</code> request	168
unlinkb () — Remove Message Block from Message Head	168
WR () — Get Pointer to the Write Queue	168
A.5. Design Guidelines	169
General Rules	169
System Calls	170
Data Structures	170
Header Files	170
Accessible Symbols and Functions	171
Rules for Put and Service Procedures	172
A.6. STREAMS Glossary	173
Index	177

Figures

Figure 1-1 Basic stream	3
Figure 1-2 STREAMS-Related Manual Pages	5
Figure 1-3 Protocol Module Portability	6
Figure 1-4 Protocol Migration	7
Figure 1-5 Module Reusability	8
Figure 1-6 stream In More Detail	9
Figure 1-7 Setting Up a stream	11
Figure 1-8 Idle stream Configuration for Example	15
Figure 1-9 Asynchronous Terminal streams	18
Figure 1-10 A Message	20
Figure 1-11 Messages on a Message Queue	21
Figure 1-12 Operational stream for Example	23
Figure 1-13 Module Put and Service Procedures	24
Figure 1-14 STREAMS Message Priority	27
Figure 1-15 Flow Control	28
Figure 1-16 Internet Multiplexing stream	29
Figure 1-17 X.25 Multiplexing stream	30
Figure 1-18 Error and Trace Logging	33
Figure 2-1 Basic Stream	39
Figure 2-2 Stream to Communications Driver	41
Figure 2-3 Case Converter Module	43
Figure 2-4 Many-to-one Multiplexor	51
Figure 2-5 One-to-many Multiplexor	51

Figure 2-6 Many-to-many Multiplexor	51
Figure 2-7 Protocol Multiplexor	52
Figure 2-8 Before Link	53
Figure 2-9 IP Multiplexor After First Link	54
Figure 2-10 IP Multiplexor	55
Figure 2-11 TP Multiplexor	56
Figure 2-12 Protocol Substitution	59
Figure 2-13 Service Interface	60
Figure 3-1 Downstream Stream Construction	73
Figure 3-2 QUEUE data structures	74
Figure 3-3 Message Form and Linkage	82
Figure 3-4 Message Queue Priority	89
Figure 3-5 Device Driver Streams	95
Figure 3-6 Loop Around Streams	105
Figure 3-7 Internet Multiplexor Before Connecting	115
Figure 3-8 Internet Multiplexor After Connecting	116
Figure 3-9 Example Multiplexor Configuration	120
Figure A-1 M_PROTO and M_PCPROTO Message Structure	151

Introduction

This manual is a guide to adding drivers for serial communication devices to the SunOS kernel. It is divided into two sections.

- *STREAMS Programming*, discusses topics relevant to the construction and installation of STREAMS drivers and modules.
- *STREAMS Supplementary Material*, summarizes the kernel structures commonly found in STREAMS development.

1.1. Manual Overview

Chapter 1 is an introduction to the STREAMS mechanism.

Chapter 2 describes the development of user-level STREAMS applications.

Chapter 3 discusses, in detail, the development of STREAMS drivers and modules.

Chapter 4 discusses those aspects of the STREAMS mechanism that are unique to SunOS. It covers the few STREAMS-specific configuration topics.

Appendix A contains supplementary STREAMS material, including kernel and message structures, message types, and utilities.

1.2. STREAMS Context

STREAMS was designed to systematize the existing UNIX character I/O mechanism and to support the development of communications services.

STREAMS consist of a set of system calls, kernel resources and kernel routines. For detailed information about the STREAMS-kernel interface, about the internal structure of STREAMS modules and about STREAMS driver programming, see the following sections.

The UNIX system was originally designed as a general-purpose, multi-user, interactive operating system for minicomputers. Initially developed in the 1970's, the system's communications environment included slow to medium speed, asynchronous terminal devices. The original design, the communications environment, and hardware state of the art influenced the character I/O mechanism but the character I/O area did not require the same emphasis on modularity and performance as other areas of the system.

Support for a broader range of devices, speeds, modes, and protocols has since been incorporated into the system, but the original character I/O mechanism, which processes one character at a time, made such development difficult. Additionally, a paucity of tools and the absence of a framework for incorporating contemporary networking protocols added to the difficulty.

Attempts to compensate for the above problems have led to diverse, ad-hoc implementations; for example, protocol drivers are often intertwined with the hardware configuration in which they were developed. As a result, functionally equivalent protocol software often cannot interface with alternate implementations of adjacent protocol layers. Portability, adaptability, and reuse of software have been hindered.

STREAMS, a general, flexible facility and a set of tools for development of UNIX system communication services, is intended to remedy these problems. STREAMS supports services ranging from complete networking protocol suites to individual device drivers.

STREAMS defines standard interfaces for character I/O within the kernel, and between the kernel and the rest of the system. The associated mechanism is simple and open-ended. It consists of a set of system calls, kernel resources, and kernel utility routines. The standard interface and open-ended mechanism enable modular, portable development and easy integration of higher performance network services and their components. STREAMS does not impose any specific network architecture. Instead, it provides a powerful framework with a consistent user interface that is compatible with the existing character I/O interface—which is still available.

STREAMS modularity and design reflect the “layers and options” characteristics of contemporary networking architectures. The basic components in a STREAMS implementation are referred to as modules. These modules, which reside in the kernel, offer a set of processing functions and associated service interfaces. From user level, modules can be dynamically selected and interconnected to provide any rational processing sequence. Kernel programming, assembly, and link editing are not required to create the interconnection. Modules can also be dynamically “plugged into” existing connections from user level. STREAMS modularity allows:

- User level programs that are independent of underlying protocols and physical communication media.
- Network architectures and higher level protocols that are independent of underlying protocols, drivers, and physical communication media.
- Higher level services that can be created by selecting and connecting lower level services and protocols.
- Enhanced portability of protocol modules resulting from STREAMS’ well-defined structure and interface standards.

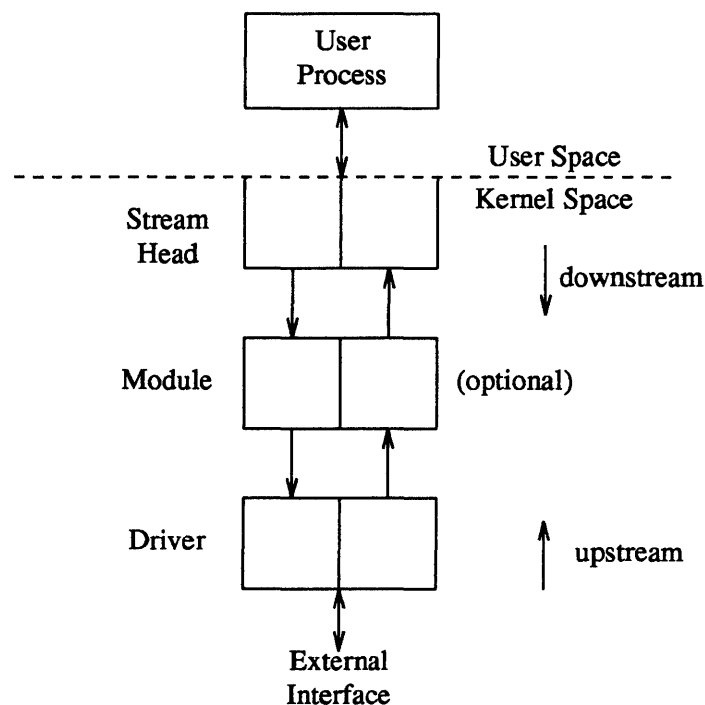
In addition to modularity, STREAMS provides developers with integral functions, a library of utility routines, and facilities that expedite software design and implementation. The principal facilities are:

- Buffer management – To maintain STREAMS' own, independent buffer pool.
- Flow control – To conserve STREAMS' memory and processing resources.
- Scheduling – To incorporate STREAMS' own scheduling mechanism.
- Multiplexing – For processing interleaved data streams, such as occur in SNA, X.25, and windows.
- Asynchronous operation of STREAMS and user processes – Allows STREAMS-related operations to be performed efficiently from user level.
- Error and trace loggers – For debugging and administrative functions.

1.3. A Basic View of a stream

“STREAMS” is a collection of system calls, kernel resources, and kernel utility routines that can create, use, and dismantle a “stream”. A stream is a full-duplex processing and data transfer path between a driver in kernel space and a process in user space (see Figure 1-1).

Figure 1-1 *Basic stream*



A stream has three parts: A stream head, module(s) (optional), and a driver (also referred to as the stream end). The stream head provides the interface between the stream and user processes. Its principal function is to process STREAMS-related user system calls. A module processes data that travel between the stream head and driver. A STREAMS driver may be a device driver, providing the services of an external I/O device, or an internal software driver, commonly called a pseudo-device driver.

Using a combination of system calls, kernel routines, and kernel utilities, STREAMS passes data between a driver and the stream head in the form of messages. Messages that are passed from the stream head toward the driver are said to travel downstream, and messages passed in the other direction travel upstream.

The stream head transfers data between the data space of a user process and STREAMS kernel data space. Data sent to a driver from a user process are packaged into STREAMS messages and passed downstream. Messages arriving at the stream head from downstream are processed by the stream head, and data are copied into user buffers. STREAMS can insert one or more modules into a stream between the stream head and driver to perform intermediate processing of data passing between the stream head and driver.

System Calls

Applications programmers can use the STREAMS facilities via a set of system calls. This system call interface is upward compatible with the existing character I/O facilities. The `open(2)` system call will recognize a STREAMS file and create a stream to the specified driver. A user process can send and receive data using `read(2)` and `write(2)` in the same manner as with character files and devices. The `ioctl(2)` system call enables application programs to perform functions specific to a particular device. In addition, a set of generic STREAMS `ioctl()` commands (see `streamio(4)`) supports a variety of functions for accessing and controlling streams. A `close(2)` will dismantle a stream.

`open()`, `close()`, `read()`, `write()`, and `ioctl()` support the basic set of operations on streams. In addition, new system calls support advanced STREAMS facilities. The `poll(2)` system call enables an application program to poll multiple streams for various events. When used with the STREAMS `I_SETSIG` `ioctl()` command, `poll()` allows an application to process I/O in an asynchronous manner. The `putmsg(2)` and `getmsg(2)` system calls enable application programs to interact with STREAMS modules and drivers through a service interface (described next).

These calls are discussed in this section and, in more detail, the sections that follow. They are precisely specified in the following manual pages:

Figure 1-2 *STREAMS-Related Manual Pages*

<i>Man Page</i>	<i>Description</i>
open.2	Open a stream
close.2	Close a stream
read.2	Read from a stream
write.2	Write to a stream
putmsg.2	Send a message on a stream
getmsg.2	Get next message off a stream
poll.2	STREAMS input/output multiplexing
clone.4	Open any minor device on a STREAMS driver
streamio.4	STREAMS ioctl commands
termio.4	General terminal interface
ttcompat.4m	V7/4BSD compatibility STREAMS module
ldterm.4m	Standard terminal STREAMS module
kbd.4s	Sun keyboard device
kb.4m	Sun keyboard STREAMS module
mouse.4s	Sun mouse device
ms.4m	Sun mouse STREAMS module
pty.4	Pseudo terminal driver
console.4s	Sun console driver and terminal emulator
mti.4s	Systech MTI-800/1600 multi-terminal interface
zs.4s	Zilog 8530 SCC serial communications drive
nit.4p	Network Interface Tap (NIT) Protocol
nit_if.4m	STREAMS NIT device interface
nit_pf.4m	STREAMS NIT packet filtering module
nit_buf.4m	STREAMS NIT buffering module

1.4. Benefits of STREAMS

STREAMS offers two major benefits for applications programmers: easy creation of modules that offer standard data communications services, and the ability to manipulate those modules on a stream.

Creating Service Interfaces

One benefit of STREAMS is that it simplifies the creation of modules that present a service interface to any neighboring application program, module, or device driver. A service interface is defined at the boundary between two neighbors. In STREAMS, a service interface is a specified set of messages and the rules for allowable sequences of these messages across the boundary. A module that implements a service interface will receive a message from a neighbor and respond with an appropriate action (for example, send back a request to retransmit) based on the specific message received and the preceding sequence of messages.

STREAMS provides features that make it easier to design various application processes and modules to common service interfaces. If these modules are written to comply with industry-standard service interfaces, they are called protocol modules.

In general, any two modules can be connected anywhere in a stream. However, rational sequences are generally constructed by connecting modules with compatible protocol service interfaces. For example, a module that implements an X.25 protocol layer, as shown in Figure 1-3, presents a protocol service interface at its input and output sides. In this case, other modules should only be connected to the input and output side if they have the compatible X.25 service interface.

Manipulating Modules

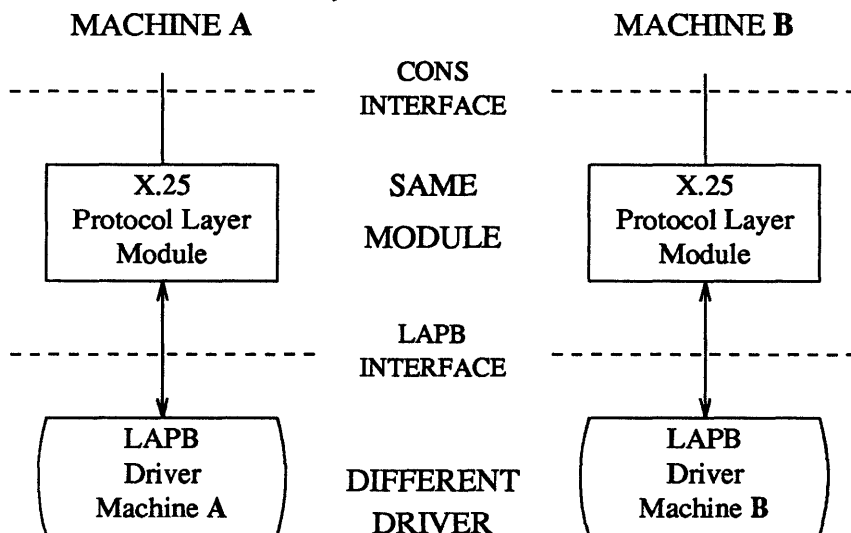
STREAMS provides the capabilities to manipulate modules from user level, to interchange modules with common service interfaces, and to present a service interface to a stream user process. As mentioned above, these capabilities yield benefits when implementing networking services and protocols, including:

- User level programs can be independent of underlying protocols and physical communication media.
- Network architectures and higher level protocols can be independent of underlying protocols, drivers and physical communication media.
- Higher level services can be created by selecting and connecting lower level services and protocols. Below are examples of the benefits of STREAMS capabilities to developers for creating service interfaces and manipulating modules.

Protocol Portability

Figure 1-3 shows how the same X.25 protocol module can be used with different drivers on different machines by implementing compatible service interfaces. The X.25 protocol module interfaces are Connection Oriented Network Service (CONS) and Link Access Protocol – Balanced (LAPB) driver.

Figure 1-3 *Protocol Module Portability*



Protocol Substitution

Alternative protocol modules (and device drivers) can be interchanged on the same machine if they are implemented to an equivalent service interface(s).

Protocol Migration

Figure 1-4 illustrates how STREAMS can migrate functions between kernel software and front end firmware. A common downstream service interface allows the transport protocol module to be independent of the number or type of modules below. The same transport module will connect without modification to either an X.25 module or X.25 driver that has the same service interface.

By shifting functions between software and firmware, developers can produce cost effective, functionally equivalent systems over a wide range of configurations. They can rapidly incorporate technological advances. The same transport protocol module can be used on a lower capacity machine, where economics may preclude the use of front-end hardware, and also on a larger scale system where a front-end is economically justified.

Figure 1-4 *Protocol Migration*

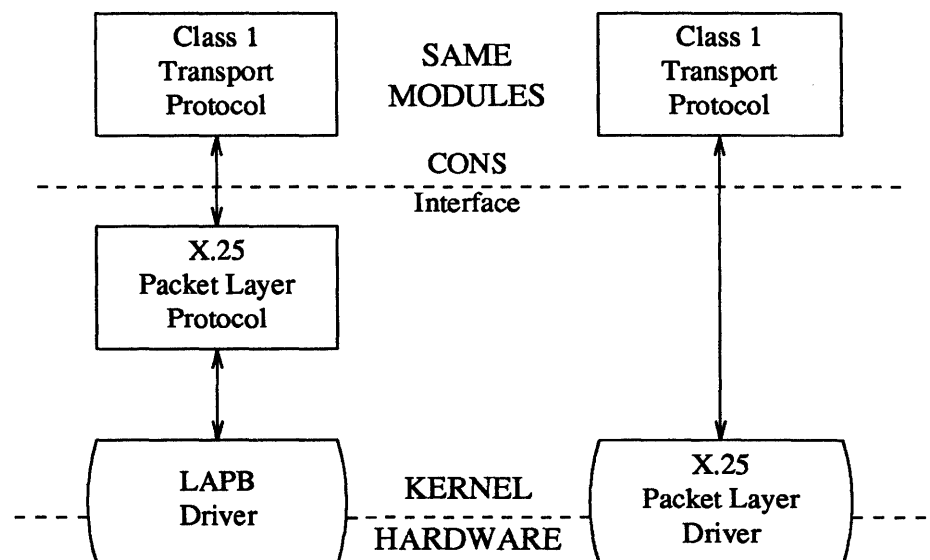
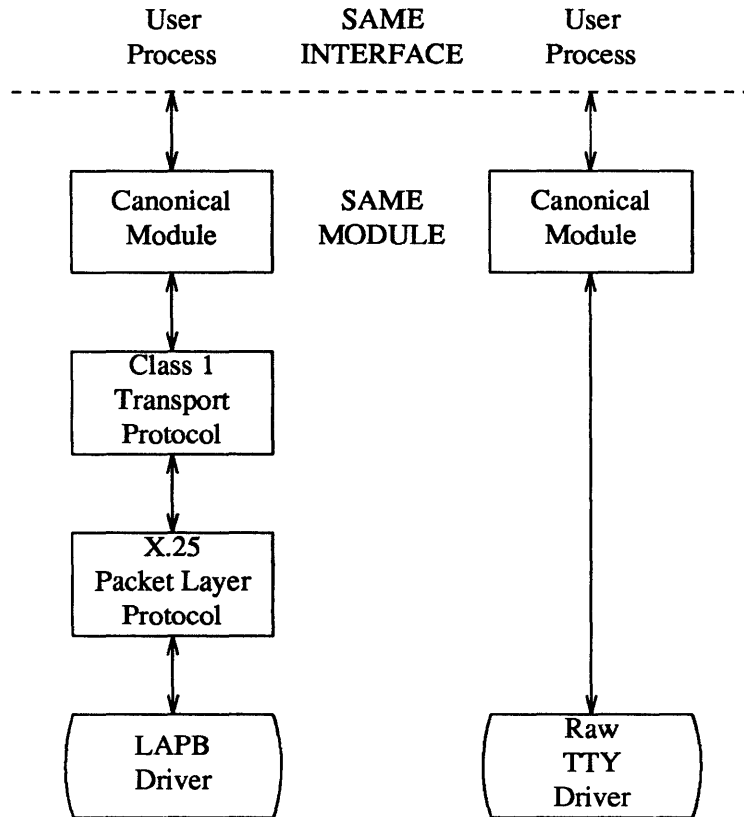
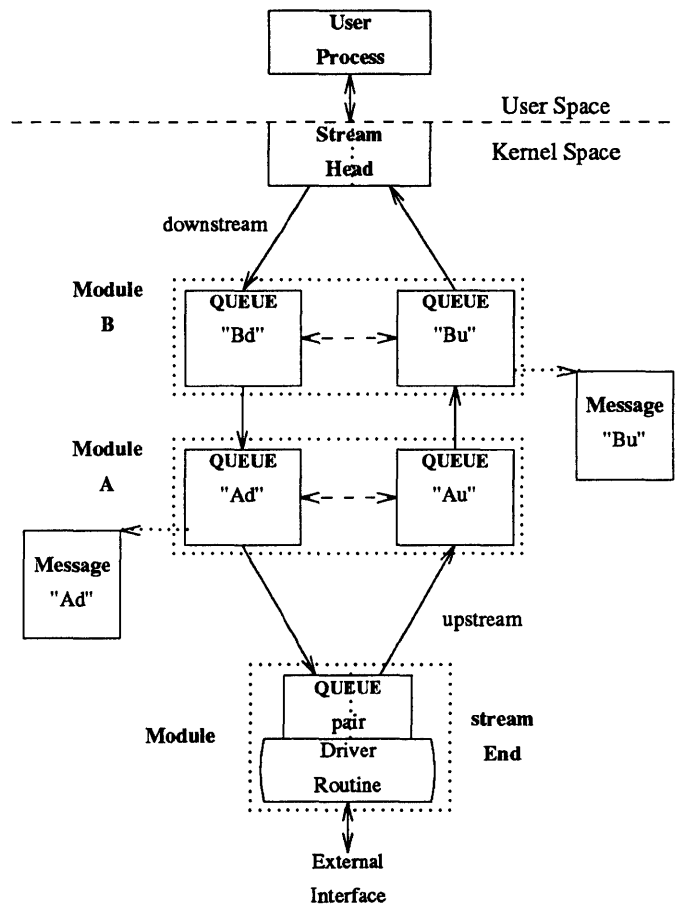
**Module Reusability**

Figure 1-5 shows the same canonical module (for example, one that provides delete and kill processing on character strings) reused in two different streams. This module would typically be implemented as a filter, with no downstream service interface. In both cases, a TTY interface is presented to the stream's user process since the module is nearest the stream head.

Figure 1-5 *Module Reusability*

1.5. An Advanced View of a stream

The STREAMS mechanism constructs a stream by serially connecting kernel resident STREAMS components, each constructed from a specific set of structures. As described earlier and shown in Figure 1-6, the primary STREAMS components are the stream head, optional module(s), and stream end.

Figure 1-6 *stream In More Detail*

Stream Head

The stream head provides the interface between the stream and an application program. The stream head processes STREAMS-related system calls from the application and performs the bidirectional transfer of data and information between the application (in user space) and messages (in STREAMS' kernel space).

Messages are the only means of transferring data and communicating within a stream. A STREAMS message contains data, status/control information, or a combination of the two. Each message includes a specified message type indicator that identifies the contents.

Modules

A module performs intermediate transformations on messages passing between stream head and driver. There may be zero or more modules in a stream (zero when the driver performs all the required character and device processing).

Each module is constructed from a pair of QUEUE structures (see Au/Ad and Bu/Bd in Figure 1-6). A pair is required to implement the bidirectional and symmetrical attributes of a stream. One QUEUE performs functions on messages passing upstream through the module (Au and Bu in Figure 1-6). The other set (Ad and Bd) performs another set of functions on downstream messages. (A

QUEUE, which is part of a module, is different from a message queue, which is described later.)

Each of the two QUEUES in a module will generally have distinct functions, that is, unrelated processing procedures and data. The QUEUES operate independently so that Au will not know if a message passes through Ad unless Ad is programmed to inform it. Messages and data can be shared only if the developer specifically programs the module functions to perform the sharing.

Each QUEUE can directly access the adjacent QUEUE in the direction of message flow (for example, Au to Bu or stream head to Bd). In addition, within a module, a QUEUE can readily locate its mate and access its messages (for example, for echoing) and data.

Each QUEUE in a module may contain or point to messages, processing procedures, or data:

- Messages – These are dynamically attached to the QUEUE on a linked list (“message queue”, see Au and Bd in Figure 1-6) as they pass through the module.
- Processing procedures – A put procedure, to process messages, must be incorporated in each QUEUE. An optional service procedure, to share the message processing with the put procedure, can also be incorporated. According to their function, the procedures can send messages upstream and/or downstream, and they can also modify the private data in their module.
- Data – Developers may provide private data if required by the QUEUE to perform message processing (for example, state information and translation tables).

In general, each of the two QUEUES in a module has a distinct set of all of these elements. Additional module elements will be described later. Although depicted as distinct from modules (see Figure 1-6), a stream head and the stream end also contain a pair of QUEUES.

stream End

A stream end is a module in which the module’s processing procedures are the driver routines. The procedures in the stream end are different from those in other modules because they are accessible from an external device and because the STREAMS mechanism allows multiple streams to be connected to the same driver.

The driver can be a device driver, providing an interface between kernel space and an external communications device, or an internal pseudo-device driver. A pseudo-device driver is not directly related to any external device, and it performs functions internal to the kernel. The multiplexing driver discussed in the *Other Facilities* section is a pseudo-device driver.

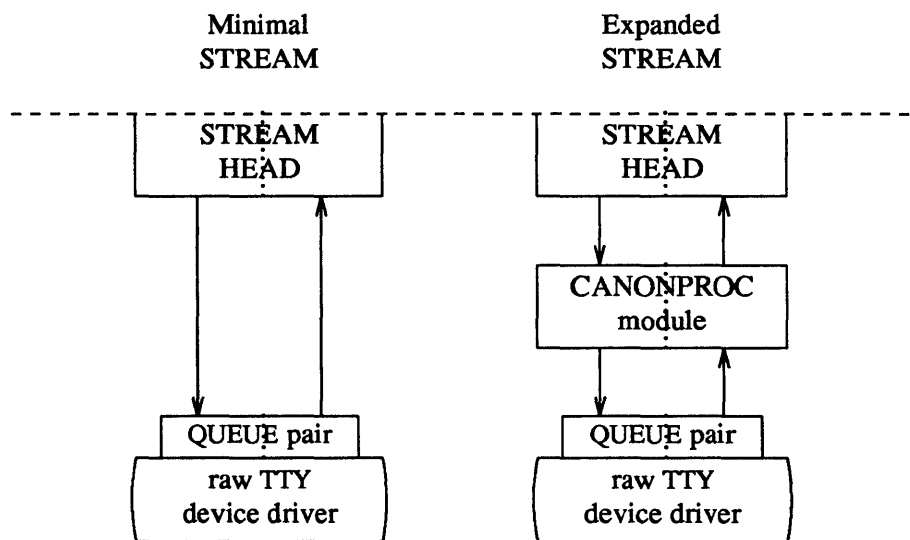
Device drivers must transform all data and status/control information between STREAMS message formats and their external representation. Differences between STREAMS and character device drivers are discussed in the *Driver Design Comparisons* section.

1.6. Building a stream

A stream is created on the first `open(2)` system call to a character special file corresponding to a STREAMS driver. A STREAMS device is distinguished from other character devices by a field contained in the associated `cdevsw` device table entry.

A stream is usually built in two steps. Step one creates a minimal stream consisting of just the stream head and device driver, and step two adds modules to produce an expanded stream (see Figure 1-7). The first step has three parts: head and driver structures are allocated and initialized; the modules in the head and end are linked to each other to form a stream; the driver open routine is called.

Figure 1-7 *Setting Up a stream*



If the driver performs all character and device processing required, no modules need be added to a stream. Examples of STREAMS drivers include a raw tty driver (one that passes along input characters without change) and a driver with multiple streams open to it (corresponding to multiple minor devices opened to a character device driver).

When the driver receives characters from the device, it places them into messages. The messages are then transferred to the next stream component, the stream head, which extracts the contents of the message and copies them to user space. Similar processing occurs for downstream character output; the stream head copies data from user space into messages and sends them to the driver.

Expanded streams

As the second step in building a stream, modules can be added to the stream. In the right-hand stream in Figure 1-7, the CANONPROC module was added to provide additional processing on the characters sent between head and driver.

Modules are added and removed from a stream in last-in-first-out (LIFO) order. They are inserted and deleted at the stream head via the `ioctl(2)` system call. In the stream on the left of Figure 1-5, the X.25 module was the first added to the stream, followed by Class 1 Transport and Canonical modules. To replace the Class 1 module with a Class 0 module, the Canonical module would have to be

removed first, then the Class 1 module, then a Class 0 module would be added and the Canonical module put back.

Because adding and removing modules resembles stack operations, the add is called a push and the remove a pop. Push and pop are two of the `ioctl()` functions included in the STREAMS subset of `ioctl()` system calls. These commands perform various manipulations and operations on streams. The modules manipulated in this manner are called pushable modules, in contrast to the modules contained in the stream head and end. This stack terminology applies only to the setup, modification, and breakdown of a stream.

The stream head processes the `ioctl()` and executes the push, which is analogous to opening the stream driver. Modules are referenced by a unique symbolic name, contained in the STREAMS `fmodsw` module table (similar to the `cdevsw` table associated with a device file). The module table and module name are internal to STREAMS and are accessible from user space only through STREAMS `ioctl()` system calls. The `fmodsw` table points to the module template in the kernel. When a module is pushed, the template is located, the module structures for both QUEUES are allocated, and the template values are copied into the structures.

In addition to the module elements described in *A Basic View of a stream*, each module contains pointers to an open routine and a close routine. The open is called when the module is pushed, and the close is called when the module is popped. Module open and close procedures are similar to a driver open and close.

As in other files, a STREAMS file is closed when the last process open to it closes the file by a `close(2)` system call. This system call causes the stream to be dismantled (modules popped and the driver close executed).

Modules are pushed onto a stream to provide special functions and/or additional protocol layers. In Figure 1-7, the stream on the left is opened in a minimal configuration with a raw tty driver and no other module added. The driver receives one character at a time from the device, places the character in a message, and sends the message upstream. The stream head receives the message, extracts the single character, and copies it into the reading process buffer to send to the user process in response to a `read(2)` system call. When the user process wants to send characters back to the driver, it issues a `write(2)` system call, and the characters are sent to the stream head. The head copies the characters into one or more multi-character messages and sends them downstream. An application program requiring no further kernel character processing would use this minimal stream.

A user requiring a more terminal-like interface would need to insert a module to perform functions such as echoing, character-erase, and line-kill. Assuming that the `CANONPROC` module in Figure 1-7 fulfills this need, the application program first opens a raw tty stream. Then, the `CANONPROC` module is pushed above the driver to create a stream of the form shown on the right of the figure. The driver is not aware that a module has been placed above it and therefore continues to send single character messages upstream. The module receives single

Subsequent use of the word "module" will refer to those pushable modules between stream head and end.

Pushable Modules

character messages from the driver, processes the characters, and accumulates them into line strings. Each line is placed into a message and sent to the stream head. The head now finds more than one character in the messages it receives from downstream.

stream head implementation accommodates this change in format automatically and transfers the multiple-character data into user space. The stream head also keeps track of messages partially transferred into user space (for example, when the current user `read()` buffer can only hold part of the current message). Downstream operation is not affected: the head sends, and the driver receives, multiple character messages.

Note that the stream head provides the interface between the stream and user process. Modules and drivers do not have to implement user interface functions other than open and close.

1.7. Basic User Level Functions

STREAMS System Calls

After a stream has been opened, STREAMS-related system calls allow a user process to insert and delete (push and pop) modules. That process can then communicate with and control the operation of the stream head, modules, and drivers, and can send and receive messages containing data and control information. This section presents an example of some of the basic functions available to STREAMS-based applications via the system calls. Additional functions are described at the end of this section and in the *Other Facilities* section.

The full set of STREAMS-related system calls is:

- open()**
Open a stream
- close()**
Close a stream
- read()**
Read data from a stream
- write()**
Write data to a stream
- ioctl()**
Control a stream
- getmsg()**
Receive the message at stream head
- putmsg()**
Send a message downstream
- poll()**
Notify the application program when selected events occur on a stream

The following two-part example describes a stream that controls the data communication characteristics of a connection between an asynchronous terminal and

a tty port. It illustrates basic user level STREAMS features, then shows how messages can be used. The *Kernel Level Functions* section discusses the kernel stream operations corresponding to the user operations described in this introduction.

An Asynchronous Protocol stream Example

In the example, our computer supports different kinds of asynchronous terminals, each logging in on its own port. The port hardware is limited in function; for example, it detects and reports line and modem status, but does not check parity.

Communications software support for these terminals is provided via a STREAMS implemented asynchronous protocol. The protocol includes a variety of options that are set when a terminal operator dials in to log on. The options are determined by a *getty*-type STREAMS user process, `getstrm()`, which analyzes data sent to it through a series of dialogs (prompts and responses) between the process and terminal operator.

The process sets the terminal options for the duration of the connection by pushing modules onto the stream or by sending control messages to cause changes in modules (or in the device driver) already on the stream. The options supported include:

- ASCII or EBCDIC character codes
- For ASCII code, the parity (odd, even or none)
- Echo or not echo input characters
- Canonical input and output processing or transparent (raw) character handling

These options are set with the following modules:

CHARPROC

Provides input character processing functions, including dynamically settable (via control messages passed to the module) character echo and parity checking. The module's default settings are to echo characters and not check character parity.

CANONPROC

Performs canonical processing on ASCII characters upstream and downstream (note that this performs some processing in a different manner from the standard UNIX character I/O tty subsystem).

ASCEBC

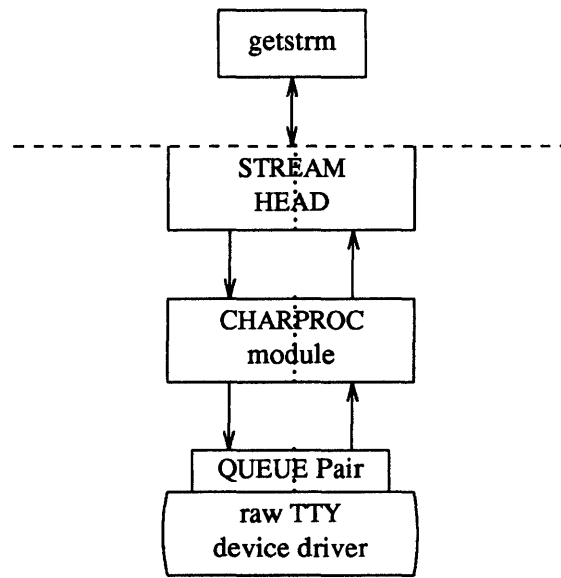
Translates EBCDIC code to ASCII upstream and ASCII to EBCDIC downstream.

Initializing the stream

At system initialization a user process, `getstrm()`, is created for each tty port. `getstrm()` opens a stream to its port and pushes the **CHARPROC** module onto the stream by use of an `ioctl()` `I_PUSH` command. Then, the process issues a `getmsg()` system call to the stream and sleeps until a message reaches the stream head. The stream is now in its idle state.

The initial idle stream, shown in Figure 1-8, contains only one pushable module, CHARPROC. The device driver is a limited function raw tty driver connected to a limited-function communication port. The driver and port transparently transmit and receive one unbuffered character at a time.

Figure 1-8 *Idle stream Configuration for Example*



Upon receipt of initial input from a tty port, `getstrm()` establishes a connection with the terminal, analyzes the option requests, verifies them, and issues STREAMS system calls to set the options. After setting up the options, `getstrm()` creates a user application process. Later, when the user terminates that application, `getstrm()` restores the stream to its idle state by use of system calls.

The next step is to analyze in more detail how the stream sets up the communications options. Before doing so, let's examine how messages are handled in STREAMS.

Message Types

All STREAMS messages are assigned message types to indicate their intended use by modules and drivers and to determine their handling by the stream head. A driver or module can assign most types to a message it generates, and a module can modify a message's type during processing. The stream head will convert certain system calls to specified message types and send them downstream, and it will respond to other calls by copying the contents of certain message types that were sent upstream. Messages exist only in the kernel, so a user process can only send and receive buffers. The process is not explicitly aware of the message type, but it may be aware of message boundaries, depending on the system call used (see the distinction between `getmsg()` and `read()` in the next section).

Most message types are internal to STREAMS and can only be passed from one STREAMS module to another. A few message types, including `M_DATA`, `M_PROTO`, and `M_PCPROTO`, can also be passed between a stream and user

processes. `M_DATA` messages carry data within a stream and between a stream and a user process. `M_PROTO` or `M_PCPROTO` messages carry both data and control information. However, the distinction between control information and data is generally determined by the developer when implementing a particular stream. Control information includes service interface information, carried between two stream entities that present service interfaces, and condition or status information, which may be sent between any two stream entities regardless of their interface. An `M_PCPROTO` message has the same general use as an `M_PROTO`, but the former moves faster through a stream (see *Message Queue Priority* in the *Other Facilities* section).

Sending and Receiving Messages

`putmsg()` is a STREAMS-related system call that sends messages; it is similar to `write()`. `putmsg()` provides a data buffer which is converted into an `M_DATA` message, and can also provide a separate control buffer to be placed into an `M_PROTO` or `M_PCPROTO` block. `write()` provides byte-stream data to be converted into `M_DATA` messages.

`getmsg()` is a STREAMS-related system call that accepts messages; it is similar to `read()`. One difference between the two calls is that `read()` accepts only data (messages sent upstream to the stream head as message type `M_DATA`), such as the characters entered from the terminal. `getmsg()` can simultaneously accept both data and control information (messages sent upstream as types `M_PROTO` or `M_PCPROTO`). `getmsg()` also differs from `read()` in that it preserves message boundaries so that the same boundaries exist above and below the stream head (that is, between a user process and a stream). `read()` generally ignores message boundaries, processing data as a byte stream.

Certain STREAMS `ioctl()` commands, such as `I_STR`, also cause messages to be sent or received on the stream. `I_STR` provides the general “ioctl” capability of the character I/O subsystem. A user process above the stream head can issue `putmsg()`, `getmsg()`, the `I_STR ioctl()` command, and certain other STREAMS related system calls. Other STREAMS `ioctl`'s perform functions that include changing the state of the stream head, pushing and popping modules, or returning special information.

In addition to message types that explicitly transfer data to a process, some messages sent upstream result in information transfer. When these messages reach the stream head, they are transformed into various forms and sent to the user process. The forms include signals, error codes, and call return values.

Using Messages in the Example

Returning to the asynchronous protocol example, the stream was in its idle configuration (see Figure 1-8). `getstrm()` had issued a `getmsg()` and was sleeping until the arrival of a message from the stream head. Such a message would result from the driver detecting activity on the associated tty port.

An incoming call arrives at port one and causes a ring detect signal in the modem. The driver receives the ring signal, answers the call, and sends upstream an `M_PROTO` message containing information indicating an incoming call. `getstrm()` is notified of all incoming calls, although it can choose to refuse the call because of system limits. In this idle state, `getstrm()` will also accept `M_PROTO` messages indicating, for example, error conditions such as detection

of line or modem problems on the idle line.

The `M_PROTO` message containing notification of the incoming call flows upstream from the driver into `CHARPROC`. `CHARPROC` inspects the message type, determines that message processing is not required, and passes the unmodified message upstream to the stream head. The stream head copies the message into the `getmsg()` buffers (one buffer for control information, the other for data) associated with `getstrm()` and wakes up the process. `getstrm()` sends its acceptance of the incoming call with a `putmsg()` system call which results in a downstream `M_PROTO` message to the driver.

Then, `getstrm()` sends a prompt to the operator with a `write()` and issues a `getmsg()` to receive the response. A `read()` could have been used to receive the response, but the `getmsg()` call allows concurrent monitoring for control (`M_PROTO` and `M_PCPROTO`) information. `getstrm()` will now sleep until the response characters, or information regarding possible error conditions detected by modules or driver, are sent upstream.

The first response, sent upstream in an `M_DATA` block, indicates that the code set is ASCII and that canonical processing is requested. `getstrm()` implements these options by pushing `CANONPROC` onto the stream, above `CHARPROC`, to perform canonical processing on the input ASCII characters.

The response to the next prompt requests even parity checking. `getstrm()` sends an `ioctl()` `I_STR` command to `CHARPROC`, requesting the module to perform even parity checking on upstream characters. When the dialog indicates protocol option setting is complete, `getstrm()` creates an application process. At the end of the connection, `getstrm()` will pop `CANONPROC` and then send a `I_STR` to `CHARPROC` requesting the module to restore the no-parity idle state (`CHARPROC` remains on the stream).

As a result of the above dialogs, the terminal at port one operates in the following configuration:

- ASCII, even parity
- Echo
- Canonical processing

In similar fashion, an operator at a different type of terminal on port two requests a different set of options, resulting in the following configuration:

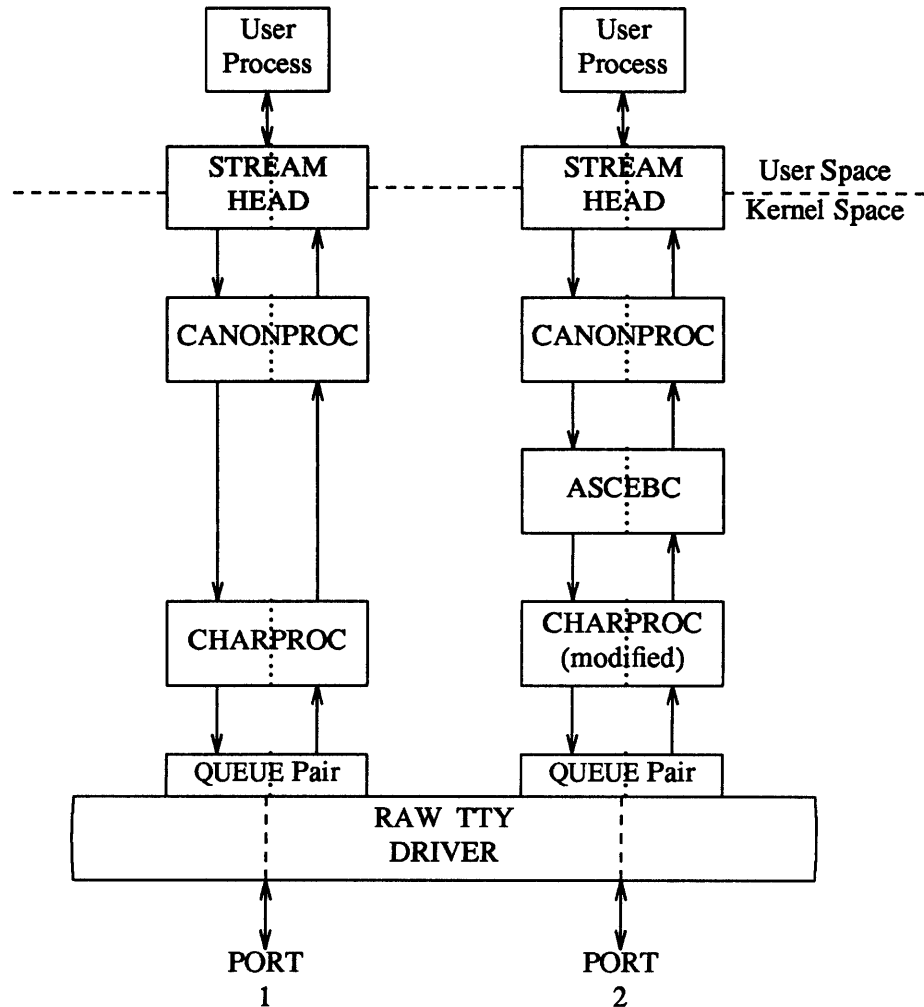
- EBCDIC
- No Echo
- Canonical processing

The resultant streams for the two ports are shown in Figure 1-9. For port one, on the left, the modules in the stream are `CANONPROC` and `CHARPROC`.

For port two, on the right, the resultant modules are `CANONPROC`, `ASCEBC` and `CHARPROC`. `ASCEBC` has been pushed on this stream to translate between the ASCII interface at the downstream side of `CANONPROC` and the EBCDIC interface of the upstream output side of `CHARPROC`. In addition, `getstrm()` has

sent an `I_STR` to the `CHARPROC` module in this stream requesting it to disable echo. The resultant modification to `CHARPROC`'s functions is indicated by the word "modified" in the right stream of Figure 1-9.

Figure 1-9 *Asynchronous Terminal streams*



Since `CHARPROC` is now performing no function for port two, it might have been popped from the stream to be reinserted by `getstrm()` at the end of connection. However, the low overhead of STREAMS does not require its removal. The module remains on the stream, passing messages unmodified between `ASCEBC` and the driver. At the end of the connection, `getstrm()` restores this stream to its idle configuration of Figure 1-8 by popping the added modules and then sending an `I_STR` to `CHARPROC` to restore the echo default.

Note that the tty driver shown in Figure 1-9 handles minor devices. Each minor device has a distinct stream connected from user space to the driver. This ability to handle multiple devices is a standard STREAMS feature, similar to the minor device mechanism in character I/O device drivers.

Other User Functions

The previous example illustrates basic STREAMS concepts. Alternate, more efficient, STREAMS calls or mechanisms could have been used in place of those described earlier. Some of the alternatives are described in the *Other Facilities* section. For details, see following chapters and the *SunOS Reference Manual*.

For example, the initialization process that created a `getstrm()` for each tty port could have been implemented as a “supergetty” by use of the STREAMS-related `poll()` system call. As described in the *Other Facilities* section, `poll()` allows a single process to efficiently monitor and control multiple streams. The “supergetty” process would handle all of the stream and terminal protocol initialization and would create application processes only for established connections.

The `M_PROTO` notification sent to `getstrm()` could have been sent by the driver as an `M_SIG` message that causes a specified signal to be sent to the process. As discussed previously under *Message Types*, error and status information can also be sent upstream from a driver or module to user processes via different message types. These messages will be transformed by the stream head into a signal or error code.

Finally, an `ioctl I_STR` command could have been used in place of a `putmsg M_PROTO` message to send information to a driver. The sending process must receive an explicit response from an `I_STR` by a specified time period or an error will be returned. A response message must be sent upstream by the destination module or driver to be translated into the user response by the stream head.

1.8. Kernel Level Functions

This section introduces the use of the STREAMS mechanism in the kernel and describes some of the tools provided by STREAMS to assist in the development of modules and drivers. In addition to the basic message passing mechanism and `QUEUE` stream linkage described previously, the STREAMS mechanism consists of various facilities including buffer management, the STREAMS scheduler, processing and message priority, flow control, and multiplexing. Over 30 STREAMS utility routines and macros are available to manipulate and utilize these facilities.

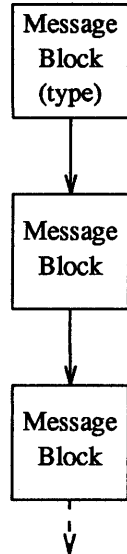
The key elements of a STREAMS kernel implementation are the processing routines in the module and drivers, and the preparation of required data structures. The structures are described in the STREAMS section of *Writing Device Drivers*. The following sections provide further information on messages and on the processing routines that operate on them. The example of the previous section is continued, associating the user-level operations described there with kernel operations.

Messages

As shown in Figure 1-10, a STREAMS message consists of one or more linked *message blocks*. That is, the first message block of a message may be attached to other message blocks that are part of the same message. Multiple blocks in a message can occur, for example, as the result of processing that adds header or trailer data to the data contained in the message, or because of message buffer size limitations which cause the data to span multiple blocks. When a message is composed of multiple message blocks, the message type of the first block determines the type of the entire message, regardless of the types of the attached

message blocks.

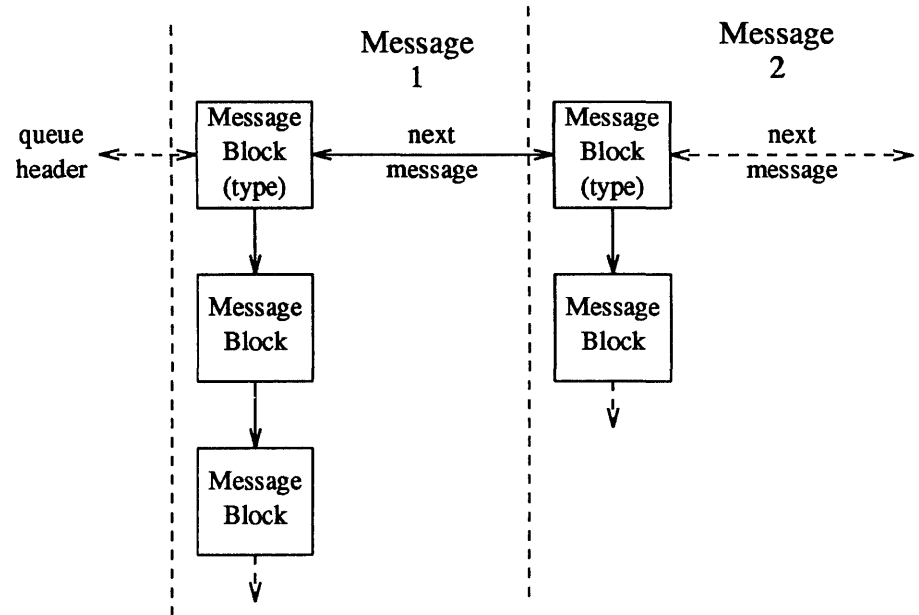
Figure 1-10 *A Message*



STREAMS allocates a message as a single block containing a buffer of a certain size (see the next section). If the data for a message exceed the size of the buffer containing the data, the procedure can allocate a new block containing a larger buffer, copy the current data to it, insert the new data and de-allocate the old block. Alternately, the procedure can allocate an additional (smaller) block, place the new data in the new message block and link it after or before the initial message block. Both alternatives yield one new message.

Messages can exist standalone, as shown in Figure 1-10 when the message is being processed by a procedure. Alternately, a message can await processing on a linked list of messages, called a message queue, in a QUEUE. In Figure 1-11, Message 1 is linked to Message 2.

Figure 1-11 Messages on a Message Queue



When a message is on a queue, the first block of the message contains links to preceding and succeeding messages on the same message queue, in addition to containing a link to the second block of the message (if present). The message queue head and tail are contained in the `QUEUE`.

STREAMS utility routines enable developers to manipulate messages and message queues.

Message Allocation

STREAMS maintains its own storage pool for messages. A procedure can request the allocation of a message of a specified size at one of three message pool priorities. The `allocb()` utility will return a message containing a single block with a buffer of at least the size requested, providing there is a buffer available at the priority requested. When requesting priority for messages, developers must weigh their process' need for resources against the needs of other processes on the same machine.

Message pool priority generally has no effect on allocation until the pool falls below internal STREAMS thresholds. When this occurs, `allocb()` may refuse a lower priority request for a message of size "x" while granting a higher priority request for the same size message. As examples of priority usage, storage for an urgent control message, such as an `M_HANGUP` or `M_PCPROTO` could be requested at high priority. An `M_DATA` buffer for holding input might be requested at medium priority, and an output buffer (presuming the output data can wait in user space) at lowest priority.

Put and Service Procedures

The procedures in the `QUEUE` are the software routines that process messages as they transit the `QUEUE`. The processing is generally performed according to the message type and can result in a modified message, new message(s) or no message. A resultant message is generally sent in the same direction in which it was received by the `QUEUE`, but may be sent in either direction. A `QUEUE` will always contain a put procedure and may also contain an associated service procedure.

Put Procedures

A put procedure is the `QUEUE` routine that receives messages from the preceding `QUEUE` in the stream. Messages are passed between `QUEUES` by a procedure in one `QUEUE` calling the put procedure contained in the following `QUEUE`. A call to the put procedure in the appropriate direction is generally the only way to pass messages between modules (unless otherwise indicated, “modules” infers “module, driver and stream head”). `QUEUES` in pushable (see *Building a stream*) modules contain a put procedure. In general, there is a separate put procedure for the read and write `QUEUES` in a module because of the “full duplex” operation of most streams.

A put procedure is associated with immediate (as opposed to deferred, see below) processing on a message. Each module accesses the adjacent put procedure as a subroutine. For example, consider that *modA*, *modB*, and *modC* are three consecutive modules in a stream, with *modC* connected to the stream head. If *modA* receives a message to be sent upstream, *modA* processes that message and calls *modB*'s put procedure, which processes it and calls *modC*'s put procedure, which processes it and calls the stream head's put procedure. Thus, the message will be passed along the stream in one continuous processing sequence. On one hand, this sequence has the benefit of completing the entire processing in a short time with low overhead (subroutine calls). On the other hand, if this sequence is lengthy and the processing is implemented on a multi-user system, then this manner of processing may be good for this stream but may be detrimental for others since they may have to wait “too long” to get their turn at bat.

In addition, there are situations where the put procedure cannot immediately process the message but must hold it until processing is allowed. The most typical examples of this are a driver which must wait until the current output completes before sending the next message and the stream head, which may have to wait until a process initiates a `read(2)` on the stream.

Service Procedures

STREAMS allows a service procedure to be contained in each `QUEUE`, in addition to the put procedure, to address the above cases and for additional purposes. A service procedure is not required in a `QUEUE` and is associated with deferred processing. If a `QUEUE` has both a put and service procedure, message processing will generally be divided between the procedures. The put procedure is always called first, from a preceding `QUEUE`. After the put procedure completes its part of the message processing, it arranges for the service procedure to be called by passing the message to the `putq()` routine. `putq()` does two things: it places the message on the message queue of the `QUEUE` (see Figure 1-11) and links the `QUEUE` to the end of the STREAMS scheduling queue. When `putq()` returns to the put procedure, the procedure typically exits. Some time later, the service procedure will be automatically called by the STREAMS scheduler.

The STREAMS scheduler is separate and distinct from the SunOS system process scheduler. It is concerned only with QUEUES linked on the STREAMS scheduling queue. The scheduler calls the service procedure of the scheduled QUEUE in a FIFO manner, one at a time.

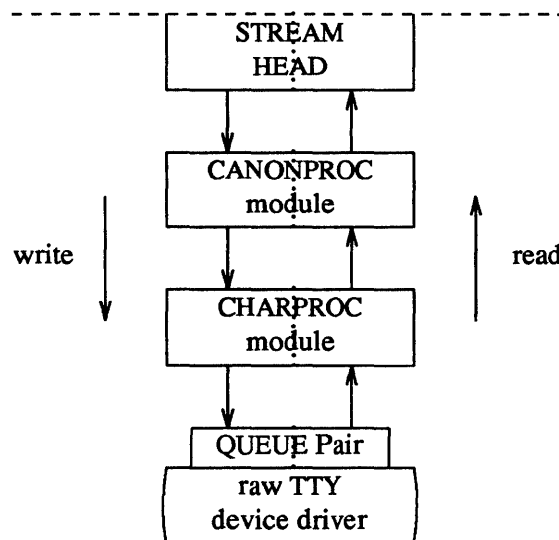
Having both a put and service procedure in a QUEUE enables STREAMS to provide the rapid response and the queuing required in multi-user systems. The put procedure allows rapid response to certain data and events, such as software echoing of input characters. Put procedures effectively have higher priority than any scheduled service procedures. When called from the preceding STREAMS component, a put procedure executes before the scheduled service procedures of any QUEUE are executed.

The service procedure implies message queuing. Queuing results in deferred processing of the service procedure, following all other QUEUES currently on the scheduling queue. For example, terminal output and input erase and kill processing would typically be performed in a service procedure because this type of processing does not have to be as timely as echoing. Use of a service procedure also allows processing time to be more evenly spread among multiple streams. As with the put procedure there will generally be a separate service procedure for each QUEUE in a module. The flow control mechanism (see the *Other Facilities* section) uses the service procedures.

Kernel Processing

The following continues the example of the previous section, describing STREAMS kernel operations and associating them, where relevant, with the user-level system calls already discussed. As a result of initializing operations and pushing a module, the stream for port one has the following configuration:

Figure 1-12 *Operational stream for Example*



As shown in Figure 1-12 the upstream QUEUE is also referred to as the read QUEUE, reflecting the message flow in response to a read() system call. Correspondingly, downstream is referred to as the write QUEUE. Read side processing is discussed first.

Read Side Processing

In our example, read side processing consists of driver processing, CHARPROC processing, and CANONPROC processing.

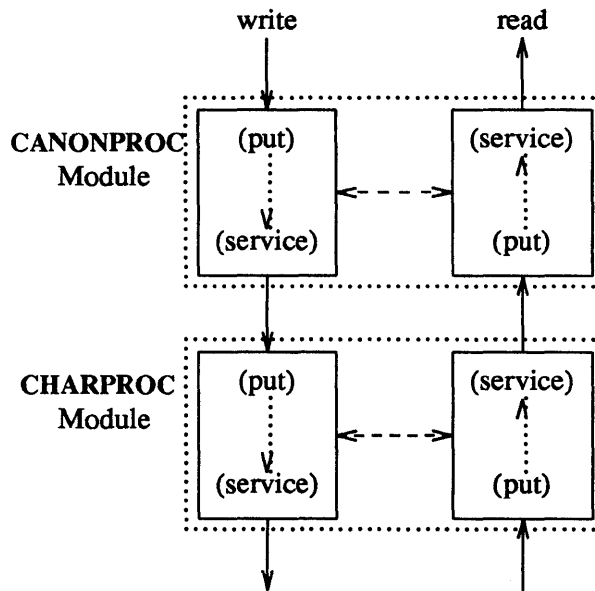
Driver Processing

In the example, the user process has blocked on the `getmsg(2)` system call while waiting for a message to reach the stream head, and the device driver independently waits for input of a character from the port hardware or for a message from upstream. Upon receipt of an input character interrupt from the port, the driver places the associated character in an `M_DATA` message, allocated previously. Then, the driver sends the message to the CHARPROC module by calling CHARPROC's upstream `put` procedure. On return from CHARPROC, the driver calls the `allocb()` utility routine to get another message for the next character.

CHARPROC

CHARPROC has both `put` and `service` procedures on its read side. In the example, the other QUEUES in the modules also have `put` and `service` procedures:

Figure 1-13 *Module Put and Service Procedures*



When the driver calls CHARPROC's read QUEUE `put` procedure, the procedure checks private data flags in the QUEUE. In this case, the flags indicate that echoing is to be performed (recall that echoing is optional and that we are working with port hardware which can not automatically echo). CHARPROC causes the echo to be transmitted back to the terminal by first making a copy of the message with a STREAMS utility. Then, CHARPROC uses another utility to obtain the address of its own write QUEUE. Finally, the CHARPROC read `put` procedure calls its write `put` procedure and passes it the message copy. The write procedure sends the message to the driver to effect the echo and then returns to the read procedure.

This part of read side processing is implemented with `put` procedures so that the entire processing sequence occurs as an extension of the driver input character

interrupt. The CHARPROC read and write put procedures appear as subroutines (nested in the case of the write procedure) to the driver. This manner of processing is intended to produce the character echo in a minimal time frame.

After returning from echo processing, the CHARPROC read put procedure checks another of its private data flags and determines that parity checking should be performed on the input character. Parity should most reasonably be checked as part of echo processing. However, for this example, parity is checked only when the characters are sent upstream. This relaxes the timing in which the checking must occur, that is, it can be deferred along with the canonical processing. CHARPROC uses `putq()` to schedule the (original) message for parity check processing by its read service procedure. When the CHARPROC read service procedure is complete, it forwards the message to the read put procedure of CANONPROC. Note that if parity checking were not required, the CHARPROC put procedure would call the CANONPROC put procedure directly.

CANONPROC

CANONPROC performs canonical processing. As implemented, all read QUEUE processing is performed in its service procedure so that CANONPROC's put procedure simply calls `putq()` to schedule the message for its read service procedure and then exits. The service procedure extracts the character from the message buffer and places it in the "line buffer" contained in another M_DATA message it is constructing. Then, the message which contained the single character is returned to the buffer pool. If the character received was not an end-of-line, CANONPROC exits. Otherwise, a complete line has been assembled and CANONPROC sends the message upstream to the stream head which unblocks the user process from the `getmsg()` call and passes it the contents of the message.

Write Side Processing

The write side of this stream carries two kinds of messages from the user process: `ioctl()` messages for CHARPROC, and M_DATA messages to be output to the terminal.

`ioctl()` messages are sent downstream as a result of an I_STR `ioctl` system call. When CHARPROC receives an `ioctl()` message type, it processes the message contents to modify internal QUEUE flags and then uses a utility to send an acknowledgement message upstream (read side) to the stream head. The stream head acts on the acknowledgement message by unblocking the user from the `ioctl()`.

For terminal output, it is presumed that M_DATA messages, sent by `write()` system calls, contain multiple characters. In general, STREAMS returns to the user process immediately after processing the `write()` call so that the process may send additional messages. Flow control, described in the next section, will eventually block the sending process. The messages can queue on the write side of the driver because of character transmission timing. When a message is received by the driver's write put procedure, the procedure will use `putq()` to place the message on its write-side service message queue if the driver is currently transmitting a previous message buffer. However, there is generally no write QUEUE service procedure in a device driver. Driver output interrupt processing takes the place of scheduling and performs the service procedure functions, removing messages from the queue.

Analysis

For reasons of efficiency, a module implementation would generally avoid placing one character per message and using separate routines to echo and parity check each character, as was done in this example. Nevertheless, even this design yields potential benefits. Consider a case where alternate, more intelligent port hardware was substituted. If the hardware processed multiple input characters and performed the echo and parity checking functions of CHARPROC, then the new driver could be implemented to present the same interface as CHARPROC. Other modules such as CANONPROC could continue to be used without modification.

1.9. Other Facilities

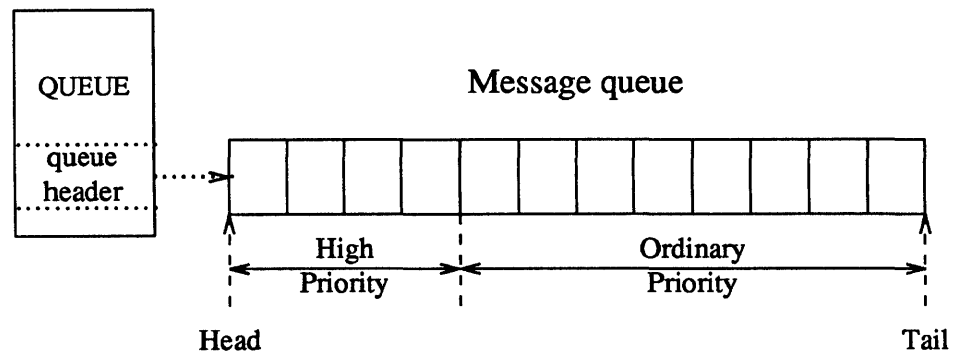
The previous sections described the basic concepts of constructing a stream and utilizing the STREAMS mechanism. Additional STREAMS features are provided to handle characteristic problems of protocol implementation, such as flow control, and to assist in development.

There are also kernel and user-level facilities that support the implementation of advanced functions, such as multiplexors, and allow asynchronous operation of a user process and STREAMS input and output.

Message Queue Priority

As mentioned in the previous section, the STREAMS scheduler operates strictly FIFO so that each QUEUE's service procedure receives control in the order it was scheduled. When a service procedure receives control, it may encounter multiple messages on its message queue. This buildup can occur if there is a long interval between the time a message is queued by a put procedure and the time that the STREAMS scheduler calls the associated service procedure. In this interval, there can be multiple calls to the put procedure causing multiple messages. The service procedure always processes all messages on its message queue unless prevented by flow control (see next section). Each message must pass through all the modules connecting its origin and destination in the stream.

If service procedures were used in all QUEUES and there was no message priority, then the most recently scheduled message would be processed after all the other scheduled messages on all streams had been processed. In certain cases, message types containing urgent information (such as a break or alarm conditions) must pass through the stream quickly. To accommodate these cases, STREAMS provides two classes of message queuing priority, ordinary and high. STREAMS prevents high-priority messages from being blocked by flow control and causes a service procedure to process them ahead of all ordinary priority messages on the procedure's queue. This results in the high-priority message transiting each module with minimal delay.

Figure 1-14 *STREAMS Message Priority*

The priority mechanism operates as shown in Figure 1-14. Message queues are generally not present in a `QUEUE` unless that `QUEUE` contains a service procedure. When a message is passed to `putq()` to schedule the message for service procedure processing, `putq()` places the message on the message queue in priority order. High priority messages are placed ahead of all ordinary priority messages, but behind any other high priority messages on the queue. STREAMS utilities deliver the messages to the processing service procedure FIFO within each priority class. The service procedure is unaware of the message priority and simply receives the next message.

Message priority is defined by the message type; once a message is created, its priority cannot be changed. Certain message types come in equivalent high/ordinary priority pairs (for example, `M_PCPROTO` and `M_PROTO`), so that a module or device driver can choose between the two priorities when sending information.

Flow Control

Even on a well-designed system, general system delays, malfunctions, and excessive message accumulation on one or more streams can cause the message buffer pools to become depleted. Additionally, processing bursts can arise when a service procedure in one module has a long message queue and processes all its messages in one pass. STREAMS provides two independent mechanisms to guard its message buffer pools from being depleted and to minimize long processing bursts at any one module.

Flow control is only applied to normal priority messages (see previous section) and not to high priority messages.

The first flow control mechanism is global and automatic. When the stream head requests a message buffer in response to a `putmsg()` or `write()` system call, it uses the lowest level of priority. Since buffer availability is based on priority and buffer pool levels, the stream head will be among the first modules refused a buffer when the pool becomes depleted. In response, the stream head will block user output until the STREAMS buffer pool recovers. As a result, output has a lower priority than input.

The second flow control mechanism is local to each stream and advisory (voluntary), and limits the number of characters that can be queued for processing at any `QUEUE` in a stream. This mechanism limits the buffers and related processing

at any one QUEUE and in any one stream, but does not consider buffer pool levels or buffer usage in other streams.

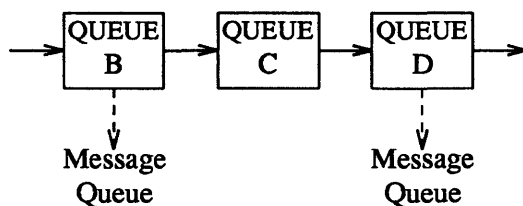
The advisory mechanism operates between the two nearest QUEUES in a stream containing service procedures (see diagram on next page). Messages are generally held on a message queue only if a service procedure is present in the associated QUEUE.

Messages accumulate at a QUEUE when its service procedure processing does not keep pace with the message arrival rate, or when the procedure is blocked from placing its messages on the following stream component by the flow control mechanism. Pushable modules contain independent upstream and downstream limits, which are set when a developer specifies high-water and low-water control values for the QUEUE. The stream head contains a preset upstream limit (which can be modified by a special message sent from downstream) and a driver may contain a downstream limit.

Flow control operates as follows:

- Each time a STREAMS message handling routine (for example, `putq()`) adds or removes a message from a message queue in a QUEUE, the limits are checked. STREAMS calculates the total size of all message blocks on the message queue.
- The total is compared to the QUEUE high-water and low-water values. If the total exceeds the high-water value, an internal full indicator is set for the QUEUE. The operation of the service procedure in this QUEUE is not affected if the indicator is set, and the service procedure continues to be scheduled.
- The next part of flow control processing occurs in the nearest preceding QUEUE that contains a service procedure. In the diagram below, if D is full and C has no service procedure, then B is the nearest preceding QUEUE.

Figure 1-15 *Flow Control*



- The service procedure in B uses a STREAMS utility routine to see if a QUEUE ahead is marked full. If messages cannot be sent, the scheduler blocks the service procedure in B from further execution. B remains blocked until the low-water mark of the full QUEUE, D, is reached.
- While B is blocked, any non-priority messages that arrive at B will accumulate on its message queue (recall that priority messages are not blocked). In turn, B can reach a full state and the full condition will propagate back to the last module in the stream.

- When the service procedure processing on D causes the message block total to fall below the low water mark, the full indicator is turned off. Then, STREAMS automatically schedules the nearest preceding blocked QUEUE (B in this case), getting things moving again. This automatic scheduling is known as back-enabling a QUEUE.

Note that to utilize flow control, a developer need only call the utility that tests if a full condition exists ahead, plus perform some housekeeping if it does. Everything else is automatically handled by STREAMS.

Multiplexing

STREAMS multiplexing supports the development of internetworking protocols such as IP and ISO CLNS, and the processing of interleaved data streams such as in SNA, X.25, and terminal window facilities.

STREAMS multiplexors (also called pseudo-device drivers) are created in the kernel by interconnecting multiple streams. Conceptually, there are two kinds of multiplexors that developers can build with STREAMS: upper and lower multiplexors. Lower multiplexors have multiple lower streams between device drivers and the multiplexor, and upper multiplexors have multiple upper streams between user processes the multiplexor.

Figure 1-16 *Internet Multiplexing stream*

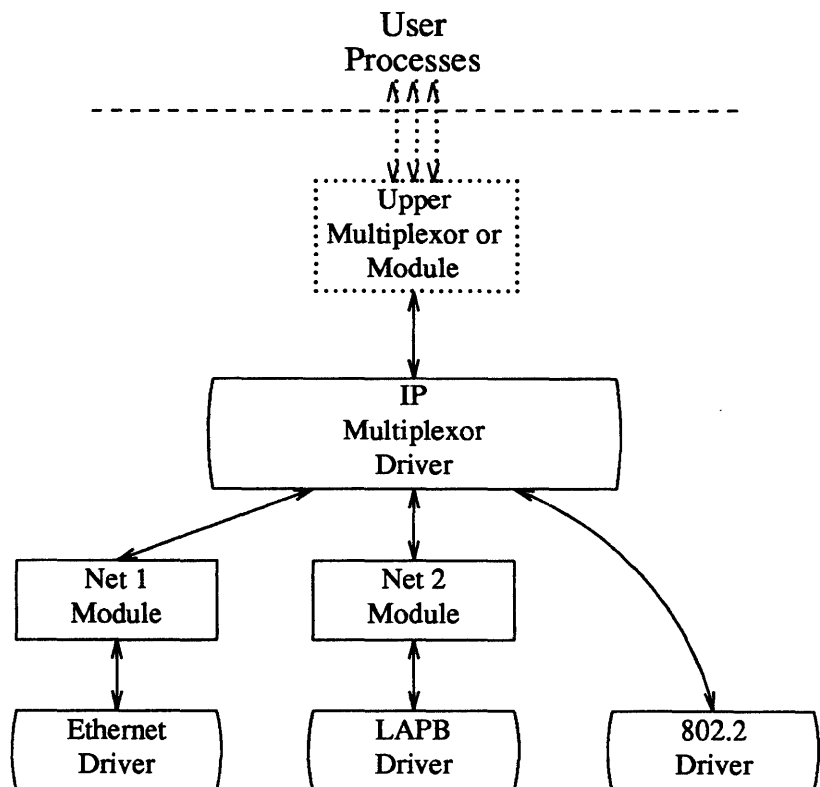


Figure 1-16 shows an example of a lower multiplexor. This configuration would typically occur where internetworking functions were included in the system. This stream contains two types of drivers: the Ethernet, LAPB, and IEEE 802.2

are hardware device drivers that terminate links to other nodes; the IP (Internet Protocol) is a multiplexor.

The IP multiplexor switches messages among the various nodes (lower streams) or sends them upstream to user processes in the system. In this example, the multiplexor expects to see an 802.2 interface downstream; for the Ethernet and LAPB drivers, the Net 1 and Net 2 modules provide service interfaces to the two non-802.2 drivers and the IP multiplexor.

Figure 1-16 depicts the IP multiplexor as part of a larger stream. The stream, as shown in the dotted rectangle, would generally have an upper TCP multiplexor and additional modules. Multiplexors could also be cascaded below the IP driver if the device drivers were replaced by multiplexor drivers.

Figure 1-17 *X.25 Multiplexing stream*

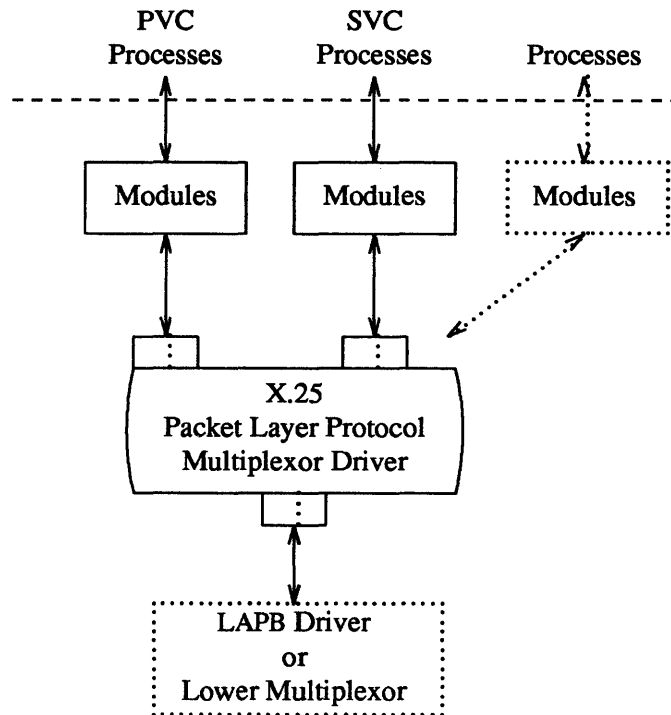


Figure 1-17 shows an upper multiplexor. In this configuration, the driver routes messages between the lower stream and one of the upper streams. This stream performs X.25 multiplexing to multiple independent SVC (Switched Virtual Circuit) and PVC (Permanent Virtual Circuit) user processes. Upper multiplexors are a specific application of standard STREAMS facilities that support multiple minor devices in a device driver. This figure also shows that more complex configurations can be built by having one or more multiplexed LAPB drivers below and multiple modules above.

Developers can choose either upper or lower multiplexing, or both, when designing their applications. For example, a window multiplexor would have a similar

configuration to the X.25 configuration of Figure 1-17, with a window driver replacing Packet Layer, a tty driver replacing LAPB, and the child processes of the terminal process replacing the user processes. Although the X.25 and window multiplexing streams have similar configurations, their multiplexor drivers would differ significantly. The IP multiplexor of Figure 1-16 has a different configuration than the X.25 multiplexor and the driver would implement its own set of processing and routing requirements.

In addition to upper and lower multiplexors, more complex configurations can be created by connecting streams containing multiplexors to other multiplexor drivers. With such a diversity of needs for multiplexors, it is not possible to provide general purpose multiplexor drivers. Rather, STREAMS provides a general purpose multiplexing facility. The facility allows users to set up the inter-module/driver plumbing to create multiplexor configurations of generally unlimited interconnection.

The connections are created from user space through specific STREAMS `ioctl()` system calls. In a lower multiplexor, multiple streams are connected below an application-specific, developer-implemented multiplexing driver. The multiplexing facility will only connect streams to a driver. The `ioctl()` call configures a multiplexor by connecting one stream at a time below the opened multiplexor driver. As each stream is connected to the driver, the connection setup procedure identifies the stream to the driver. The driver will generally store this setup information in a private data structure for later use.

Subsequently, when messages flow into the driver on the various connected streams, the identity of the associated stream is passed to the driver as part of the standard procedure call. The driver then has available the stream identification, the previously stored setup information for this stream, and any internal routing information contained in the message. These data are used, according to the application implemented, to process the incoming message and route the output to the appropriate outgoing stream.

Additionally, new streams can be dynamically connected to an operating multiplexor without interfering with ongoing traffic, and existing streams can be disconnected with similar ease.

Monitoring

STREAMS allows user processes to monitor and control streams so that system resources (such as CPU cycles and process slots) can be used effectively. Monitoring is especially useful to user-level multiplexors, in which a user process can create multiple streams and switch messages among them (similar to STREAMS kernel-level multiplexing, described previously).

User processes can efficiently monitor and control multiple streams with two STREAMS system calls: `poll(2)` and the `ioctl(2)` `I_SETSIG` command. These calls allow a user process to detect events that occur at the stream head on one or more streams, including receipt of a data or protocol message on the read queue and cessation of flow control.

Synchronous monitoring is provided by use of `poll()` alone; in this case, the user process cannot continue processing until after the system call completes. When the calls are used together, they allow asynchronous, or concurrent,

operation of the process and STREAMS input/output. This allows the user process to monitor the stream while carrying on other activities.

To monitor streams with `poll()`, a user process issues that system call and specifies the streams to be monitored, the events to look for, and the amount of time to wait for an event. `poll()` will block the process until the time expires or until an event occurs. If an event occurs, `poll()` will return the type of event and the stream on which the event occurred.

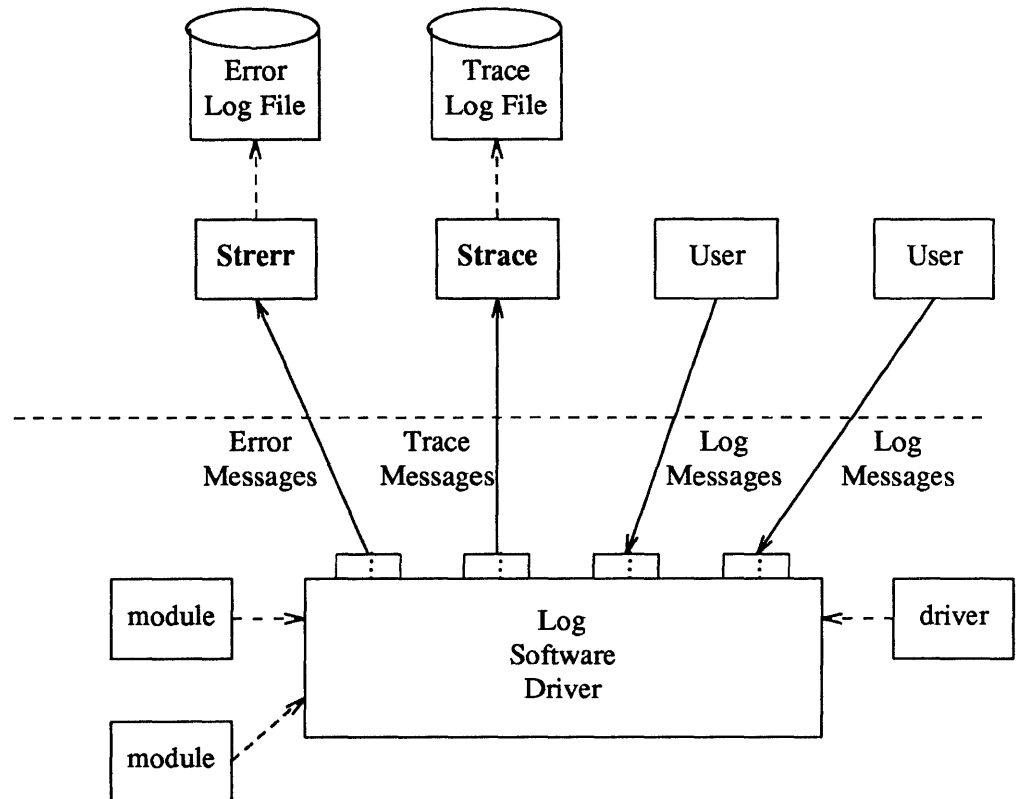
Instead of waiting for an event to occur, a user process may want to monitor one or more streams while processing other data. It can do so by issuing the `ioctl I_SETSIG` command, specifying one or more streams and events (as with `poll()`). Unlike a `poll()`, this `ioctl()` does not force the user process to wait for the event but returns immediately and will issue a signal when an event occurs. The process must also request `signal(2)` or `sigset(2)` to catch the resultant `SIGPOLL` signal.

If any selected event occurs on any of the selected streams, STREAMS will cause the `SIGPOLL` catching function to be executed in all associated requesting processes. However, the process(es) will not know which event occurred, nor on what stream the event occurred. A process that issues the `I_SETSIG` can get more detailed information by issuing a `poll()` after it detects the event.

Error and Trace Logging

STREAMS includes error and trace loggers useful for debugging and administering modules and drivers.

Any module or driver in any stream can call the STREAMS logging functions described in `log(4)`. Formatted text can be sent to the error logger `strerr(8V)`, the trace logger `strace(8V)`, or both.

Figure 1-18 *Error and Trace Logging*

`strerr()` is intended to operate as a daemon process initiated at system startup. `strerr()` formats the contents and places them in a daily file. The utility `strclean(8V)` is provided to periodically purge aged, unreferenced daily log files.

A user process can submit its own `M_PROTO` messages to the log driver for inclusion in the logger of its choice through `putmsg(2)`. The messages must be in the same format required by the logging processes and will be switched to the logger(s) requested in the message.

The output to the log files is formatted, ASCII text. The files can be processed by standard system commands such as `grep(1)` or `ed(1)`, or by developer-provided routines.

1.10. Driver Design Comparisons

This section compares operational features of character I/O device drivers with STREAMS drivers and modules. It is intended for experienced developers of UNIX system character device drivers. Details are provided in the STREAMS section of *Writing Device Drivers*.

Environment

No user environment is generally available to STREAMS module procedures and drivers. The exception is the module and driver open and close routines, both of which have access to the `user` structure of the calling process and can sleep. Otherwise, a STREAMS driver, module put procedure, and module service procedure has no user context and can neither sleep nor access any `user` structure.

Multiple streams can use a copy of the same module (that is, the same `fmodsw`), each containing the same processing procedures. This means that module code is reentrant, so care must be exercised when using global data in a module. Put and service procedures are always passed the address of the `QUEUE` (for example, in Figure 1-6 Au calls Bu's put procedure with Bu as a parameter). The processing procedure establishes its environment solely from the `QUEUE` contents, typically the private data (for example, state information).

Drivers

At the interface to hardware devices, character I/O drivers have interrupt entry points; at the system interface, those same drivers generally have direct entry points (routines) to process `probe()`, `open()`, `close()`, `read()`, `write()` and `ioctl()` system calls.

STREAMS device drivers have similar interrupt entry points at the hardware device interface and have direct entry points only for `open()` and `close()` system calls. These entry points are accessed via STREAMS, and the call formats differ from character device drivers. The put procedure is a driver's third entry point, but it is a message (not system) interface. The stream head translates `write()` and `ioctl()` calls into messages and sends them downstream to be processed by the driver's write `QUEUE` put procedure. `read()` is seen directly only by the stream head, which contains the functions required to process system calls. A driver does not know about system interfaces other than `open()` and `close()`, but it can detect absence of a `read()` indirectly if flow control propagates from the stream head to the driver and affects the driver's ability to send messages upstream.

For input processing, when the driver is ready to send data or other information to a user process, it does not wake up the process. It prepares a message and sends it to the read `QUEUE` of the appropriate (minor device) stream. The driver's open routine generally stores the `QUEUE` address corresponding to this stream.

For output processing, the driver receives messages in place of a `write()` call. If the message can not be sent immediately to the hardware, it may be stored on the driver's write message queue. Subsequent output interrupts can remove messages from this queue.

Drivers and modules can pass signals, error codes, and return values to processes via message types provided for that purpose.

Modules

As described above, modules have user context available only during the execution of their open and close routines. Otherwise, the `QUEUES` forming the module are not associated with the user process at the end of the stream, nor with any other process. Because of this, `QUEUE` procedures must not sleep when they cannot proceed; instead, they must explicitly return control to the system. The system saves no state information for the `QUEUE`. The `QUEUE` must store this

information internally if it is to proceed from the same point on a later entry.

When a module or driver that requires private working storage (for example, for state information) is pushed, the open routine must obtain the storage from external sources. STREAMS copies the module template from `fmodsw` for the `I_PUSH`, so only fixed data can be contained in the module template. STREAMS has no automatic mechanism to allocate working storage to a module when it is opened. The sources for the storage typically include a module-specific kernel array, installed when the system is configured, or the STREAMS buffer pool. When using an array as a module storage pool, the maximum number of copies of the module that can exist at any one time must be determined. For drivers, this is typically determined from the physical devices connected, such as the number of ports on a multiplexor. However, certain types of modules may not be associated with a particular external physical limit. For example, the CANONICAL module shown in Figure 1-5 could be used on different types of streams.

1.11. Glossary

<i>Downstream</i>	The direction from stream head to driver.
<i>Driver</i>	The end of the stream closest to an external interface. The principal functions of the driver are handling any associated device, and transferring data and information between the external interface and stream. It can also be a pseudo-driver, not directly associated with a device, which performs functions internal to a stream, such as a multiplexor or log driver.
<i>Message</i>	One or more linked blocks of data or information, with associated STREAMS control structures containing a message type. Messages are the only means of transferring data and communicating within a stream.
<i>Message Queue</i>	A linked list of messages connected to a <code>QUEUE</code> .
<i>Message Type</i>	A defined set of values identifying the contents of a message.
<i>Module</i>	Software that performs functions on messages as they flow between stream head and driver. A module is the STREAMS counterpart to the commands in a Shell pipeline except that a module contains a pair of functions which allow independent bidirectional (downstream and upstream) data flow and processing.
<i>Multiplexor</i>	A mechanism for connecting multiple streams to a multiplexing driver. The mechanism supports the processing of interleaved data streams and the processing of internetworking protocols. The multiplexing driver routes messages among the connected streams. The other end of a stream connected to a multiplexing driver is typically connected to a device driver.
<i>pushable module</i>	A module between the stream head and driver. A driver is a non-pushable module and a stream head includes a non-pushable module.
<i>QUEUE</i>	The set of structures that forms a module. A module is composed of two <code>QUEUES</code> , a read (upstream) <code>QUEUE</code> and a write (downstream) <code>QUEUE</code> .
<i>Read Queue</i>	The message queue in a module or driver containing messages moving upstream. Associated with input from a driver.

- stream* The kernel aggregate created by connecting STREAMS components, resulting from an application of the STREAMS mechanism. The primary components are a stream head, a driver and zero or more pushable modules between the stream head and driver. A stream forms a full duplex processing and data transfer path in the kernel, between a user process and a driver. A stream is analogous to a Shell pipeline except that data flow and processing are bidirectional.
- Stream Head* The end of the stream closest to the user process. The stream head provides the interface between the stream and the user process. The principal functions of the stream head are processing STREAMS-related system calls, and bidirectional transfer of data and information between a user process and messages in STREAMS' kernel space.
- STREAMS* A kernel mechanism that supports development of network services and data communication drivers. It defines interface standards for character input/output within the kernel, and between the kernel and user level. The STREAMS mechanism comprises integral functions, utility routines, kernel facilities and a set of structures.
- Upstream* The direction from driver to stream head.
- Write Queue* The message queue in a module or driver containing messages moving downstream. Associated with output from a user process.

STREAMS Applications Programming

2.1. Introduction

This chapter provides detailed information about the STREAMS mechanism and system call interface. It includes the following topics.

- The *STREAMS Overview*, below, reintroduces the STREAMS mechanism.
- *Basic Operations* describes the basic operations available for constructing, using, and dismantling streams. These operations are performed using `open(2)`, `close(2)`, `read(2)`, `write(2)`, and `ioctl(2)`.
- *Advanced Operations* presents advanced facilities provided by STREAMS, including: `poll(2)`, a user level I/O polling facility; asynchronous I/O processing support; and a method to sample drivers for available resources.
- *Multiplexed STREAMS* describes the construction of sophisticated, multiplexed stream configurations.
- *Message Handling* describes how users can process STREAMS messages using `putmsg(2)` and `getmsg(2)` in the context of a service interface example.

The following *STREAMS Module and Driver Programming* chapter is the companion to this chapter—it provides an analogous discussion of system-level STREAMS. Both chapters assume a working knowledge of UNIX[†] system programming, data communication facilities, and the material covered in the previous *Introduction to STREAMS* chapter.

STREAMS Overview

This section reviews the STREAMS mechanism, a general, flexible facility and a set of tools for development of SunOS and UNIX system communication services. It supports the implementation of services ranging from complete networking protocol suites to individual device drivers. The STREAMS mechanism defines standard interfaces for character I/O within the kernel, and between the kernel and the rest of the system. The associated mechanism is simple and open-ended. It consists of a set of system calls, kernel resources, and kernel routines.

The standard mechanism enables modular, portable development and easy integration of higher performance network services and their components. STREAMS provide a framework; they do not impose any specific network

[†] UNIX is a registered trademark of AT&T.

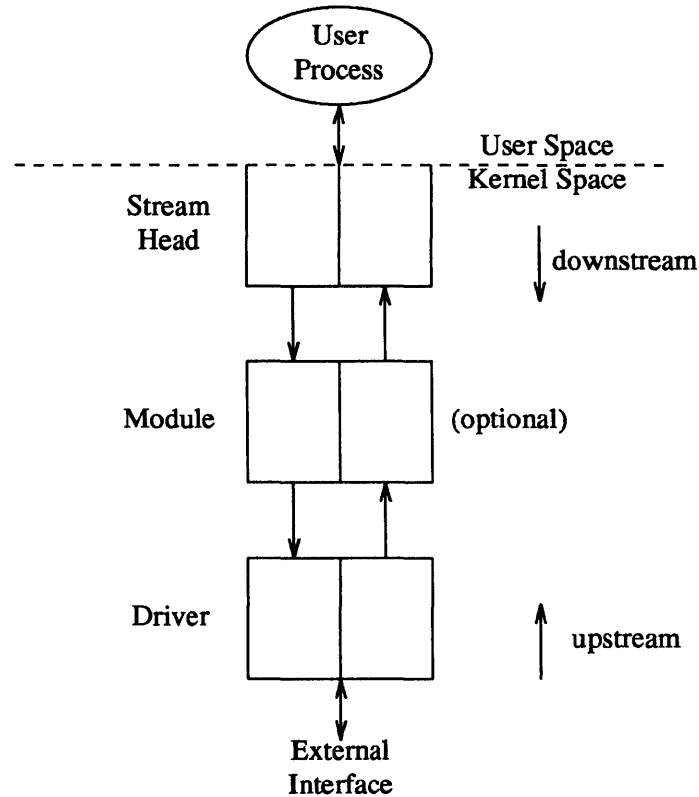
architecture. The STREAMS user interface is upward compatible with the character I/O user interface, and both user interfaces are available.

A stream is a full-duplex processing and data transfer path between a STREAMS driver in kernel space and a process in user space (see the figure below). In the kernel, a stream is constructed by linking a stream head, a driver, and zero or more modules between the stream head and driver. The stream head is the end of the stream closest to the user process. Throughout this guide, the word “STREAMS” refers to the mechanism, and the word *stream* refers to the data path between a user and a driver.¹

A STREAMS driver may be a device driver that provides the services of an external I/O device, or a software driver, commonly referred to as a pseudo-device driver, that performs functions internal to a stream. The stream head provides the interface between the stream and user processes. Its principal function is to process STREAMS-related user system calls.

Data are passed between a driver and the stream head in messages. Messages that are passed from the stream head toward the driver are said to travel downstream. Similarly, messages passed in the other direction travel upstream. The stream head transfers data between the data space of a user process and STREAMS kernel data space. Data to be sent to a driver from a user process are packaged into STREAMS messages and passed downstream. When a message containing data arrives at the stream head from downstream, the message is processed by the stream head, which copies the data into user buffers.

¹ The word “stream” is also used by 4.x BSD to refer to a nonseekable data source such as a pipe or socket. A STREAMS *stream* need not be restricted in this way.

Figure 2-1 *Basic Stream*

Within a stream, messages are distinguished by a type indicator. Certain message types sent upstream may cause the stream head to perform specific actions, such as sending a signal to a user process. Other message types are intended to carry information within a stream and are not directly seen by a user process.

One or more kernel-resident modules may be inserted into a stream between the stream head and driver to perform intermediate processing of data as it passes between the stream head and driver. STREAMS modules are dynamically interconnected in a stream by a user process. No kernel programming, assembly, or link editing is required to create the interconnection.

Development Facilities

General and STREAMS-specific system calls provide the user level facilities required to implement application programs. This system call interface is upwardly compatible with the character I/O facilities. The `open(2)` system call will recognize a STREAMS file and create a stream to the specified driver. A user process can receive and send data on STREAMS files using `read(2)` and `write(2)` in the same manner as with character files. The `ioctl(2)` system call enables users to perform functions specific to a particular device. A set of generic STREAMS `ioctl()` commands, described by `streamio(4)`, support a variety of functions for accessing and controlling streams. A `close(2)` dismantles a stream.

In addition to the generic `ioctl()` commands, there are STREAMS-specific system calls to support unique STREAMS facilities. The `poll(2)` system call enables a user to poll multiple streams for various events. The `putmsg(2)` and `getmsg(2)` system calls enable users to send and receive STREAMS messages, and are suitable for interacting with STREAMS modules and drivers through a service interface.

STREAMS provide kernel facilities and utilities to support development of modules and drivers. The stream head handles most system calls so that the related processing does not have to be incorporated in a module and driver. The configuration mechanism allows modules and drivers to be incorporated into the system.

Examples are used throughout both parts of this document to highlight the most important and common capabilities of STREAMS. The descriptions are not meant to be exhaustive. For simplicity, the examples reference fictional drivers and modules.

2.2. Basic Operations

This section describes the basic set of operations for manipulating STREAMS.

A Simple Stream

A STREAMS driver is similar to a character I/O driver in that it has one or more nodes associated with it in the file system and it is accessed using the `open()` system call. Typically, each file system node corresponds to a separate minor device for that driver. Opening different minor devices of a driver will cause separate streams to be connected between a user process and the driver. The file descriptor returned by the `open()` call is used for further access to the stream. If the same minor device is opened more than once, only one stream will be created; the first `open()` call will create the stream, and subsequent `open()` calls will return a file descriptor that references that stream. Each process that opens the same minor device will share the same stream to the device driver.

Once a device is opened, a user process can send data to the device using the `write()` system call and receive data from the device using the `read()` system call. Access to STREAMS drivers using `read()` and `write()` is compatible with the character I/O mechanism.

The `close()` system call will close a device and dismantle the associated stream.

The following example shows how a simple stream is used. In the example, the user program interacts with a generic communications device that provides point-to-point data transfer between two computers. Data written to the device is transmitted over the communications line, and data arriving on the line can be retrieved by reading it from the device.

```

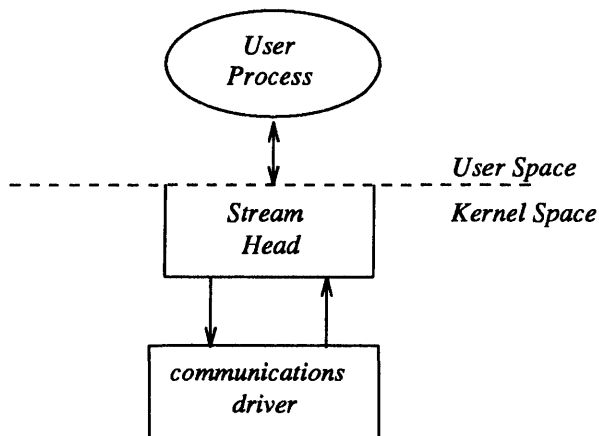
#include <fcntl.h>
main()
{
    char buf[1024];
    int fd, count;

    if ((fd = open("/dev/comm01", O_RDWR)) < 0) {
        perror("open failed");
        exit(1);
    }
    while ((count = read(fd, buf, 1024)) > 0) {
        if (write(fd, buf, count) != count) {
            perror("write failed");
            break;
        }
    }
    exit(0);
}

```

In the example, `/dev/comm01` identifies a minor device of the communications device driver. When this file is opened, the system recognizes the device as a STREAMS device and connects a stream to the driver. The figure below shows the state of the stream following the call to `open()`.

Figure 2-2 *Stream to Communications Driver*



This example illustrates a user reading data from the communications device and then writing the input back out to the same device. In short, this program echoes all input back over the communications line. The example assumes that a user is sending data from the other side of the communications line. The program reads up to 1024 bytes at a time, and then writes the number of bytes just read.

The `read()` call returns the available data, which may contain fewer than 1024 bytes. If no data are currently available at the stream head, the `read()` call blocks until data arrive.

Similarly, the `write()` call attempts to send *count* bytes to `/dev/comm01`. However, STREAMS implements a flow control mechanism that prevents a user from flooding a device driver with data, thereby exhausting system resources. If the stream exerts flow control on the user, the `write()` call blocks until the flow control has been relaxed. The call will not return until it has sent *count* bytes to the device. `exit(2)` is called to terminate the user process. This system call also closes all open files, thereby dismantling the stream in this example.

Inserting Modules

An advantage of STREAMS over the existing character I/O mechanism stems from the ability to insert various modules into a stream to process and manipulate data that passes between a user process and the driver. The following example extends the previous communications device echoing example by inserting a module in the stream to change the case of certain alphabetic characters. The case converter module is passed an input string and an output string by the user. Any incoming data (from the driver) is inspected for instances of characters in the module's input string and the alphabetic case of all matching characters is changed. Similar actions are taken for outgoing data using the output string. The necessary declarations for this program are shown below:

```
#include <string.h>
#include <fcntl.h>
#include <stropts.h>
/*
 * These defines would typically be
 * found in a header file for the module
 */
#define OUTPUT_STRING 1
#define INPUT_STRING 2

main()
{
    char buf[1024];
    int fd, count;
    struct strioctl strioctl;
```

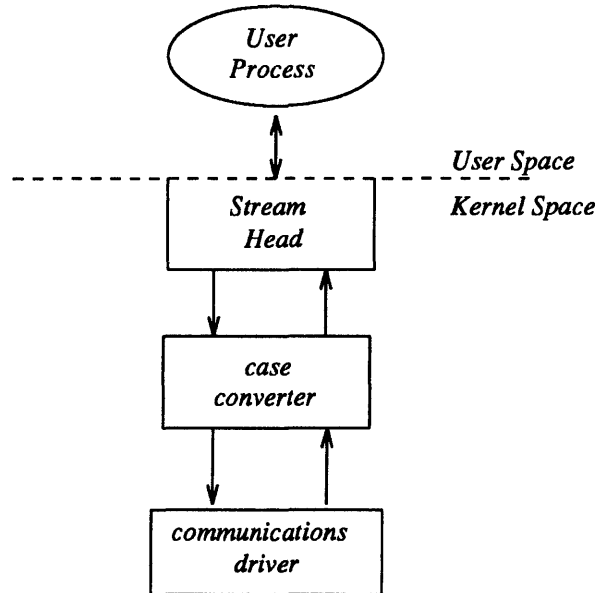
The first step is to establish a stream to the communications driver and insert the case converter module. The following sequence of system calls accomplishes this:

```
if ((fd = open("/dev/comm01", O_RDWR)) < 0) {
    perror("open failed");
    exit(1);
}
if (ioctl(fd, I_PUSH, "case_converter") < 0) {
    perror("ioctl I_PUSH failed");
    exit(2);
}
```

The `I_PUSH ioctl()` call directs the stream head to insert the case converter module between the driver and the stream head, creating the stream shown in the

figure below. As with any driver, this module resides in the kernel and must have been configured into the system before it was booted. `I_PUSH` is one of several generic STREAMS `ioctl()` commands that enable a user to access and control individual streams (see the `streamio(4)` man page).

Figure 2-3 *Case Converter Module*



An important difference between STREAMS drivers and modules is illustrated here. Drivers are accessed through a node or nodes in the file system and may be opened just like any other device. Modules, on the other hand, do not occupy a file system node. Instead, they are identified through a separate naming convention, and are inserted into a stream using `I_PUSH`. The name of a module is defined by the module developer, and is typically included on the manual page describing the module (manual pages describing STREAMS drivers and modules are found in section 7 of the *SunOS Reference Manual*).

Modules are pushed onto a stream and removed from a stream in Last-In-First-Out (LIFO) order. Therefore, if a second module was pushed onto this stream, it would be inserted between the stream head and the case converter module.

Module and Driver Control

The next step in this example is to pass the input string and output string to the case converter module. This can be accomplished by issuing `ioctl()` calls to the case converter module as follows:

```

/* Set input conversion string */
striocntl.ic_cmd = INPUT_STRING;      /* Command type */
striocntl.ic_timeout = 0;              /* Default = 15 sec */
striocntl.ic_dp = "ABCDEFGHJIJ";
striocntl.ic_len = strlen(striocntl.ic_dp);

if (ioctl(fd, I_STR, &striocntl) < 0) {
    perror("ioctl I_STR failed");
    exit(3);
}

/* Set output conversion string */
striocntl.ic_cmd = OUTPUT_STRING;     /* Command type */
striocntl.ic_dp = "abcdefghij";
striocntl.ic_len = strlen(striocntl.ic_dp);

if (ioctl(fd, I_STR, &striocntl) < 0) {
    perror("ioctl I_STR failed");
    exit(4);
}

```

`ioctl()` requests are issued to STREAMS drivers and modules indirectly, using the `I_STR ioctl()` call (see the `streamio(4)` man page). The argument to `I_STR` must be a pointer to a `striocntl` structure, which specifies the request to be made to a module or driver. This structure is defined in `<stropts.h>` and has the following format:

```

struct striocntl {
    int ic_cmd;          /* ioctl request */
    int ic_timeout;     /* ACK/NAK timeout */
    int ic_len;         /* Length of data argument */
    char *ic_dp;        /* Ptr to data argument */
}

```

where `ic_cmd` identifies the command intended for a module or driver, `ic_timeout` specifies the number of seconds an `I_STR` request should wait for an acknowledgement before timing out, `ic_len` is the number of bytes of data to accompany the request, and `ic_dp` points to that data.

`I_STR` is intercepted by the stream head, which packages it into a message, using information contained in the `striocntl` structure, and sends the message downstream. The request will be processed by the module or driver closest to the stream head that understands the command specified by `ic_cmd`. The `ioctl()` call will block up to `ic_timeout` seconds, waiting for the target module or driver to respond with either a positive or negative acknowledgement message. If an acknowledgement is not received in `ic_timeout` seconds, the `ioctl()` call will fail.

`I_STR` is actually a nested request; the stream head intercepts `I_STR` and then sends the driver or module request (as specified in the `striocntl` structure) downstream. Any module that does not understand the command in `ic_cmd` will pass the message further downstream. Eventually, the request will reach the target module or driver, where it is processed and acknowledged. If no module or

driver understands the command, a negative acknowledgement will be generated and the `ioctl()` call will fail.

In the example, two separate commands are sent to the case converter module. The first contains the conversion string for input data, and the second contains the conversion string for output data. The `ic_cmd` field is set to indicate whether the command is setting the input or output conversion string. For each command, the value of `ic_timeout` is set to zero, which specifies the system default timeout value of 15 seconds. Also, a data argument that contains the conversion string accompanies each command. The `ic_dp` field points to the beginning of each string, and `ic_len` is set to the length of the string.

NOTE *Only one I_STR request can be active on a stream at one time. Further requests will block until the active I_STR request is acknowledged and the system call completes.*

The `strioc1` structure is also used to retrieve the results, if any, of an `I_STR` request. If data are returned by the target module or driver, `ic_dp` must point to a buffer large enough to hold that data, and `ic_len` will be set on return to indicate the amount of data returned.

The remainder of this example is identical to the previous example:

```
while ((count = read(fd, buf, 1024)) > 0) {
    if (write(fd, buf, count) != count) {
        perror("write failed");
        break;
    }
}
exit(0);
}
```

The case converter module will convert the specified input characters to lower case, and the corresponding output characters to upper case. Notice that the case conversion processing was realized with *no* change to the communications driver.

As with the previous example, the `exit()` system call will dismantle the stream before terminating the process. The case converter module will be removed from the stream automatically when it is closed. Alternatively, modules may be removed from a stream using the `I_POP` `ioctl()` call described in `streamio(4)`. This call removes the topmost module on the stream, and enables a user process to alter the configuration of a stream dynamically, by pushing and popping modules as needed.

A few of the important `ioctl()` requests supported by STREAMS have been discussed. Several other requests are available to support operations such as determining if a given module exists on the stream, or flushing the data on a stream. These requests are described fully in the `streamio(4)` man page).

2.3. Advanced Operations

This section introduces advanced features provided by STREAMS, such as an I/O polling facility, asynchronous I/O processing support, and a method to sample drivers for available resources.

Advanced Input/Output Facilities

The traditional input/output `open()`, facilities—`close()`, `read()`, `write()`, and `ioctl()`—have been discussed, but STREAMS supports new user capabilities that will be described in the remaining sections of this guide. This section describes a facility that enables a user process to poll multiple streams simultaneously for various events. Also discussed is a signaling feature that supports asynchronous I/O processing. Finally, this section presents a new mechanism for finding available minor devices, called *clone open*.

Input/Output Polling

The `poll(2)` system call provides users with a mechanism for monitoring input and output on a set of file descriptors that reference open streams. It identifies those streams over which a user can send or receive data. For each stream of interest users can specify one or more events about which they should be notified. These events include the following:

POLLIN

Input data are available on the stream associated with the given file descriptor.

POLLPRI

A priority message is available on the stream associated with the given file descriptor. Priority messages are described in the section of Chapter 4 entitled "Accessing the Datagram Provider."

POLLOUT

The stream associated with the given file is writable. That is, the stream has relieved the flow control that would prevent a user from sending data over that stream.

`poll()` will examine each file descriptor for the requested events and, on return, will indicate which events have occurred for each file descriptor. If no event has occurred on any polled file descriptor, `poll()` blocks until a requested event or timeout occurs. The specific arguments to `poll()` are the following:

- an array of file descriptors and events to be polled
- the number of file descriptors to be polled
- the number of milliseconds `poll()` should wait for an event if no events are pending (-1 specifies wait forever)

The following example shows the use of `poll()`. Two separate minor devices of the communications driver presented earlier are opened, thereby establishing two separate streams to the driver. Each stream is polled for incoming data. If data arrives on either stream, it is read and then written back to the other stream. This program extends the previous echoing example by sending echoed data over a separate communications line (minor device). The steps needed to establish each stream are as follows:

```

#include <fcntl.h>
#include <poll.h>

#define NPOLL 2  /* Number of file descriptors to poll */

main()
{
    struct pollfd pollfds[NPOLL];
    char buf[1024];
    int count, i;

    if ((pollfds[0].fd =
        open("/dev/comm01", O_RDWR|O_NDELAY)) < 0) {
        perror("open failed for /dev/comm01");
        exit(1);
    }
    if ((pollfds[1].fd =
        open("/dev/comm02", O_RDWR|O_NDELAY)) < 0) {
        perror("open failed for /dev/comm02");
        exit(2);
    }
}

```

The variable *pollfds* is declared as an array of *pollfd* structures, where this structure is defined in *<poll.h>* and has the following format:

```

struct pollfd {
    int     fd;          /* File descriptor */
    short   events;     /* Requested events */
    short   revents;    /* Returned events */
}

```

For each entry in the array, *fd* specifies the file descriptor to be polled and *events* is a bitmask that contains the bitwise inclusive OR of events to be polled on that file descriptor. On return, the *revents* bitmask will indicate which of the requested events has occurred.

The example opens two separate minor devices of the communications driver and initializes the *pollfds* entry for each. The remainder of the example uses *poll()* to process incoming data as follows:

```

/* Set events to poll for incoming data */
pollfds[0].events = POLLIN;
pollfds[1].events = POLLIN;

while (1) {
    /* Poll and use -1 timeout (infinite) */
    if (poll(pollfds, NPOLL, -1) < 0) {
        perror("poll failed");
        exit(3);
    }
    for (i = 0; i < NPOLL; i++) {
        switch (pollfds[i].revents) {
            default: /* Default err case */

```

```

        perror("error event");
        exit(4);
    case 0:                /* No events */
        break;
    case POLLIN:
        /* Echo incoming data on "other" stream */
        while ((count =
            read(pollfds[i].fd, buf, 1024)) > 0)
            /*
             * write loses data if flow control
             * prevents transmit at this time.
             */
            if (write((i == 0 ?
                pollfds[1].fd: pollfds[0].fd),
                buf, count) != count)
                fprintf(stderr, "write lost data\n");
            break;
        }
    }
}

```

The user specifies the polled events by setting the *events* field of the `pollfd` structure to `POLLIN`. This requested event directs `poll()` to notify the user of any incoming data on each stream. The bulk of the example is an infinite loop, where each iteration will poll both streams for incoming data.

The second argument to `poll()` specifies the number of entries in the `pollfds` array (2 in this example). The third argument is a timeout value indicating the number of milliseconds `poll()` should wait for an event if none has occurred. On a system where millisecond accuracy is not available, *timeout* is rounded up to the nearest legal value available on that system. Here, the value of *timeout* is -1, specifying that `poll()` should block indefinitely until a requested event occurs or until the call is interrupted.

If `poll()` succeeds, the program looks at each entry in `pollfds`. If *revents* is set to 0, no event has occurred on that file descriptor. If *revents* is set to `POLLIN`, incoming data are available. In this case, all available data are read from the polled minor device and written to the other minor device.

If *revents* is set to a value other than 0 or `POLLIN`, an error event must have occurred on that stream, because the only requested event was `POLLIN`. The following error events are defined for `poll()`. These events may not be polled for by the user, but will be reported in *revents* whenever they occur. As such, they are only valid in the *revents* bitmask:

POLLERR

A fatal error has occurred in some module or driver on the stream associated with the specified file descriptor. Further system calls will fail.

POLLHUP

A hangup condition exists on the stream associated with the specified file descriptor.

POLLNVAL

The specified file descriptor is not associated with an open stream.

The example attempts to process incoming data as quickly as possible. However, when writing data to a stream, the `write()` call may block if the stream is exerting flow control. To prevent the process from blocking, the minor devices of the communications driver were opened with the `O_NDELAY` flag set. If flow control is exerted and `O_NDELAY` is set, `write()` will not be able to send all the data. This can occur if the communications driver is unable to keep up with the user's rate of data transmission. If the stream becomes full, the number of bytes `write()` sends will be less than the requested *count*. For simplicity, the example ignores the data if the stream becomes full, and a warning is printed to `stderr`.

This program will continue until an error occurs on a stream, or until the process is interrupted.

Asynchronous Input/Output

The `poll()` system call described above enables a user to monitor multiple streams in a synchronous fashion. The `poll()` call normally blocks until an event occurs on any of the polled file descriptors. In some applications, however, it is desirable to process incoming data asynchronously. For example, an application may wish to do some local processing and be interrupted when a pending event occurs. Some time-critical applications cannot afford to block, but must have immediate indication of success or failure.

A new facility is available for use with STREAMS that enables a user process to request a signal when a given event occurs on a stream. When used with `poll()`, this facility enables applications to asynchronously monitor a set of file descriptors for events.

The `I_SETSIG ioctl()` call (see the `streamio(4)` man page) is used to request that a `SIGPOLL` signal be sent to a user process when a specific event occurs. Listed below are the events for which an application may be signaled:

S_INPUT

Data has arrived at the stream head, and no data existed at the stream head when it arrived.

S_HIPRI

A priority STREAMS message has arrived at the stream head.

S_OUTPUT

The stream is no longer full and can accept output. That is, the stream has relieved the flow control that would prevent a user from sending data over that stream.

S_MSG

A special STREAMS signal message that contains a `SIGPOLL` signal has reached the front of the stream head input queue. This message may be sent by modules or drivers to generate immediate notification of data or events to follow.

The polling example could be written to process input from each communications driver minor device by issuing `I_SETSIG` to request a signal for the `S_INPUT` event on each stream. The signal catching routine could then call `poll()` to determine on which stream the event occurred. The default action for `SIGPOLL` is to terminate the process. Therefore, the user process must catch the signal using `signal(2)`. `SIGPOLL` will only be sent to processes that request the signal using `I_SETSIG`.

Clone Open

In the earlier examples, each user process connected a stream to a driver by opening a particular minor device of that driver. Often, however, a user process wants to connect a new stream to a driver regardless of which minor device is used to access the driver.

In the past, this typically forced the user process to poll the various minor device nodes of the driver for an available minor device. To alleviate this task, a facility called *clone open* is supported for STREAMS drivers. If a STREAMS driver is implemented as a cloneable device, a single node in the file system may be opened to access any unused minor device. This special node guarantees that the user will be allocated a separate stream to the driver on every `open()` call. Each stream will be associated with an unused minor device, so the total number of streams that may be connected to a cloneable driver is limited by the number of minor devices configured for that driver.

The clone device may be useful, for example, in a networking environment where a protocol pseudo-device driver requires each user to open a separate stream over which it will establish communication. Typically, the users would not care which minor device they used to establish a stream to the driver. Instead, the clone device can find an available minor device for each user and establish a unique stream to the driver. Chapter 3 describes this type of transport protocol driver.

NOTE A user program has no control over whether a given driver supports the *clone open*. The decision to implement a STREAMS driver as a cloneable device is made by the designers of the device driver.

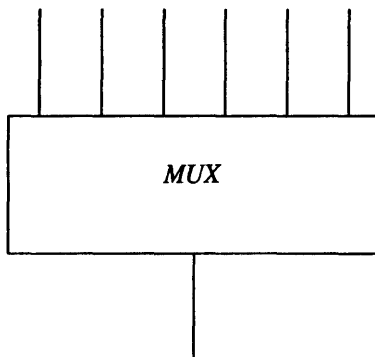
2.4. Multiplexed Streams

This section describes the construction of multiplexed stream configurations.

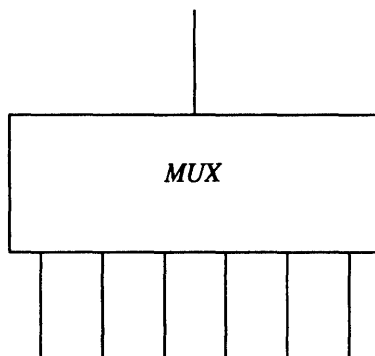
Multiplexor Configurations

In the earlier sections, streams were described as linear connections of modules, where each invocation of a module is connected to at most one upstream module and one downstream module. While this configuration is suitable for many applications, others require the ability to multiplex streams in a variety of configurations. Typical examples are terminal window facilities, and internet-working protocols (which might route data over several subnetworks).

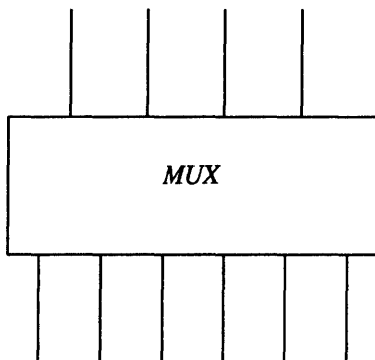
An example of a multiplexor is one that multiplexes data from several upper streams over a single lower stream, as shown in the figure below. An upper stream is one that is upstream from a multiplexor, and a lower stream is one that is downstream from a multiplexor. A terminal windowing facility might be implemented in this fashion, where each upper stream is associated with a separate window.

Figure 2-4 *Many-to-one Multiplexor*

A second type of multiplexor might route data from a single upper stream to one of several lower STREAMS, as shown in the figure below. An internetworking protocol could take this form, where each lower stream links the protocol to a different physical network.

Figure 2-5 *One-to-many Multiplexor*

A third type of multiplexor might route data from one of many upper streams to one of many lower streams, as shown in the figure below.

Figure 2-6 *Many-to-many Multiplexor*

A STREAMS mechanism is available that supports the multiplexing of streams through special pseudo-device drivers. Using a linking facility, users can

dynamically build, maintain, and dismantle each of the above multiplexed stream configurations. In fact, these configurations can be further combined to form complex, multi-level multiplexed stream configurations.

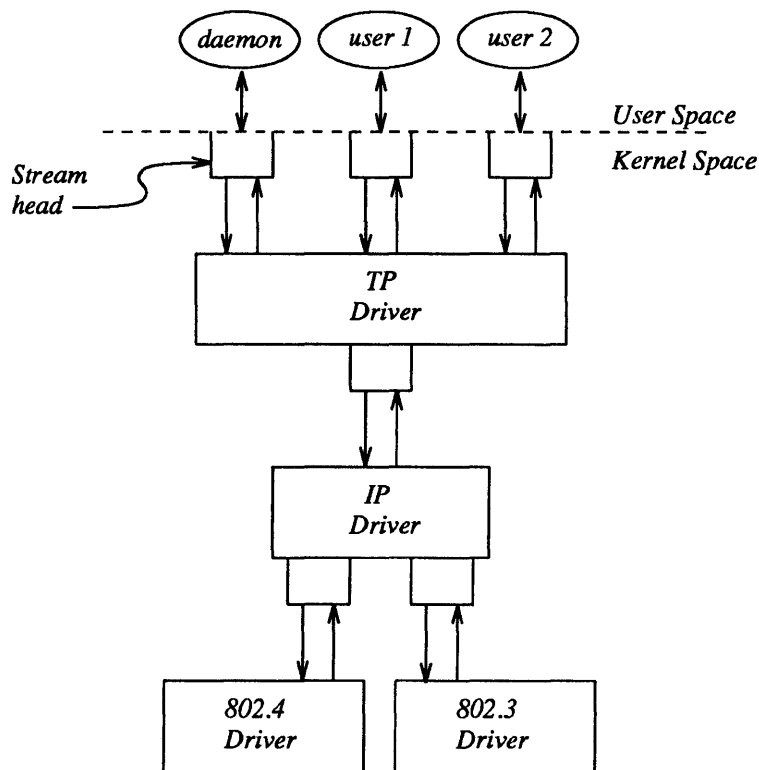
The remainder of this section describes multiplexed stream configurations in the context of an example (see figure below). In this example, an internetworking protocol pseudo-device driver (IP) is used to route data from a single upper stream to one of two lower streams. This driver supports two STREAMS connections beneath it to two distinct sub-networks. One sub-network supports the IEEE 802.3 standard for the CSMA/CD medium access method. The second sub-network supports the IEEE 802.4 standard for the token-passing bus medium access method.

The example also presents a transport protocol pseudo-device driver (TP) that multiplexes multiple virtual circuits (upper streams) over a single stream to the IP pseudo-device driver.

Building a Multiplexor

The figure below shows the multiplexing configuration to be created. This configuration will enable users to access the services of the transport protocol. To free users from the need to know about the underlying protocol structure, a user-level daemon process will build and maintain the multiplexing configuration. Users can then access the transport protocol directly by opening the TP driver device node.

Figure 2-7 *Protocol Multiplexor*



The following example shows how this daemon process sets up the protocol

multiplexor. The necessary declarations and initialization for the daemon program are as follows:

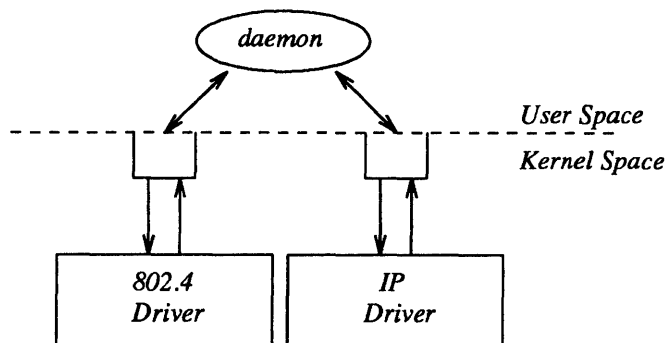
```
#include <fcntl.h>
#include <stropts.h>

main()
{
    int fd_802_4,
        fd_802_3,
        fd_ip,
        fd_tp;
    /* Daemon-ize this process */
    switch (fork()) {
    case 0:
        break;
    case -1:
        perror("fork failed");
        exit(2);
    default:
        exit(0);
    }
    setpgrp();
}
```

This multi-level multiplexed stream configuration will be built from the bottom up. Therefore, the example begins by constructing the IP multiplexor. This multiplexing pseudo-device driver is treated like any other software driver. It owns a node in the file system and is opened just like any other STREAMS device driver.

The first step is to open the multiplexing driver and the 802.4 driver, creating separate streams above each driver as shown in the figure below. The stream to the 802.4 driver may now be connected below the multiplexing IP driver using the `I_LINK` `ioctl()` call.

Figure 2-8 *Before Link*



The sequence of instructions to this point is:

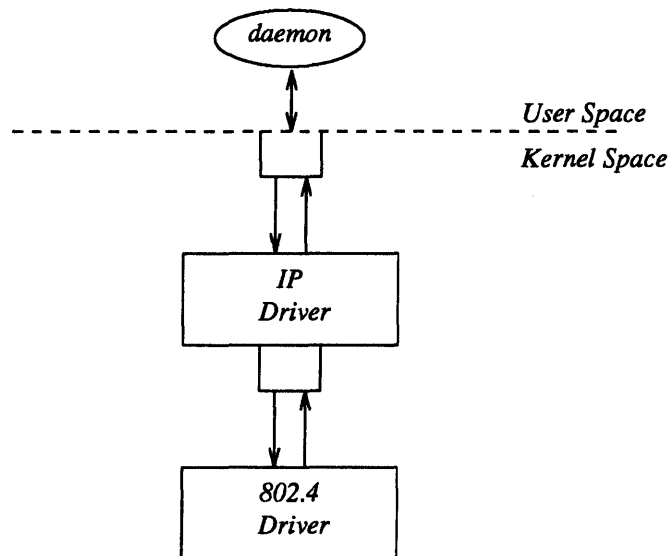
```

if ((fd_802_4 = open("/dev/802_4", O_RDWR)) < 0) {
    perror("open of /dev/802_4 failed");
    exit(1);
}
if ((fd_ip = open("/dev/ip", O_RDWR)) < 0) {
    perror("open of /dev/ip failed");
    exit(2);
}
/* Now link 802.4 to underside of IP */
if (ioctl(fd_ip, I_LINK, fd_802_4) < 0) {
    perror("I_LINK ioctl failed");
    exit(3);
}

```

`I_LINK` takes two file descriptors as arguments. The first file descriptor, `fd_ip`, must reference the stream connected to the multiplexing driver, and the second file descriptor, `fd_802_4`, must reference the stream to be connected below the multiplexor. The figure below shows the state of these streams following the `I_LINK` call. The complete stream to the 802.4 driver has been connected below the IP driver, including the stream head. The stream head of the 802.4 driver will be used by the IP driver to manage the multiplexor.

Figure 2-9 *IP Multiplexor After First Link*



`I_LINK` will return an integer value, called a mux id, which is used by the multiplexing driver to identify the stream just connected below it. This mux id is ignored in the example, but may be useful for dismantling a multiplexor or routing data through the multiplexor. Its significance is discussed later.

The following sequence of system calls is used to continue building the internet-working multiplexor (IP):

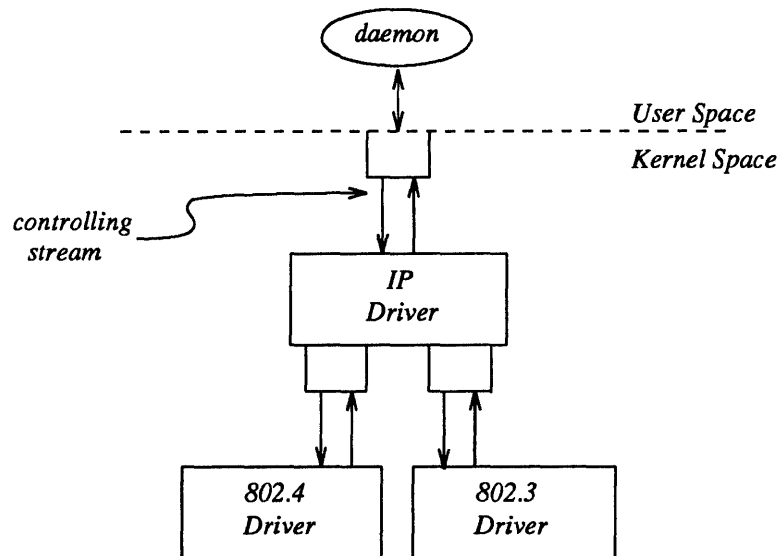
```

if ((fd_802_3 = open("/dev/802_3", O_RDWR)) < 0) {
    perror("open of /dev/802_3 failed");
    exit(4);
}
if (ioctl(fd_ip, I_LINK, fd_802_3) < 0) {
    perror("I_LINK ioctl failed");
    exit(5);
}

```

All links below the IP driver have now been established, giving the configuration in the figure below.

Figure 2-10 *IP Multiplexor*



The stream above the multiplexing driver used to establish the lower connections is the controlling stream and has special significance when dismantling the multiplexing configuration, as will be illustrated later in this section. The stream referenced by *fd_ip* is the controlling stream for the IP multiplexor.

NOTE *The order in which the streams in the multiplexing configuration are opened is unimportant. If, however, it is necessary to have intermediate modules in the stream between the IP driver and media drivers, these modules must be added to the streams associated with the media drivers (using *I_PUSH*) before the media drivers are attached below the multiplexor.*

The number of streams that can be linked to a multiplexor is restricted by the design of the particular multiplexor. The manual page describing each driver (typically found in section 4 of the *SunOS Reference Manual*) should describe such restrictions. However, only one *I_LINK* operation is allowed for each lower stream; a single stream cannot be linked below two multiplexors simultaneously.

Continuing with the example, the IP driver will now be linked below the transport protocol (TP) multiplexing driver. As seen earlier in the figure below, only one link will be supported below the transport driver. This link is formed by the following sequence of system calls:

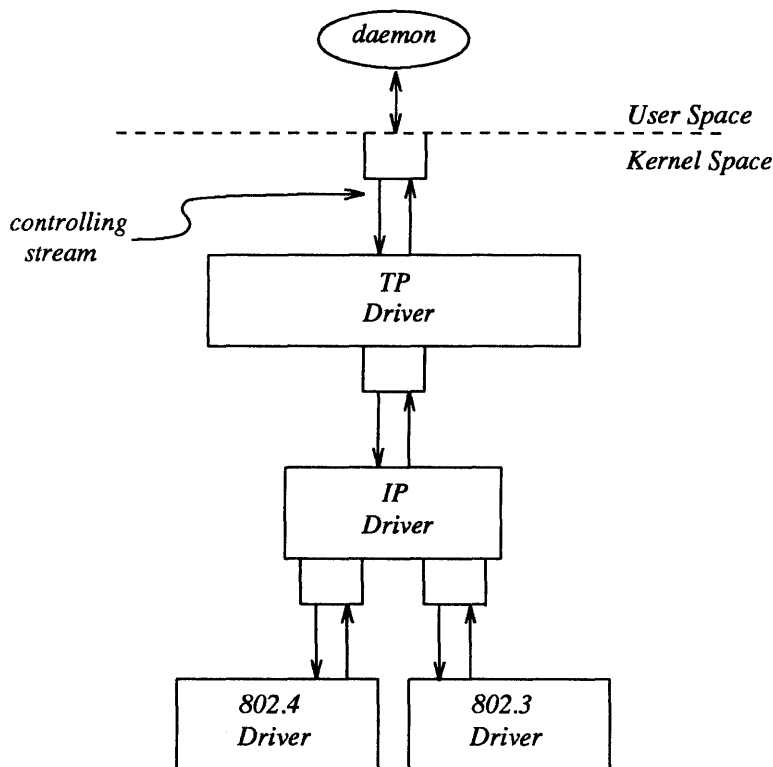
```

if ((fd_tp = open("/dev/tp", O_RDWR)) < 0) {
    perror("open of /dev/tp failed");
    exit(6);
}
if (ioctl(fd_tp, I_LINK, fd_ip) < 0) {
    perror("I_LINK ioctl failed");
    exit(7);
}

```

The multi-level multiplexing configuration shown in the figure below has now been created.

Figure 2-11 *TP Multiplexor*



Because the controlling stream of the IP multiplexor has been linked below the TP multiplexor, the controlling stream for the new multi-level multiplexor configuration is the stream above the TP multiplexor.

At this point the file descriptors associated with the lower drivers can be closed without affecting the operation of the multiplexor. Closing these file descriptors may be necessary when building large multiplexors, so that many devices can be linked together without exceeding the system limit on the number of

simultaneously open files per process. If these file descriptors are not closed, all subsequent `read()`, `write()`, `ioctl()`, `poll()`, `getmsg()`, and `putmsg()` system calls issued to them will fail. That is because `I_LINK` associates the stream head of each linked stream with the multiplexor, so the user may not access that stream directly for the duration of the link.

The following sequence of system calls will complete the multiplexing daemon example:

```
close(fd_802_4);
close(fd_802_3);
close(fd_ip);
/* Hold multiplexor open forever */
pause();
}
```

The figure below shows the complete picture of the multi-level protocol multiplexor. The transport driver is designed to support several, simultaneous virtual circuits, where these virtual circuits map one-to-one to streams opened to the transport driver. These streams will be multiplexed over the single stream connected to the IP multiplexor. The mechanism for establishing multiple streams above the transport multiplexor is actually a by-product of the way in which streams are created between a user process and a driver. By opening different minor devices of a STREAMS driver, separate streams will be connected to that driver. Of course, the driver must be designed with the intelligence to route data from the single lower stream to the appropriate upper stream.

Notice in the figure below that the daemon process maintains the multiplexed stream configuration through an open stream (the controlling stream) to the transport driver. Meanwhile, other users can access the services of the transport protocol by opening new streams to the transport driver; they are freed from the need for any unnecessary knowledge of the underlying protocol configurations and sub-networks that support the transport service.

Multi-level multiplexing configurations, such as the one presented in the above example, should be assembled from the bottom up. That is because STREAMS does not allow `ioctl()` requests (including `I_LINK`) to be passed through higher multiplexing drivers to reach the desired multiplexor; they must be sent directly to the intended driver. For example, once the IP driver is linked under the TP driver, `ioctl()` requests cannot be sent to the IP driver through the TP driver.

Dismantling a Multiplexor

streams connected to a multiplexing driver from above with `open()`, can be dismantled by closing each stream with `close()`. In the protocol multiplexor, these streams correspond to the virtual circuit streams above the TP multiplexor. The mechanism for dismantling streams that have been linked below a multiplexing driver is less obvious, and is described below in detail.

The `I_UNLINK` `ioctl()` call is used to disconnect each multiplexor link below a multiplexing driver individually. This command takes the following form:

```
ioctl(fd, I_UNLINK, mux_id);
```

where *fd* is a file descriptor associated with a stream connected to the multiplexing driver from above, and *mux_id* is the identifier that was returned by `I_LINK` when a driver was linked below the multiplexor. Each lower driver may be disconnected individually in this way, or a special *mux_id* value of -1 may be used to disconnect all drivers from the multiplexor simultaneously.

In the multiplexing daemon program presented earlier, the multiplexor is never explicitly dismantled. That is because all links associated with a multiplexing driver are automatically dismantled when the controlling stream associated with that multiplexor is closed. Because the controlling stream is open to a driver, only the final call of `close()` for that stream will close it. In this case, the daemon is the only process that has opened the controlling stream, so the multiplexing configuration will be dismantled when the daemon exits.

For the automatic dismantling mechanism to work in the multi-level, multiplexed stream configuration, the controlling stream for each multiplexor at each level must be linked under the next higher level multiplexor. In the example, the controlling stream for the IP driver was linked under the TP driver. This resulted in a single controlling stream for the full, multi-level configuration. Because the multiplexing program relied on closing the controlling stream to dismantle the multiplexed stream configuration instead of using explicit `I_UNLINK` calls, the *mux id* values returned by `I_LINK` could be ignored.

An important side effect of automatic dismantling on `close()` is that it is not possible for a process to build a multiplexing configuration and then exit. That is because `exit(2)` will close all files associated with the process, including the controlling stream. To keep the configuration intact, the process must exist for the life of that multiplexor. That is the motivation for implementing the example as a daemon process.

Routing Data Through a Multiplexor

As demonstrated, STREAMS has provided a mechanism for building multiplexed stream configurations. However, the criteria on which a multiplexor routes data is driver dependent. For example, the protocol multiplexor shown in the last example might use address information found in a protocol header to determine over which sub-network a given packet should be routed. It is the multiplexing driver's responsibility to define its routing criteria.

One routing option available to the multiplexor is to use the *mux id* value to determine to which stream data should be routed (remember that each multiplexor link is associated with a *mux id*). `I_LINK` passes the *mux id* value to the driver and returns this value to the user. The driver can therefore specify that the *mux id* value must accompany data routed through it. For example, if a multiplexor routed data from a single upper stream to one of several lower streams (as did the IP driver), the multiplexor could require the user to insert the *mux id* of the desired lower stream into the first four bytes of each message passed to it. The driver could then match the *mux id* in each message with the *mux id* of each lower stream, and route the data accordingly.

2.5. Message Handling

This section describes how to process STREAMS messages in a service interface.

Service Interface Messages

A STREAMS message format has been defined to simplify the design of service interfaces. Also, two new system calls, `getmsg(2)` and `putmsg(2)` are available for sending these messages downstream and receiving messages that are available at the stream head. This section describes these system calls in the context of a service interface example. First, a brief overview of STREAMS service interfaces is presented.

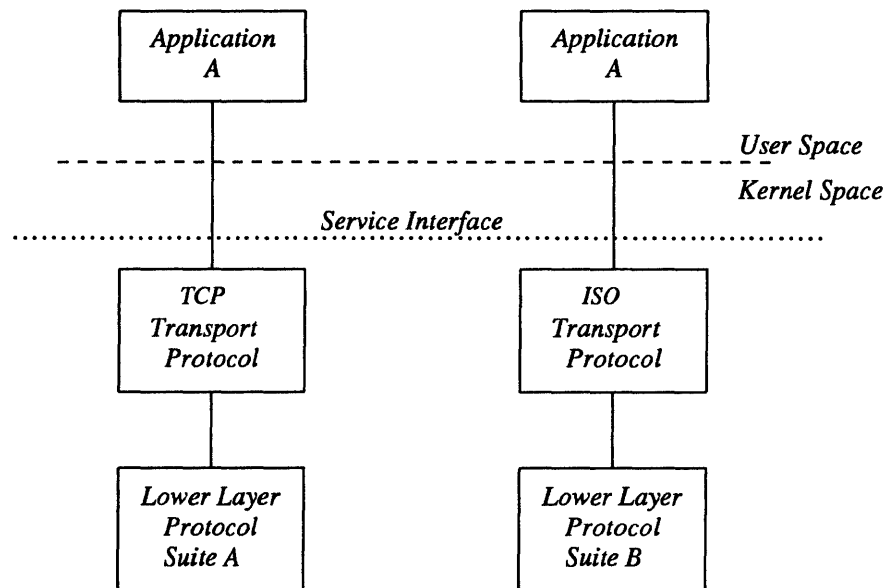
Service Interfaces

A principal advantage of the STREAMS mechanism is its modularity. From user level, kernel-resident modules can be dynamically interconnected to implement any reasonable processing sequence. This modularity reflects the layering characteristics of contemporary network architectures.

One benefit of modularity is the ability to interchange modules of like function. For example, two distinct transport protocols, implemented as STREAMS modules, may provide a common set of services. An application or higher layer protocol that requires those services can use either module. This ability to substitute modules enables user programs and higher level protocols to be independent of the underlying protocols and physical communication media.

Each STREAMS module provides a set of processing functions, or services, and an interface to those services. The service interface of a module defines the interaction between that module and any neighboring modules, and therefore is a necessary component for providing module substitution. By creating a well-defined service interface, applications and STREAMS modules can interact with any module that supports that interface. The figure below demonstrates this.

Figure 2-12 *Protocol Substitution*

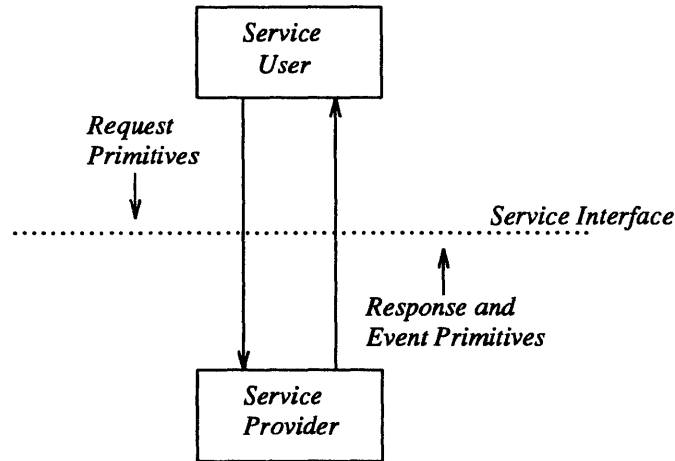


By defining a service interface through which applications interact with a transport protocol, it is possible to substitute a different protocol below that service interface in a manner completely transparent to the application. In this example,

the same application can run over the Transmission Control Protocol (TCP) and the ISO transport protocol. Of course, the service interface must define a set of services common to both protocols.

The three components of any service interface are the service user, the service provider, and the service interface itself, as seen in the figure below.

Figure 2-13 *Service Interface*



Typically, a user makes a request of a service provider using some well-defined service primitive. Responses and event indications are also passed from the provider to the user using service primitives. The service interface is defined as the set of primitives that define a service and the allowable state transitions that result as these primitives are passed between the user and provider.

The Message Interface

A message format has been defined to simplify the design of service interfaces using STREAMS. Each service interface primitive is a distinct STREAMS message that has two parts: a control part and a data part. The control part contains information that identifies the primitive and includes all necessary parameters. The data part contains user data associated with that primitive.

An example of a service interface primitive is a transport protocol connect request. This primitive requests the transport protocol service provider to establish a connection with another transport user. The parameters associated with this primitive may include a destination protocol address and specific protocol options to be associated with that connection. Some transport protocols also allow a user to send data with the connect request. A STREAMS message would be used to define this primitive. The control part would identify the primitive as a connect request and would include the protocol address and options. The data part would contain the associated user data.

STREAMS enables modules to create these messages and pass them to neighbor modules. However, the `read()` and `write()` system calls are not sufficient to enable a user process to generate and receive such messages. First, `read()` and `write()` are byte-stream oriented, with no concept of message boundaries. To support service interfaces, the message boundary of each service primitive must be preserved so that the beginning and end of each primitive can be located.

Also, `read()` and `write()` offer only one buffer to the user for transmitting and receiving STREAMS messages. If control information and data were placed in a single buffer, the user would have to parse the contents of the buffer to separate the data from the control information.

Two new STREAMS system calls are available that enable user processes to create STREAMS messages and send them to neighboring kernel modules and drivers or receive the contents of such messages from kernel modules and drivers. These system calls preserve message boundaries and provide separate buffers for the control and data parts of a message.

The `putmsg()` system call enables a user to create STREAMS messages and send them downstream. The user supplies the contents of the control and data parts of the message in two separate buffers. Likewise, the `getmsg()` system call retrieves such messages from a stream and places the contents into two user buffers.

The syntax of `putmsg()` is as follows:

```
int putmsg (fd, ctlptr, dataptr, flags)
int fd;
struct strbuf *ctlptr;
struct strbuf *dataptr;
int flags;
```

fd identifies the stream to which the message will be passed, *ctlptr* and *dataptr* identify the control and data parts of the message, and *flags* may be used to specify that a priority message should be sent.

The `strbuf` structure is used to describe the control and data parts of a message, and has the following format:

```
struct strbuf {
    int maxlen;           /* Maximum buffer length */
    int len;             /* Length of data */
    char *buf;          /* Pointer to buffer */
}
```

buf points to a buffer containing the data and *len* specifies the number of bytes of data in the buffer. *maxlen* specifies the maximum number of bytes the given buffer can hold, and is only meaningful when retrieving information into the buffer using `getmsg()`.

The `getmsg()` system call retrieves messages available at the stream head, and has the following syntax:

```
int getmsg (fd, ctlptr, dataptr, flags)
int fd;
struct strbuf *ctlptr;
struct strbuf *dataptr;
int *flags;
```

The arguments to `getmsg()` are the same as those for `putmsg()`.

The remainder of this section presents an example that demonstrates how `putmsg()` and `getmsg()` may be used to interact with the service interface of a simple datagram protocol provider. A potential provider of such a service might be the IEEE 802.2 Logical Link Control Protocol Type 1. The example implements a user level library that would free the user from knowledge of the underlying STREAMS system calls. The Transport Interface of the Network Services Library in UNIX System V Release 3.0 provides a similar function for transport layer services. The example here illustrates how a service interface might be defined, and is not an example of a complete IEEE 802.2 service interface.

Datagram Service Interface Example

The example datagram service interface library presented below includes four functions that enable a user to do the following:

- establish a stream to the service provider and bind a protocol address to the stream
- send a datagram to a remote user
- receive a datagram from a remote user
- close the stream connected to the provider

First, the structure and constant definitions required by the library are shown. These typically will reside in a header file associated with the service interface.

```

/*
 * Primitives initiated by the service user.
 */
#define BIND_REQ      1    /* Bind request */
#define UNITDATA_REQ  2    /* Unitdata request */
/*
 * Primitives initiated by the service provider.
 */
#define OK_ACK        3    /* Bind acknowledgment */
#define ERROR_ACK     4    /* Error acknowledgment */
#define UNITDATA_IND  5    /* Unitdata indication */
/*
 * The following structure definitions define the format
 * of the control part of the service interface message
 * of the above primitives.
 */
struct bind_req {          /* Bind request */
    long    PRIM_type;     /* Always BIND_REQ */
    long    BIND_addr;     /* Addr to bind */
};
struct unitdata_req {     /* Unitdata request */
    long    PRIM_type;     /* Always UNITDATA_REQ */
    long    DEST_addr;     /* Destination addr */
};
struct ok_ack {           /* Positive acknowledgment */
    long    PRIM_type;     /* Always OK_ACK */
};
struct error_ack {        /* Error acknowledgment */
    long    PRIM_type;     /* Always ERROR_ACK */
    long    UNIX_error;    /* UNIX error code */
};
struct unitdata_ind {     /* Unitdata indication */
    long    PRIM_type;     /* Always UNITDATA_IND */
    long    SRC_addr;      /* Source addr */
};
union primitives {        /* Union of all primitives */
    long    type;
    struct bind_req    bind_req;
    struct unitdata_req    unitdata_req;
    struct ok_ack    ok_ack;
    struct error_ack    error_ack;
    struct unitdata_ind    unitdata_ind;
};
/* Header files needed by library */
#include <stropts.h>
#include <stdio.h>
#include <errno.h>

```

Five primitives have been defined. The first two represent requests from the service user to the service provider. These are:

BIND_REQ

This request asks the provider to bind a specified protocol address. It requires an acknowledgement from the provider to verify that the contents of the request were syntactically correct.

UNITDATA_REQ

This request asks the provider to send a datagram to the specified destination address. It does not require an acknowledgement from the provider.

The three other primitives represent acknowledgements of requests, or indications of incoming events, and are passed from the service provider to the service user. These are:

OK_ACK

This primitive informs the user that a previous bind request was received successfully by the service provider.

ERROR_ACK

This primitive informs the user that a non-fatal error was found in the previous bind request. It indicates that no action was taken with the primitive that caused the error.

UNITDATA_IND

This primitive indicates that a datagram destined for the user has arrived.

The structures defined above describe the contents of the control part of each service interface message passed between the service user and service provider. The first field of each control part defines the type of primitive being passed.

Accessing the Datagram Provider

The first routine presented below, *inter_open*, opens the protocol driver device file specified by *path* and binds the protocol address contained in *addr* so that it may receive datagrams. On success, the routine returns the file descriptor associated with the open stream; on failure, it returns -1 and sets *errno* to indicate the appropriate error value.

```

inter_open(path, oflags, addr)
char *path;
{
    int fd;
    struct bind_req bind_req;
    struct strbuf ctlbuf;
    union primitives rcvbuf;
    struct error_ack *error_ack;
    int flags;

    if ((fd = open(path, oflags)) < 0)
        return(-1);
    /* Send bind request msg down stream */
    bind_req.PRIM_type = BIND_REQ;
    bind_req.BIND_addr = addr;
    ctlbuf.len = sizeof(struct bind_req);
    ctlbuf.buf = (char *)&bind_req;
    if (putmsg(fd, &ctlbuf, NULL, 0) < 0) {
        close(fd);
        return(-1);
    }
}

```

After opening the protocol driver, *inter_open* packages a bind request message to send downstream. `putmsg()` is called to send the request to the service provider. The bind request message contains a control part that holds a *bind_req* structure, but it has no data part. *ctlbuf* is a structure of type `strbuf`, and it is initialized with the primitive type and address. Notice that the *maxlen* field of *ctlbuf* is not set before calling `putmsg()`. That is because `putmsg()` ignores this field. The *dataptr* argument to `putmsg()` is set to `NULL` to indicate that the message contains no data part. Also, the *flags* argument is 0, which specifies that the message is not a priority message.

After *inter_open* sends the bind request, it must wait for an acknowledgement from the service provider, as follows:

```

/* Wait for ack of request */
ctlbuf.maxlen = sizeof(union primitives);
ctlbuf.len = 0;
ctlbuf.buf = (char *)&rcvbuf;
flags = RS_HIPRI;
if (getmsg(fd, &ctlbuf, NULL, &flags) < 0) {
    close(fd);
    return(-1);
}
/* Did we get enough to determine type */
if (ctlbuf.len < sizeof(long)) {
    close(fd);
    errno = EPROTO;
    return(-1);
}
/* Switch on type (first long in rcvbuf) */
switch(rcvbuf.type) {
    default:
        errno = EPROTO;
        close(fd);
        return(-1);
    case OK_ACK:
        return(fd);
    case ERROR_ACK:
        if (ctlbuf.len < sizeof(struct error_ack)) {
            errno = EPROTO;
            close(fd);
            return(-1);
        }
        error_ack = (struct error_ack *)&rcvbuf;
        errno = error_ack->UNIX_error;
        close(fd);
        return(-1);
}
}

```

`getmsg()` is called to retrieve the acknowledgement of the bind request. The acknowledgement message consists of a control part that contains either an `ok_ack` or `error_ack` structure, and no data part.

The acknowledgement primitives are defined as priority messages. Two classes of messages can arrive at the stream head: priority and normal. Normal messages are queued in a first-in-first-out manner at the stream head, while priority messages are placed at the front of the stream head queue. The STREAMS mechanism allows only one priority message per stream at the stream head at one time; any further priority messages are discarded until the first message is processed. Priority messages are particularly suitable for acknowledging service requests when the acknowledgement should be placed ahead of any other messages at the stream head.

NOTE *These messages are not intended to support the expedited data capabilities of many communication protocols, as evidenced by the one-at-a-time restriction*

just described.

Before calling `getmsg()`, this routine must initialize the `strbuf` structure for the control part. `buf` should point to a buffer large enough to hold the expected control part, and `maxlen` must be set to indicate the maximum number of bytes this buffer can hold.

Because neither acknowledgement primitive contains a data part, the `dataptr` argument to `getmsg()` is set to `NULL`. The `flags` argument points to an integer containing the value `RS_HIPRI`. This flag indicates that `getmsg()` should wait for a STREAMS priority message before returning, and is set because the acknowledgement primitives are priority messages. Even if a normal message is available, `getmsg()` will block until a priority message arrives.

On return from `getmsg()`, the `len` field is checked to ensure that the control part of the retrieved message is an appropriate size. The example then checks the primitive type and takes appropriate actions. An `OK_ACK` indicates a successful bind operation, and `inter_open()` returns the file descriptor of the open stream. An `ERROR_ACK` indicates a bind failure, and `errno` is set to identify the problem with the request.

Closing the Service

The next routine in the datagram service library is `inter_close`, which closes the stream to the service provider.

```
inter_close(fd)
{
    close(fd);
}
```

The routine simply closes the given file descriptor. This will cause the protocol driver to free any resources associated with that stream. For example, the driver may unbind the protocol address that had previously been bound to that stream, thereby freeing that address for use by some other service user.

Sending a Datagram

The third routine, `inter_snd`, passes a datagram to the service provider for transmission to the user at the address specified in `addr`. The data to be transmitted is contained in the buffer pointed to by `buf` and contains `len` bytes. On successful completion, this routine returns the number of bytes of data passed to the service provider; on failure, it returns -1 and sets `errno` to an appropriate system error value.

```

inter_snd(fd, buf, len, addr)
char *buf;
long addr;
{
    struct strbuf ctlbuf;
    struct strbuf databuf;
    struct unitdata_req unitdata_req;

    unitdata_req.PRIM_type = UNITDATA_REQ;
    unitdata_req.DEST_addr = addr;
    ctlbuf.len = sizeof(struct unitdata_req);
    ctlbuf.buf = (char *)&unitdata_req;
    databuf.len = len;
    databuf.buf = buf;
    if (putmsg(fd, &ctlbuf, &databuf, 0) < 0)
        return(-1);
    return(len);
}

```

In this example, the datagram request primitive is packaged with both a control part and a data part. The control part contains a *unitdata_req* structure that identifies the primitive type and the destination address of the datagram. The data to be transmitted is placed in the data part of the request message.

Unlike the bind request, the datagram request primitive requires no acknowledgment from the service provider. In the example, this choice was made to minimize the overhead during data transfer. Since datagram services are inherently unreliable, this is a valid design choice. If the `putmsg()` call succeeds, this routine assumes all is well and returns the number of bytes passed to the service provider.

Receiving a Datagram

The final routine in this example, *inter_rcv*, retrieves the next available datagram. *buf* points to a buffer where the data should be stored, *len* indicates the size of that buffer, and *addr* points to a long integer where the source address of the datagram will be placed. On successful completion, *inter_rcv* returns the number of bytes in the retrieved datagram; on failure, it returns -1 and sets the appropriate system error value.


```

inter_rcv(fd, buf, len, addr)
char *buf;
long *addr;
{
    struct strbuf ctlbuf;
    struct strbuf databuf;
    struct unitdata_ind unitdata_ind;
    int retval;
    int flags;

    ctlbuf.maxlen = sizeof(struct unitdata_ind);
    ctlbuf.len = 0;
    ctlbuf.buf = (char *)&unitdata_ind;
    databuf.maxlen = len;
    databuf.len = 0;
    databuf.buf = buf;
    flags = 0;
    if ((retval = getmsg(fd, &ctlbuf, &databuf, &flags)) < 0)
        return(-1);
    if (unitdata_ind.PRIM_type != UNITDATA_IND) {
        errno = EPROTO;
        return(-1);
    }
    if (retval) {
        errno = EIO;
        return(-1);
    }
    *addr = unitdata_ind.SRC_addr;
    return(databuf.len);
}

```

`getmsg()` is called to retrieve the datagram indication primitive, where that primitive contains both a control and data part. The control part consists of a *unitdata_ind* structure that identifies the primitive type and the source address of the datagram sender. The data part contains the data itself.

In *ctlbuf*, *buf* must point to a buffer where the control information will be stored, and *maxlen* must be set to indicate the maximum size of that buffer. Similar initialization is done for *databuf*.

The *flags* argument to `getmsg()` is set to zero, indicating that the next message should be retrieved from the stream head, regardless of its priority. Datagrams will arrive in normal priority messages. If no message currently exists at the stream head, `getmsg()` will block until a message arrives.

The user's control and data buffers should be large enough to hold any incoming datagram. If both buffers are large enough, `getmsg()` will process the datagram indication and return 0, indicating that a full message was retrieved successfully. However, if either buffer is not large enough, `getmsg()` will only retrieve the part of the message that fits into each user buffer. The remainder of the message is saved for subsequent retrieval, and a positive, non-zero value is returned to the user. A return value of `MORECTL` indicates that more control

information is waiting for retrieval. A return value of MOREDATA indicates that more data are waiting for retrieval. A return value of MORECTL | MOREDATA indicates that data from both parts of the message remain. In the example, if the user buffers are not large enough (that is, `getmsg()` returns a positive, non-zero value), the function will set *errno* to EIO and fail.

The type of the primitive returned by `getmsg()` is checked to make sure it is a datagram indication. The source address is then set and the number of bytes of data in the datagram is returned.

The above example presented a simplified service interface. The state transition rules for such an interface were not presented for the sake of brevity. The intent was to show typical uses of the `putmsg()` and `getmsg()` system calls. See `putmsg(2)` and `getmsg(2)` for further details.

STREAMS Module and Driver Programming

3.1. Introduction

This chapter provides detailed information on the use of the STREAMS mechanism at the kernel level, including examples, information on development methods and design philosophy. It describes the use of STREAMS kernel facilities for developing and installing modules and drivers, and is intended for system programmers with knowledge of UNIX kernel programming, device driver development, networking, and other data communication facilities.

Examples are used throughout this chapter to highlight the most important and common capabilities of STREAMS. The descriptions are not meant to be exhaustive. For simplicity, the examples reference fictional drivers and modules.

The preceding *STREAMS Application Programming* chapter is the companion to this chapter—it provides an analogous discussion of the STREAMS applications level.

Both of these chapters assumes a working knowledge of the material covered in the preceding *Introduction to STREAMS* chapter (hereafter simply called the *Introduction to STREAMS*). This introduction includes a useful Glossary of STREAMS-related terms. STREAMS kernel utilities, functions, and macros are summarized in the *Utilities* section of the *Supplementary STREAMS Material* chapter of this manual. STREAMS system calls are specified in Section 2 of the *SunOS Reference Manual*. The STREAMS modules and drivers available with SunOS are described in section 4 of the *SunOS Reference Manual*. STREAMS-specific `ioctl()` calls are specified in `streamio(4)`.

The STREAMS facility was incorporated into SunOS to augment the existing character input/output (I/O) mechanism and to support the development of communication services. A STREAMS driver may be a device driver that provides the services of an external I/O device, or a software driver, commonly referred to as a pseudo-device driver, that performs functions internal to a stream. The Stream Head provides the interface between the stream and user processes. Its principal function is to process STREAMS-related user system calls so that this processing does not have to be incorporated in a module and driver.

Data is passed between a driver and the Stream Head in messages. Messages that are passed from the Stream Head toward the driver are said to travel downstream. Similarly, messages passed in the other direction travel upstream. The Stream Head transfers data between the data space of a user process and kernel data space. Data to be sent to a driver from a user process are packaged into STREAMS messages and passed downstream. When a message containing data

arrives at the Stream Head from downstream, the message is processed by the Stream Head, which copies the data into user buffers.

Within a stream, messages are distinguished by a type indicator. Certain message types sent upstream may cause the Stream Head to perform specific actions, such as sending a signal to a user process. Other message types are intended to carry information within a stream and are not directly seen by a user process.

One or more kernel-resident modules may be inserted into a stream between the Stream Head and driver to perform intermediate processing of data as it passes between the Stream Head and driver. STREAMS modules are dynamically interconnected in a stream by a user process. No kernel programming, assembly, or link editing is required to create the interconnection.

Development Facilities

General and STREAMS-specific system calls provide the user level facilities required to implement application programs. This system call interface is upwardly compatible with the character I/O facilities. The `open(2)` system call recognizes a STREAMS file and creates a stream to the specified driver. A user process can receive and send data on STREAMS files using `read(2)` and `write(2)` in the same manner as with character files. The `ioctl(2)` system call enables users to perform functions specific to a particular device and a set of generic STREAMS `ioctl()` commands (see `streamio(4)`) support a variety of functions for accessing and controlling streams. A `close(2)` will dismantle a stream.

In addition to the generic `ioctl()` commands, there are STREAMS-specific system calls to support unique STREAMS facilities. The `putmsg(2)` and `getmsg(2)` system calls enable users to send and receive STREAMS messages, and are suitable for interacting with STREAMS modules and drivers through a service interface.

STREAMS provides module and driver developers with integral functions, a set of utility routines, and facilities that expedite design and implementation. The principal development facilities are:

- Message storage management – to maintain STREAMS' own memory resources for message storage
- Flow control – to conserve STREAMS memory and processing resources
- Scheduling – to control the execution of service procedures
- Multiplexing – to switch data among multiple streams

3.2. STREAMS Mechanism

A stream implements a connection within the kernel between a driver in kernel space and a process in user space. It provides a general character input/output (I/O) interface for user processes that is upwardly compatible with the interface of the preexisting character I/O facilities. A stream is analogous to a shell pipeline except that data flow and processing are bidirectional to support concurrent input and output.

The components that form a stream are the Stream Head, driver, and optional modules. A stream is initially constructed as the result of a user process

`open(2)` system call referencing a STREAMS file. The call causes a kernel-resident driver to be connected with a Stream Head to form a stream. Subsequent `ioctl(2)` calls select kernel-resident modules and cause them to be inserted in the stream. A module represents intermediate processing on messages flowing between the Stream Head and driver. A module can function as, for example, a communication protocol, line discipline, or data filter. STREAMS allows a user to connect a module with any other module. The user determines the module connection sequences that result in useful configurations.

A process can send and receive characters on a stream using `write(2)` and `read(2)`, as on character files. When user data enters the Stream Head or external data enters the driver, the data is placed into messages for transmission on the stream. All data passed on a stream is carried in messages, each having a defined message type identifying the message contents. Internal control and status information is transmitted among modules or between the stream and user process as messages of certain types interleaved on the stream. Modules and drivers can send certain message types to the Stream Head to cause the generation of signals or errors to be received by the user process.

A module is composed of two identical sets of data structures called QUEUES. One QUEUE is for upstream processing and the other is for downstream processing. The processing performed by the two QUEUES is generally independent so that a stream operates in a full-duplex manner. The interface between modules is uniform and simple. Messages flow from module to module. A message from one module is passed to the single entry point of its neighboring module.

The last `close(2)` system call dismantles the stream and closes the file, semantically identical to character I/O drivers.

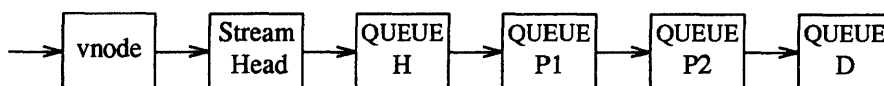
STREAMS supports implementation of user-level applications with extensions to the above general system calls and STREAMS specific system calls:

`putmsg(2)`, `getmsg(2)` and a set of STREAMS generic `ioctl(2)` functions.

Stream Construction

STREAMS constructs a stream as a linked list of kernel-resident data structures. In a STREAMS file, the `vnode` points to the stream header structure. The header is used by STREAMS kernel routines to perform operations on this stream that are generally related to system calls. Figure 3-1 depicts the downstream (write) portion of a stream (see *Building a Stream*, in the *Introduction to STREAMS*) connected to a header. There is one header per stream. From the header onward, a stream is constructed of QUEUES. The upstream (read) portion of the stream (not shown here) parallels the downstream portion in the opposite direction and terminates at the stream header structure.

Figure 3-1 *Downstream Stream Construction*



At the same relative location in each QUEUE is the address of the entry point, a procedure to be executed on any message received by that QUEUE. The

procedure for QUEUE H, at one end of the stream, is the STREAMS-provided Stream Head routine. QUEUE H is the downstream half of the Stream Head. The procedure for QUEUE D, at the other end, is the driver routine. QUEUE D is the downstream half of the stream end. P1 and P2 are pushable modules, each containing their own unique procedures. That is, all STREAMS components are of similar organization.

This similarity results in a uniform manner of navigating in either direction on a stream: messages move from one end to the other, from QUEUE to the next linked QUEUE, executing the procedure specified in the QUEUE.

Figure 3-2 shows the data structures forming each QUEUE: `queue_t`, `qinit`, `module_info`, and `module_stat`. `queue_t` contains various modifiable values for this QUEUE, generally used by STREAMS. `qinit` contains a pointer to the processing procedures, `module_info` contains limit values, and `module_stat` is used for statistics. The two QUEUES in a module will generally each contain a different set of these structures. The contents of these structures are described in following sections.

Figure 3-2 *QUEUE data structures*

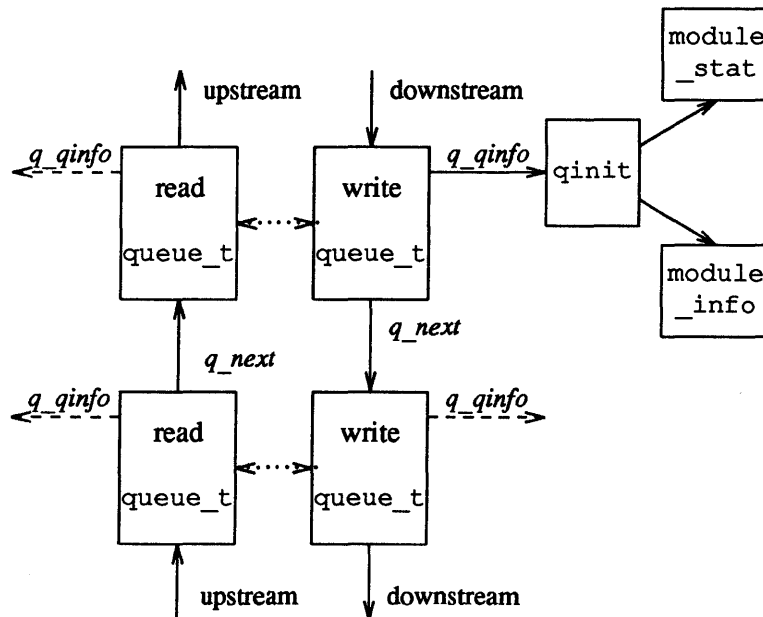


Figure 3-1 shows QUEUE linkage in one direction while figure 3-2 shows two neighboring modules with links (solid vertical arrows) in both directions. When a module is pushed onto a stream, STREAMS creates two QUEUES and links each QUEUE in the module to its neighboring QUEUE in the upstream and downstream direction. The linkage allows each QUEUE to locate its next neighbor. The next relation is implemented between `queue_t`s in adjacent modules by the `q_next` pointer. Within a module, each `queue_t` locates its mate (see dotted arrows in figure 3-2) by use of STREAMS macros, since there is no pointer between the two `queue_t`s. The existence of the Stream Head and driver is known to the QUEUE procedures only as destinations towards which messages

are sent.

Opening a Stream

When a file is opened (see `open(2)`), a STREAMS file is recognized by a non-null value in the `d_str` field of the associated `cdevsw` entry. `d_str` points to a `streamtab` structure:

```
struct streamtab {
    struct qinit    *st_rdinit;    /* defines read QUEUE */
    struct qinit    *st_wrinit;    /* defines write QUEUE */
    struct qinit    *st_muxrinit; /* for multiplexing drivers only */
    struct qinit    *st_muxwinit; /* for multiplexing drivers only */
    char            **st_modlist; /* NULL-terminated list of
                                   modules to be pushed */
}
```

`streamtab` defines a module or driver and points to the read and write `qinit` structures for the driver.

If this `open()` call is the initial file open, a stream is created. First, the single header structure and the Stream Head (see figure 3-1) `queue_t` structure pair are allocated. Their contents are initialized with predetermined values including, as noted above (see `QUEUE H`), the Stream Head processing routines.

Then, a `queue_t` structure pair is allocated for the driver. The STREAMS framework takes the responsibility for maintaining the fields of the `queue_t` structure. It is unusual for a driver or module to need to examine any of these fields, especially with the intention of modifying them. The exception, though, is the `q_ptr` field, which is the means by which a module accesses its own private data.

A single, common `qinit` structure pair is shared among all the streams opened from the same `cdevsw` entry, as are the associated `module_info` and `module_stat` structures (see figure 3-2.)

Next, the `q_next` values are set so that the Stream Head write `queue_t` points to the driver write `queue_t` and the driver read `queue_t` points to the Stream Head read `queue_t`. The `q_next` values at the ends of the stream are set to `NULL`. Then, the driver open procedure (located via `qinit`) is called.

If the `st_modlist` pointer is not `NULL`, it is assumed to point to the first member of a `NULL`-terminated array of pointers to module names. After the driver's `open()` procedure has been called, the modules whose names are pointed to by the members of that array are pushed onto the stream, in the order that they appear in the array. (See *Adding and Removing Modules*, below). If one of these modules cannot be pushed, the `open()` fails.

If this `open()` is not the initial open of this stream, the only actions performed are to call the driver open procedure and the open procedures of all pushable modules on the stream.

Adding and Removing Modules

As part of constructing a stream, a module can be added with an `ioctl()` `I_PUSH` (see the `streamio(4)` man page) system call. The push inserts a module beneath the Stream Head. Because of the similarity of STREAMS components, the push operation is similar to the driver open. First, the address of the `qinit` structure for the module is obtained via an `fmodsw` entry.

`fmodsw` is an array, analogous to `cdevsw`. Each `fmodsw` entry corresponds to a unique module and contains the name of the module (used by `I_PUSH` and certain other STREAMS `ioctl()`s) and a pointer to the module's `streamtab`. Next, STREAMS allocates `queue_t` structures and initializes their contents as in the driver open, above. As with the driver, the read and write `qinit` structures are shared among all the modules opened from this `fmodsw` entry (see figure 3-2).

Then, `q_next` values are set and modified so that the module is interposed between the Stream Head and the driver or module previously connected to the head. Finally, the module open procedure (located via `qinit`) is called. Unlike `open()`, no other module or driver open procedure is called.

Each push of a module is independent, even in the same stream. If the same module is pushed more than once onto a stream, there will be multiple occurrences of that module in the stream. The total number of pushable modules that may be contained on any one stream is limited by the kernel parameter `NSTRPUSH` (see the *SunOS STREAMS Topics* chapter).

An `ioctl()` `I_POP` (see the `streamio(4)` man page) system call (`pop`) removes the module immediately below the Stream Head. The `pop` calls the module's close procedure. On return from the module close, any messages left on the module's message queues are freed (deallocated). Then, STREAMS connects the Stream Head to the component previously below the popped module and deallocates the module's two `queue_t` structures. `I_POP` enables a user process to dynamically alter the configuration of a stream by pushing and popping modules as required. For example, a module may be removed or a new one inserted below a module. In the latter case, the original module is popped and pushed back after the new module has been pushed.

An `I_POP` cannot be used on a driver.

Closing

The last `close()` system call to a STREAMS file dismantles the stream. Dismantling consists of popping any modules on the stream, closing the driver and closing the file. Before a module is popped by `close()`, it may delay to allow any messages on the write message queue of the module to be drained by module processing. If `O_NDELAY` [see `open(2)`] is clear, `close()` will wait up to 15 seconds for each module to drain. If `O_NDELAY` is set, the pop is performed immediately. Note that a module's close routine may also choose to sleep waiting for output to drain. The system call `close()` will then wait for the driver's write queue to drain. Messages can remain queued, for example, if flow control (see *Other Facilities*, in the *Introduction to STREAMS*). is inhibiting execution of the write `QUEUE`. When all modules are popped and any wait for the driver to drain is completed, the driver close routine is called. On return from the driver close, any messages left on the driver's message queues are freed, and the

queue_t and header structures are deallocated.

NOTE *STREAMS frees only the messages contained on a message queue. Any messages used internally by the driver or module must be freed by the driver or module close procedure.*

Finally, the file is closed.

3.3. Modules

Module Declarations

A module and driver will contain, as a minimum, declarations of the following form:

```
#include <sys/types.h>      /* required in all modules and drivers */
#include <sys/stream.h>     /* required in all modules and drivers */
#include <sys/param.h>

static struct module_info rminfo = {0, "mod", 0, INFPSZ, 0, 0};
static struct module_info wminfo = {0, "mod", 0, INFPSZ, 0, 0};
static int modopen(), modrput(), modwput(), modclose();

static struct qinit rinit = {
    modrput, NULL, modopen, modclose, NULL, &rminfo, NULL
};
static struct qinit winit = {
    modwput, NULL, NULL, NULL, NULL, &wminfo, NULL
};
struct streamtab modinfo = { &rinit, &winit, NULL, NULL };
```

The contents of these declarations are constructed for the null module example in this section. This module performs no processing: Its only purpose is to show linkage of a module into the system. The descriptions in this section are general to all STREAMS modules and drivers unless they specifically reference the example.

The declarations shown are: the header set; the read and write QUEUE (rminfo and wminfo) module_info structures (see figure 3-2); the module open, readput, writeput, and close procedures; the read and write (rinit and winit) qinit structures; and the streamtab structure.

The minimum header set for modules and drivers is types.h and stream.h. param.h contains definitions for NULL and other values for STREAMS modules and drivers as shown in the *Accessible Symbols and Functions* section of the *Supplementary STREAMS Material* chapter.

NOTE *Configuring a STREAMS module or driver (see the SunOS STREAMS Topics chapter) does not require any procedures to be externally accessible, only streamtab. The streamtab structure name must be the prefix used in configuring, appended with "info."*

As described in the previous section, `streamtab` contains `qinit` values for the read and write QUEUES, pointing to a `module_info` and an optional `module_stat` structure. The two required structures, shown in figure 3-2), are these:

```

struct qinit {
    int (*qi_putp)(); /* put procedure */
    int (*qi_srvp)(); /* service procedure */
    int (*qi_qopen)(); /* called on each open or a push */
    int (*qi_qclose)(); /* called on last close or a pop */
    int (*qi_qadmin)(); /* reserved for future use */
    struct module_info *qi_minfo; /* information structure */
    struct module_stat *qi_mstat; /* optional statistics structure */
};

struct module_info {
    ushort mi_idnum; /* module ID number */
    char *mi_idname; /* module name */
    short mi_minpsz; /* min packet size accepted, for developer use */
    short mi_maxpsz; /* max packet size accepted, for developer use */
    short mi_hiwat; /* hi-water mark, for flow control */
    ushort mi_lowat; /* lo-water mark, for flow control */
};

```

`qinit` contains the QUEUE procedures. All modules and drivers with the same `streamtab` (i.e., the same `fmodsw` or `cdevsw` entry) point to the same upstream and downstream `qinit` structure(s). The structure is meant to be software read-only, as any changes to it affect all instantiations of that module in all streams. Pointers to the open and close procedures must be contained in the read `qinit`. These fields are ignored in the write side. The example has no service procedure on the read or write side.

`module_info` contains identification and limit values. All modules and drivers with the same `streamtab` point to the same upstream and downstream `module_info` structure(s). As with `qinit`, this structure is intended to be software read-only. However, the four limit values are copied to `queue_t` (see the *Message Queues and Service Procedures* section) where they are modifiable. In the example, the flow control high and low water marks (see the *Drivers* section) are zero since there are no service procedures and messages are not queued in the module.

Three names are associated with a module: the character string in `fmodsw`; the prefix for `streamtab`, used in configuring the module; and the module name field in the `module_info` structure. The module name value used in the `I_PUSH` or other STREAMS `ioctl()` commands is contained in `fmodsw`. Each module ID and module name should be unique in the system. The module ID is currently used only in logging and tracing (see *Other Facilities*, in the *Introduction to STREAMS*). For the example in this section, the module ID is zero.

Minimum and maximum packet size are intended to limit the total number of bytes contained in all (if any) of the `M_DATA` blocks in each message passed to

this QUEUE. These limits are advisory except for the Stream Head. For certain system calls that write to a stream, the Stream Head will observe the packet sizes set in the write QUEUE of the module immediately below it. Otherwise, the use of packet size is developer dependent. In the example, INFPSZ indicates unlimited size on the read (input) side.

`module_stat` is optional, intended for future use. Currently, there is no STREAMS support for statistical information gathering. The structure is described in *Kernel Structures* in the *Supplementary STREAMS Material* chapter.

Module Procedures

The null module procedures are as follows:

```
static int modopen(q, dev, flag, sflag)
    queue_t *q;      /* pointer to read queue */
    dev_t dev;      /* major/minor device number -- zero for modules */
    int flag;       /* file open flags -- zero for modules */
    int sflag;      /* stream open flags */
{
    /* return success */
    return (0);
}

static int modwput(q, mp)      /* write put procedure */
    queue_t *q;      /* pointer to the write queue */
    mblk_t *mp;      /* message pointer */
{
    putnext(q, mp); /* pass message through */
}

static int modrput(q, mp)      /* read put procedure */
    queue_t *q;      /* pointer to the read queue */
    mblk_t *mp;      /* message pointer */
{
    putnext(q, mp); /* pass message through */
}

static int modclose(q, flag)
    queue_t *q;      /* pointer to the read queue */
    int flag;       /* file open flags - zero for modules */
{
}
```

The form and arguments of these four procedures are the same in all modules and all drivers. Modules and drivers can be used in multiple streams and their procedures must be reentrant.

`modopen()` illustrates the open call arguments and return value. The arguments are the read queue pointer (`q`), the major/minor device number (`dev`, in drivers only), the file open flags (`flag`, defined in `sys/file.h`), and the stream open flag (`sflag`). For a module, the values of `flag` and `dev` are always zero. The stream open flag can take on the following values:

MODOPEN

normal module open

0 normal driver open (see the *Drivers* section).**CLONEOPEN**clone driver open (see the *Example Driver* section).

The return value from `open` is ≥ 0 for success and `OPENFAIL` for error. When a driver is opened with `CLONEOPEN` set, the return value should be the minor device number of the device instance that the driver actually opened. The open procedure is called on the first `I_PUSH` and on all subsequent `open()` calls to the same stream. During a push, a return value of `OPENFAIL` causes the `I_PUSH` to fail and the module to be removed from the stream. If `OPENFAIL` is returned by a module during an `open()` call, the `open()` fails, but the stream remains intact. For example, it can be returned by a module/driver that only wishes to be opened by a superuser:

```
if (!suser()) return (OPENFAIL);
```

In the example, `modopen()` simply returns successfully.

`modrput()` and `modwput()` illustrate the common interface to put procedures. The arguments are the read or write `queue_t` pointer, as appropriate, and the message pointer. The put procedure in the appropriate side of the `QUEUE` is called when a message is passed from upstream or downstream. The put procedure has no return value. In the example, no message processing is performed. All messages are forwarded using the `putnext()` macro (see *Utilities* in the *Supplementary STREAMS Material* chapter. `putnext()` calls the put procedure of the next `QUEUE` in the proper direction.

The close procedure is only called on an `I_POP` or on the last `close()` call of the stream (see the last two sections of *STREAMS Mechanism*). The arguments are the read `queue_t` pointer and the file open flags as in `modopen()`. For a module, the value of `flag` is always zero. There is no return value. In the example, `modclose()` does nothing.

Module and Driver Environment

User context is not generally available to `STREAMS` module procedures and drivers. The exception is during execution of the open and close routines. Driver and module open and close routines have user context and may access the `user` structure (defined in `user.h`, see *Accessible Symbols and Functions* in the *Supplementary STREAMS Material*) Appendix. These routines are allowed to sleep, but must always return to the caller. That is, if they sleep, it must be at priority \leq `PZERO`, or with `PCATCH` set in the sleep priority. (A process that is sleeping at priority $>$ `PZERO`, with no `PCATCH`, and is sent a signal via `kill(2)`, never returns from the sleep call. Instead, the system call is aborted.)

`STREAMS` driver and module put procedures and service procedures have no user context. They cannot access the `user` structure of a process and must not sleep.

3.4. Messages

Message Format

Messages are the means of communication within a stream. A message contains data or information identified by one of 18 message types (see *Message Types* in the *Supplementary STREAMS Material* chapter. Messages may be generated by a driver, a module, or the Stream Head. The contents of certain message types can be transferred between a process and a stream by use of system calls. STREAMS maintains its own pools for allocation of message storage.

All messages are composed of one or more message blocks. A message block is a linked triplet; two structures and a variable length buffer block. These structures include the fields described below, along with others that are not documented. The reason some elements have been left undocumented is because they are used by the STREAMS mechanism internally and may change from one release to the next. Note also that programs that count on these structures to have a specific offset or size are not portable and may not work in future releases.

The two structures in particular are `msgb` (`mblk_t`), the message block, and `datab` (`dblck_t`), the data block:

```

/*
 * Fields marked with "RO" are read-only; STREAMS
 * modules and drivers are allowed to examine them, but not to change
 * them (except indirectly through the appropriate utility routines).
 */

struct msgb {
    struct msgb    *b_next; /* RO: next message on queue */
    struct msgb    *b_prev; /* RO: previous message on queue */
    struct msgb    *b_cont; /* next message block of message */
    unsigned char  *b_rptr; /* first unread byte in buffer */
    unsigned char  *b_wptr; /* first unwritten byte in buffer */
    struct datab   *b_datap; /* RO: data block */
};
typedef struct msgb mblk_t;

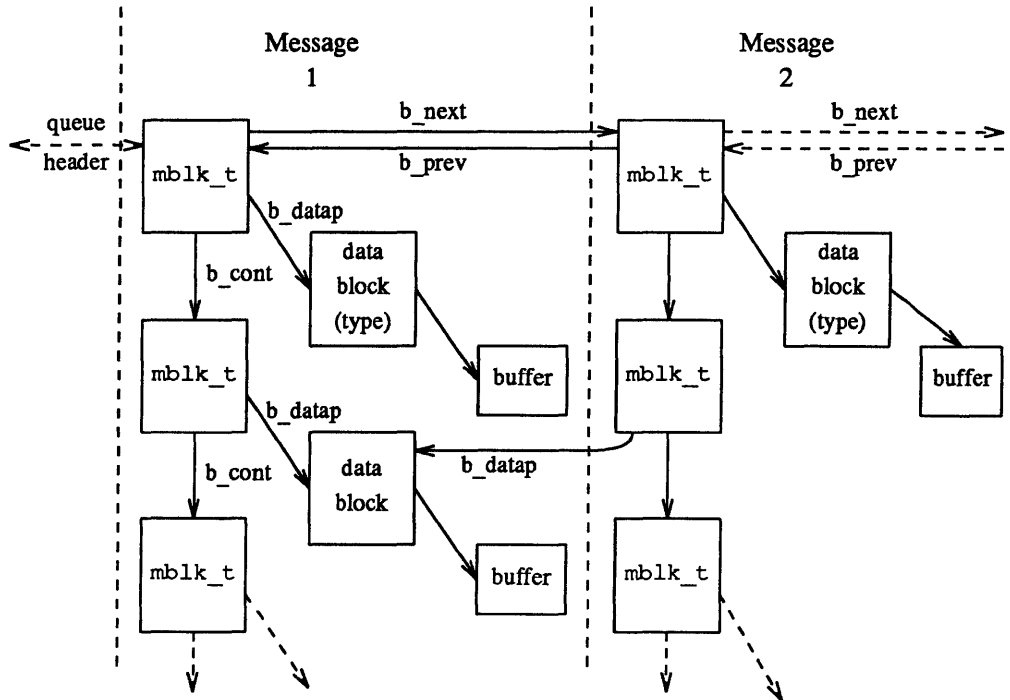
struct datab {
    unsigned char  *db_base; /* RO: first byte of buffer */
    unsigned char  *db_lim; /* RO: last byte+1 of buffer */
    unsigned char  db_ref; /* RO: cnt of msg's pointing to this block */
    unsigned char  db_type; /* message type */
};
typedef struct datab dblck_t;

```

`mblk_t` is used to link messages on a message queue, link the blocks in a message, and manage the reading and writing of the associated buffer. `b_rptr` and `b_wptr` are used to locate the data currently contained in the buffer. As shown in figure 3-3, `mblk_t` points to the data block of the triplet. The data block contains the message type, buffer limits, and control variables. `db_base` and `db_lim` are the fixed beginning and end (+1) of the buffer.

A message consists of one or more linked message blocks. Multiple message blocks in a message can occur, for example, because of buffer size limitations, or as the result of processing that expands the message. When a message is composed of multiple message blocks, the type associated with the first message block determines the overall message type, regardless of the types of the attached message blocks.

Figure 3-3 *Message Form and Linkage*



A message may occur singly, as when it is processed by a put procedure, or it may be linked on the message queue in a QUEUE, generally waiting to be processed by the service procedure. Message 1, as shown in figure 3-3, links to message 2. In the first message on a queue, `b_prev` points back to the header in the QUEUE. The last `b_next` points to the tail.

Note that a data block in message 1 is shared between message 1 and another message. Multiple message blocks can point to the same data block to conserve storage and to avoid copying overhead. For example, the same data block, with associated buffer, may be referenced in two messages, from separate modules that implement separate protocol levels. (Figure 3-3 illustrates the concept, but data blocks would not typically be shared by messages on the same queue). The buffer can be retransmitted, if required by errors or timeouts, from either protocol level without replicating the data. Data block sharing is accomplished by means of a utility routine (see `dupmsg()` in the *Utilities* section of the *Supplementary STREAMS Material* chapter.) STREAMS maintains a count of the message blocks sharing a data block in the `db_ref` field, which modules may examine but not

change.

STREAMS provides utility functions and macros, specified in the *Utilities* section of the *Supplementary STREAMS Material* chapter, to assist in managing messages and message queues, and to assist in other areas of module and driver development. A utility routine should always be used when operating on a message queue or accessing the message storage pool.

Message Generation and Reception

As discussed in the *Introduction to STREAMS*, most message types can be generated by modules and drivers. A few are reserved for the Stream Head. The most commonly used types are `M_DATA`, `M_PROTO`, and `M_PCPROTO`. These, and certain other message types, can also be passed between a process and the topmost module in a stream, with the same message boundary alignment maintained on both sides of the kernel. This allows a user process to function, to some degree, as a module above the stream and maintain a service interface (see the *Service Interface* section). `M_PROTO` and `M_PCPROTO` messages are intended to carry service interface information among modules, drivers, and user processes. Some message types can only be used within a stream and cannot be sent or received from user level.

As discussed previously, modules and drivers do not interact directly with any system calls except `open()` and `close()`. The Stream Head handles all message translation and passing. Message transfer between process and Stream Head can occur in different forms. For example, `M_DATA`, `M_PROTO`, or `M_PCPROTO` messages can be transferred in their direct form by `getmsg(2)` and `putmsg(2)` system calls (see the *Service Interface* section). Alternatively, a `write()` causes one or more `M_DATA` messages to be created from the data buffer supplied in the call. `M_DATA` messages received from downstream at the Stream Head will be consumed by `read(2)` and copied into the user buffer. As another example, `M_SIG` causes the Stream Head to send a signal to a process (see the *Advanced Topics* section).

Any module or driver can send any message type in either direction on a stream. However, based on their intended use in STREAMS and their treatment by the Stream Head, certain message types can be categorized as upstream, downstream, or bidirectional. `M_DATA`, `M_PROTO`, or `M_PCPROTO` messages, for example, can be sent in both directions. Other message types are intended to be sent upstream to be processed only by the Stream Head. Downstream messages are silently discarded if received by the Stream Head.

Filter Module Declarations

The module shown below, `crmod`, is an asymmetric filter. On the write side, newline is converted to carriage return followed by newline. On the read side, no conversion is done. The declarations are essentially the same as the null module of the preceding section:

```

/* Simple filter - converts newline -> carriage return, newline */

#include <sys/types.h>
#include <sys/param.h>
#include <sys/stream.h>

static struct module_info minfo = {0, "crmod", 0, INFPSZ, 0, 0};

static int modopen(), modrput(), modwput(), modclose();
static struct qinit rinit = {
    modrput, NULL, modopen, modclose, NULL, &minfo, NULL
};
static struct qinit winit = {
    modwput, NULL, NULL, NULL, NULL, &minfo, NULL
};
struct streamtab crmdinfo = { &rinit, &winit, NULL, NULL };

```

Note that, in contrast to the null module example, a single `module_info` structure is shared by the read and write sides. The flag `INFPSZ` is defined in `sys/stream.h`.

`modopen()`, `modrput()`, and `modclose()` are the same as in the null module of the preceding section.

bappend() Subroutine

The module makes use of a subroutine, `bappend()`, which appends a character to a message block:

```

*
* Append a character to a message block. If (*bpp) is null, it will allocate
* a new block. Returns 0 when the message block is full, 1 otherwise
*/

#define MODBLKSZ    128 /* convenient buffer size */

static int bappend(bpp, ch)
mblk_t **bpp;
int ch;
{
    mblk_t *bp;

    if (bp = *bpp) {
        if (bp->b_wptr >= bp->b_datap->db_lim)
            return (0);
    } else {
        if ((*bpp=bp=allocb(MODBLKSZ, BPRI_MED)) == NULL)
            return (1);
        bp->b_datap->db_type = M_DATA;
    }
    *bp->b_wptr++ = ch;
    return (1);
}

```


`bappend()` receives a pointer to a message block pointer and a character as arguments. If a message block is supplied (`*bpb != NULL`), `bappend()` checks if there is room for more data in the block. If not, it fails. If there is no message block, a block with capacity `MODBLKSZ` is allocated through `allocb()`, described below.

If the `allocb()` fails, `bappend()` returns success, silently discarding the character. This may or may not be acceptable. For tty-type devices, it is generally acceptable. If the original message block is not full or the `allocb()` is successful, `bappend()` stores the character in the block.

Message Allocation

The `allocb()` utility (see the *Utilities* section of the *Supplementary STREAMS Material* chapter) is used to allocate message storage from the STREAMS pool. Its declaration is:

```
mblk_t *allocb(bufferize, priority)
```

`allocb()` will return a message block containing a buffer of at least the size requested, providing there is a buffer available at the message pool priority specified, or it will return `NULL` on failure.

Three levels of message pool priority can be specified (see the *Utilities* section of the *Supplementary STREAMS Material* chapter). Priority generally does not affect `allocb()` until the pool approaches depletion. In this case, for the same internal level of pool resources, `allocb()` will reject low priority requests while granting higher priority requests. This allows module and driver developers to use STREAMS memory resources to their best advantage and for the common good of the system. (Note that the current implementation disregards priority information, but this may change in future implementations.) Message pool priority does not affect subsequent handling of the message by STREAMS. `BPRI_HI` is intended for special situations, such as the transmission of urgent messages relating to time sensitive events, conditions that could result in loss of state, loss of data, or inability to recover. `BPRI_MED` might be used, for example, when requesting an `M_DATA` buffer for holding input, and `BPRI_LO` might be used for an output buffer (presuming the output data can wait in user space). The Stream Head uses `BPRI_LO` to allocate messages to contain output from a process (e.g., by `write()` or `putmsg()`).

Note that `allocb()` will always return a message of type `M_DATA`. The type may then be changed if required. `b_rptr` and `b_wptr` are set to `db_base` (see `mblk_t` and `dbl_t`).

`allocb()` may return a buffer larger than the size requested. In `bappend()`, if the message block contents were intended to be limited to `MODBLKSZ`, a check would have to be inserted.

If `allocb()` indicates buffers are not available, the `bufcall()` utility can be used to defer processing in the module or the driver until a buffer becomes available (`bufcall()` is described in the *Advanced Topics* section).

Put Procedure

modwput () processes all the message blocks in any downstream data (type M_DATA) messages.

```

/* Write side put procedure */
static modwput(q, mp)
queue_t *q;
mblk_t *mp;
{
    switch (mp->b_datap->db_type) {
    default:
        putnext(q, mp); /* Don't do these, pass them along */
        break;

    case M_DATA: {
        register mblk_t *bp;
        struct mblk_t *nmp = NULL, *nbp = NULL;

        for (bp = mp; bp != NULL; bp = bp->b_cont) {
            while (bp->b_rptr < bp->b_wptr) {
                if (*bp->b_rptr == '\n')
                    if (!bappend(&nbp, '\r'))
                        goto newblk;
                if (!bappend(&nbp, *bp->b_rptr))
                    goto newblk;

                bp->b_rptr++;
                continue;

            newblk:
                if (nmp == NULL)
                    nmp = nbp;

                /* link message block to tail of nmp */
                else
                    linkb(nmp, nbp);

                nbp = NULL;
            }
        }

        if (nmp == NULL)
            nmp = nbp;
        else linkb(nmp, nbp);
        freemsg(mp); /* de-allocate message */
        if (nmp)
            putnext(q, nmp);
        break;
    }
}

```

Data messages are scanned and filtered. modwput () copies the original message into a new block(s), modifying as it copies. nbp points to the current new

message block. `nmp` points to the new message being formed as multiple `M_DATA` message blocks. The outer `for()` loop goes through each message block of the original message. The inner `while()` loop goes through each byte. `bappend()` is used to add characters to the current or new block. If `bappend()` fails, the current new block is full. If `nmp` is `NULL`, `nmp` is pointed at the new block. If `nmp` is non-`NULL`, the new block is linked to the end of `nmp` by use of the `linkb` utility.

At the end of the loops, the final new block is linked to `nmp`. The original message (all message blocks) is returned to the pool by `freemsg()`. If a new message exists, it is sent downstream.

3.5. Message Queues and Service Procedures

The `queue_t` Structure

Service procedures, message queues and priority, and basic flow control are all intertwined in STREAMS. A `QUEUE` will generally not use its message queue if there is no service procedure in the `QUEUE`. The function of a service procedure is to process messages on its queue. Message priority and flow control are associated with message queues.

As discussed previously, there are structure elements that may be undocumented for reasons of privacy. Any elements that are used internally will not be shown here. The operation of a `QUEUE` revolves around the `queue_t` structure, which contains the following fields, and possibly others that are undocumented:

```

/*
 * Fields marked with "RO" are read-only; STREAMS
 * modules and drivers are allowed to examine them, but not to change
 * them (except indirectly through the appropriate utility routines).
 */
struct queue {
    struct qinit *q_qinfo; /* RO: procedures and limits for queue */
    struct msgb *q_first; /* RO: head of msg queue for QUEUE */
    struct msgb *q_last; /* RO: tail of msg queue for QUEUE */
    struct queue *q_next; /* RO: next QUEUE in stream */
    caddr_t q_ptr; /* to private data structure */
    ushort q_count; /* RO: count of bytes on message queue */
    ushort q_flag; /* RO: QUEUE state */
    short q_minpsz; /* min packet size accepted by this QUEUE */
    short q_maxpsz; /* max packet size accepted by this QUEUE */
    ushort q_hiwat; /* msg queue high water mark, for flow control */
    ushort q_lowat; /* msg queue low water mark, for flow control */
};
typedef struct queue queue_t;

```

As described previously, two of these structures form a module. When a `queue_t` pair is allocated, their contents are zero unless specifically initialized. The following fields are initialized by STREAMS:

- `q_qinfo` - from `streamtab`
- `q_minpsz`, `q_maxpsz`, `q_hiwat`, `q_lowat` - from `module_info`
Copying values from `module_info` allows them to be changed in the `queue_t` without modifying the template (i.e., `streamtab` and `module_info`) values.

`q_count` is used in flow control calculations and is the aggregate number of bytes in the messages currently on the queue.

Service Procedures

Put procedures are generally required in pushable modules. Service procedures are optional. The general processing flow when both procedures are present is as follows. A message is received by the put procedure in a QUEUE, where some processing may be performed on the message. The put procedure transfers the message to the service procedure by use of the `putq()` utility. `putq()` places the message on the tail (see `q_last` in `queue_t`) of the message queue. Then, `putq()` will generally schedule the QUEUE for execution by the STREAMS scheduler following all other QUEUES currently scheduled. After some indeterminate delay (intended to be short), the scheduler calls the service procedure. The service procedure gets the first message (`q_first`) from the message queue with the `getq()` utility. The service procedure processes the message and passes it to the put procedure of the next QUEUE with `putnext()`. The service procedure gets the next message and processes it. This FIFO processing continues until the queue is empty or flow control blocks further processing. The service procedure returns to caller.

A service routine must never sleep and it has no user context. It must always return to its caller.

If no processing is required in the put procedure, the procedure does not have to be explicitly declared. Rather, `putq()` can be placed in the `qinit` structure declaration for the appropriate QUEUE side, to queue the message for the service procedure, e.g.:

```
static struct qinit winit = { putq, modwsrv, ... };
```

More typically, put procedures will, as a minimum, process priority messages (see below) to avoid queueing them.

The key attribute of a service procedure in the STREAMS architecture is delayed processing. When a service procedure is used in a module, the module developer is implying that there are other, more time-sensitive activities to be performed elsewhere in this stream, in other streams, or in the system in general. The presence of a service procedure is mandatory if the flow control mechanism is to be utilized by the QUEUE.

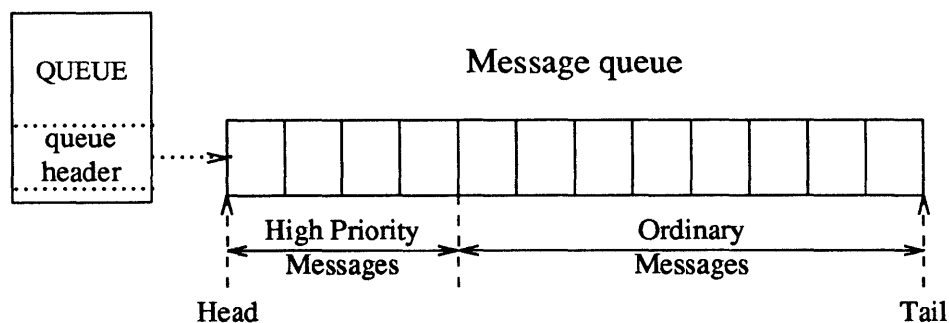
The delay for STREAMS to call a service procedure will vary with implementation and system activity. However, once the service procedure is scheduled, it is guaranteed to be called before user level activity is resumed.

Also see the *Put and Service Procedures* section of the *Introduction to STREAMS*.

Message Queues and Message Priority

Figure 3-3 depicts a message queue linked by `b_next` and `b_prev` pointers. As discussed in the *Introduction to STREAMS*, message queues grow when the STREAMS scheduler is delayed from calling a service procedure because of system activity, or when the procedure is blocked by flow control. When it is called by the scheduler, the service procedure processes enqueued messages in FIFO order. However, certain conditions require that the associated message (e.g., an `M_ERROR`) reach its stream destination as rapidly as possible. STREAMS does this by assigning all message types to one of the two levels of message queuing priority — high priority and ordinary. As shown in figure 3-4, when a message is queued, the `putq()` utility will place high priority messages at the head of the message queue, FIFO within their order of queuing.

Figure 3-4 Message Queue Priority



High priority messages are not subject to flow control. When they are queued by `putq()`, the associated QUEUE is always scheduled (in the same manner as any QUEUE; following all other QUEUES currently scheduled). When the service procedure is called by the scheduler, the procedure uses `getq()` to retrieve the first message on queue, which will be a high priority message, if present. Service procedures must be implemented to act on high priority messages immediately (see next section). The above mechanisms—priority message queuing, absence of flow control, and immediate processing by a procedure—result in rapid transport of priority messages between the originating and destination components in the stream.

The priority level for each message type is shown in the *Message Types* section of the *Supplementary STREAMS Material* chapter. Message queue management utilities are provided for use in service procedures (see the *Utilities* section of the *Supplementary STREAMS Material* chapter).

Flow Control

The elements of flow control are discussed in the *Other Facilities* section of the *Introduction to STREAMS*. Flow control is only used in a service procedure. Module and driver coding should observe the following guidelines for message priority. High priority messages, determined by the type of the first block in the message,

```
(bp->b_datap->db_type > QPCTL),
```

are not subject to flow control. They should be processed immediately and forwarded, as appropriate.

For ordinary messages, flow control must be tested before any processing is performed. The `canput()` utility determines if the forward path from the QUEUE is blocked by flow control. The manner in which STREAMS determines flow control status for modules and drivers is described under *Driver Flow Control* in the *Drivers* section.

This is the general processing for flow control. Retrieve the message at the head of the queue with `getq()`. Determine if the type is high priority and not to be processed here. If both are true, pass the message to the put procedure of the following QUEUE with `putnext()`. If the type is ordinary, use `canput()` to determine if messages can be sent onward. If `canput()` indicates messages should not be forwarded, put the message back on the queue with `putbq()` and return from the procedure. In all other cases, process the message.

The canonical representation of this processing within a service procedure is as follows:

```
while (getq != NULL)
    if (high priority message || canput)
        process message
        putnext
    else
        putbq
        return
```

NOTE *A service procedure must process all messages on its queue unless flow control prevents this.*

When an ordinary message is enqueued by `putq()`, `putq()` will cause the service procedure to be scheduled only if the queue was previously empty. If there are messages on the queue, `putq()` presumes the service procedure is blocked by flow control and the procedure will be automatically rescheduled by STREAMS when the block is removed. If the service procedure cannot complete processing as a result of conditions other than flow control (e.g., no buffers), it must assure it will return later (e.g., by use of `bufcall()`, see the *Advanced Topics* section) or it must discard all messages on the queue. If this is not done, STREAMS will never schedule the service procedure to be run unless the QUEUE's put procedure queues a priority message with `putq()`.

`putbq()` replaces messages at the beginning of the appropriate section of the message queue in accordance with their message type priority (see figure 3-4). This might not be the same position at which the message was retrieved by the preceding `getq()`. A subsequent `getq()` might return a different message.

Example

The filter module example of the *Messages* section is here modified to have a service procedure. The declarations from the example are unchanged except for the following lines (changes are shown in **bold**):

```

#include <sys/stropts.h>

static struct module_info minfo = {
    0, "ps_crmod", 0, INFPSZ, 512, 128
};
static int modopen(), modrput(), modwput();
static int modwsrv(), modclose();

static struct qinit winit = {
    modwput, modwsrv, NULL, NULL, NULL, &minfo, NULL
};

```

`stropts.h` is generally intended for user level. However, it includes definitions of flush message options common to user level, modules, and drivers. `module_info` now includes the flow control high- and low-water marks (512 and 128) for the write QUEUE (even though the same `module_info` is used on the read QUEUE side, the read side has no service procedure so flow control is not used). `qinit` now contains the service procedure pointer. `modopen()`, `modclose()`, and `modrput()` (read side put procedure) are unchanged from the *Modules* and *Messages* sections. The `bappend()` subroutine is also unchanged from the *Messages* section.

Procedures

The write side put procedures and the beginning of the service procedure are shown below:

```

static int modwput(q, mp)
queue_t *q;
register mblk_t *mp;
{
    if (mp->b_datap->db_type > QPCTL &&
        mp->b_datap->db_type != M_FLUSH)
        putnext(q, mp);
    else
        putq(q, mp);    /* Put it on the queue */
}

static int modwsrv(q)
queue_t *q;
{
    mblk_t *mp;

    while ((mp = getq(q)) != NULL) {
        switch (mp->b_datap->db_type) {

            default:
                /* always putnext high priority messages */
                if (mp->b_datap->db_type > QPCTL ||
                    canput(q->q_next)) {
                    putnext(q, mp);
                }
            }
        }
    }
}

```

```

        continue;
    }
    else {
        putbq(q, mp);
        return;
    }

case M_FLUSH:
    if (*mp->b_rptr & FLUSHW)
        flushq(q, FLUSHDATA);
    putnext(q, mp);
    continue;

```

`ps_crmod` performs a similar function to `crmod` of the previous section, but it uses a service routine.

`modwput()`, the write put procedure, switches on the message type. High priority messages that are not type `M_FLUSH` are `putnext()` to avoid scheduling. The others are queued for the service procedure. An `M_FLUSH` message is a request to remove all messages on one or both QUEUES. It can be processed in the put or service procedure.

`modwsrv()` is the write service procedure. It takes a single argument, a pointer to the write `queue_t`. `modwsrv()` processes only one high priority message, `M_FLUSH`. All other high priority messages are passed through. Actually, no other high priority messages should reach `modwsrv()`. The check is included to show the canonical form when high priority messages are queued by the put procedure.

For an `M_FLUSH` message, `modwsrv()` checks the first data byte. If `FLUSHW` (defined in `stropts.h`) is set in the byte, the write queue is flushed by use of `flushq()`. `flushq()` takes two arguments, the queue pointer and a flag. The flag indicates what should be flushed, data messages (`FLUSHDATA`) or everything (`FLUSHALL`). In this case, data includes `M_DATA`, `M_PROTO`, and `M_PCPROTO` messages. The choice of what types of messages to flush is module-specific. As a general rule, `FLUSHDATA` should be used.

Ordinary messages will be returned to the queue if

```
canput(q->q_next)
```

returns false, indicating the downstream path is blocked.

In the remaining part of `modwsrv()`, `M_DATA` messages are processed simi-

larly to the previous example:

```

case M_DATA: {
    mblk_t *nbp = NULL;
    mblk_t *next;

    if (!canput(q->q_next)) {
        putbq(q, mp);
        return;
    }
    /* Filter data, appending to queue */
    for (; mp != NULL; mp = next) {
        while (mp->b_rptr < mp->b_wptr) {

            if (*mp->b_rptr == '\n')
                if (!bappend(&nbp, '\r'))
                    goto push;
            if (!bappend(&nbp, *mp->b_rptr))
                goto push;
            mp->b_rptr++;
            continue;

        push:
            putnext(q, nbp);
            nbp = NULL;
            if (!canput(q->q_next)) {
                if (mp->b_rptr >= mp->b_wptr) {
                    next = mp->b_cont;
                    freeb(mp);
                    mp=next;
                }
                if (mp)
                    putbq(q, mp);
                return;
            }
        }
        next = mp->b_cont;
        freeb(mp);
    }
    if (nbp)
        putnext(q, nbp);
}
}
}

```

The differences in M_DATA processing between this and the previous example relate to the manner in which the new messages are forwarded and to flow control. For the purpose of demonstrating alternative means of processing messages, this version creates individual new messages rather than a single message containing multiple message blocks. When a new message block is full, it is immediately forwarded with `putnext()` rather than being linked into a single,

large message (as was done in the previous example). This alternative may not be desirable because message boundaries will be altered and because of the additional overhead of handling and scheduling multiple messages.

When the filter processing is performed (following `push()`), flow control is checked (`canput()`) after, rather than before, each new message is forwarded. This is done because there is no provision to hold the new message until the QUEUE becomes unblocked. If the downstream path is blocked, the remaining part of the original message is returned to the queue. Otherwise, processing continues.

Another difference between the two examples is that each message block of the original message is returned to the pool with `freeb()` when its processing is completed.

3.6. Drivers

Overview of Drivers

This section describes the organization of a STREAMS driver, and discusses some of the processing typically required in drivers. Certain elements of driver flow control are discussed. Procedures for handling user ioctls, common to modules and drivers, are described.

As discussed under *Stream Construction* in the *STREAMS Mechanism* section, driver and module organization are very similar. The interfaces to all the driver procedures are identical to module interfaces, and driver procedures must be reentrant. As described under *Environment* in the *Modules* section, the driver put and service procedures have no user environment and cannot sleep. Other than with `open()` and `close()`, a driver interfaces with a user process by messages, and indirectly, through flow control.

There are two significant differences between modules and drivers. First, a device driver must also be accessible from an interrupt as well as from the stream, and second, a driver can have multiple streams connected to it. Multiple connections occur when more than one minor device uses the same driver and in the case of multiplexors (see the *Multiplexing* section). However, these particular differences are not recognized by the STREAMS mechanism; they are handled by developer-provided code included in the driver procedures.

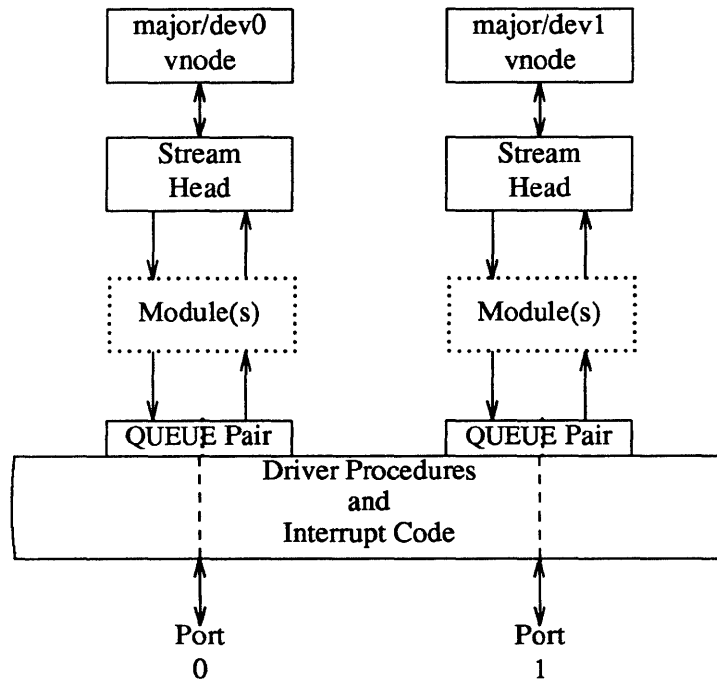
Figure 3-5 shows multiple streams (corresponding to minor devices), to a common driver. This depiction of two streams connected to a single driver is somewhat misleading. These are really two distinct streams opened from the same `cdevsw` (i.e., same major device). Consequently, they have the same `streamtab` and the same driver procedures. Modules opened from the same `fmodsw` might be depicted similarly if they had any reason to be cognizant, as do drivers, of common resources or alternate instantiations.

Multiple instantiations (minor devices) of the same driver are handled during the initial `open` for each device. Typically, the `queue_t` address is stored in a driver-private structure indexed by the minor device number. The structure is typically pointed at by `q_ptr` (see the *Message Queues and Service Procedures* section). When the messages are received by the QUEUE, the calls to the driver

put and service procedures pass the address of the `queue_t`, allowing the procedures to determine the associated device.

In addition to these differences, a driver is always at the end of a stream. As a result, drivers must include standard processing for certain message types that a module might simply be able to pass to the next component.

Figure 3-5 *Device Driver Streams*



Driver Flow Control

The same utilities (described in the *Message Queues and Service Procedures* section) and mechanisms used for module flow control are used by drivers. However, they are typically used in a different manner in drivers, because a driver generally does not have a service procedure. The developer sets flow control values (`mi_hiwat` and `mi_lowat`) in the write side `module_info` structure, which STREAMS will copy into `q_hiwat` and `q_lowat` in the `queue_t` structure of the QUEUE. A device driver typically has no write service procedure, but does maintain a write message queue. When a message is passed to the driver write side put procedure, the procedure will determine if device output is in progress. In the event output is busy, the put procedure cannot immediately send the message and calls the `putq()` utility (see the *Utilities* section of the *Supplementary STREAMS Material* chapter) to queue the message. (Note that the driver might have elected to queue the message in all cases.) `putq()` recognizes the absence of a service procedure and does not schedule the QUEUE.

When the message is queued, `putq()` increments the value of `q_count` (the enqueued byte count, see the beginning of the *Message Queues and Service Procedures* section) by the size of the message and compares the result against the driver's write high water limit (`q_hiwat`) value. If the count exceeds `q_hiwat`, `putq()` will set the internal FULL (see *Flow Control* in the

Introduction to STREAMS) indicator for the driver write QUEUE. This will cause messages from upstream to be halted (`canput()` returns FALSE) until the write queue count reaches `q_lowat`. The driver messages waiting to be output are dequeued by the driver output interrupt routine with `getq()`, which decrements the count. If the resulting count is below `q_lowat`, `getq()` will back-enable any upstream QUEUE that had been blocked. The above STREAMS processing also applies to modules on both write and read sides of the stream.

Device drivers typically discard input when unable to send it to a user process. However, STREAMS allows flow control to be used on the driver read side, possibly to handle temporary upstream blocks. This is described in the *Advanced Topics* section in the *Advanced Flow Control* section.

To some extent, a driver or module can control when its upstream transmission will become blocked. Control is available through the `M_SETOPTS` message (see the *Advanced Topics* section, here, and the *Message Types* section of *Supplementary STREAMS Material*) to modify the Stream Head read side flow control limits.

Driver Programming

The example below shows how a simple interrupt-per-character line printer driver could be written. The driver is unidirectional and has no read side processing. It demonstrates some differences between module and driver programming, including the following:

Open handling

A driver is passed a minor device number or is asked to select one (see next section).

Flush handling

A driver must loop `M_FLUSH` messages back upstream.

Ioctl handling

A driver must acknowledge, in the form of a nak, `ioctl` messages it does not understand. This is discussed under *Driver and Module Ioctls*, below. Write side flow control is also illustrated as described above.

Driver Declarations

The driver declarations are as follows:

```

/* Simple line printer driver. */

#include "lp.h"
#include <sys/types.h>
#include <sys/param.h>
#include <sys/sysmacros.h>
#include <sys/stream.h>
#include <sys/stropts.h>
#include <sys/dir.h>          /* required for user.h */
#include <sys/signal.h>      /* required for user.h */
#include <sys/user.h>
#include <sys/errno.h>

static struct module_info minfo = {

```

```

    0, "lp", 0, INFPSZ, 150, 50
};

static int lpopen(), lpclose(), lpwput();

static struct qinit rinit = {
    NULL, NULL, lpopen, lpclose, NULL, &minfo, NULL
};
static struct qinit winit = {
    lpwput, NULL, NULL, NULL, NULL, &minfo, NULL
};
struct streamtab lpinfo = { &rinit, &winit, NULL, NULL };

#define SET_OPTIONS (('l' << 8) | 1) /* really must be in a .h file,
                                     see _IO(x,y) in sysiocom.h */
/*
 * This is a private data structure, one per minor device number.
 */
struct lp {
    short flags; /* flags -- see below */
    mblk_t *msg; /* current message being output */
    queue_t *qptr; /* back pointer to write queue */
};
/* Flags bits */
#define BUSY 1 /* device is running and interrupt is pending */

struct lp lp_lp[NLP]; /* per device lp structure array */
int lp_cnt = NLP; /* number of valid minor devices */

```

As noted for modules in the *Modules* section, configuring a STREAMS driver does not require the driver procedures to be externally accessible; only `streamtab` must be. All STREAMS driver procedures would typically be declared `static`.

There is no read side put or service procedure. The flow control limits for use on the write side are 50 and 150 characters. The private `lp` structure is indexed by the minor device number and contains these elements:

flags

A set of flags. Only one bit is used: `BUSY` indicates that output is active and a device interrupt is pending.

msg

A pointer to the current message being output.

qptr

A back pointer to the write queue. This is needed to find the write queue

during interrupt processing.

Driver Open

The driver open, `lpopen()`, has the same interface as the module open:

```
static int lpopen(q, dev, flag, sflag)
queue_t *q;    /* read queue */
dev_t dev;
int flag;
int sflag;
{
    struct lp *lp;

    /* Check if non-driver open */
    if (sflag)
        return (OPENFAIL);

    /* Dev is major/minor */
    dev = minor(dev);
    if (dev >= lp_cnt)
        return (OPENFAIL);

    /* Check if open already. q_ptr is assigned below */
    if (q->q_ptr) {
        u.u_error = EBUSY; /* only 1 user of the printer at a time */
        return (OPENFAIL);
    }

    lp = &lp_lp[dev];
    lp->q_ptr = WR(q);
    q->q_ptr = (char *) lp;
    WR(q)->q_ptr = (char *) lp;
    return (dev);
}
```

The stream flag, `sflag`, must have the value 0, indicating a normal driver open. `dev` holds both the major and minor device numbers for this port. After checking `sflag`, the open flag, `lpopen()` extracts the minor device from `dev`, using the `minor()` macro defined in `sysmacros.h`.

Remember that the use of major devices, minor devices, and the `minor()` macro may be machine dependent.

The minor device number selects a printer and must be less than `lp_cnt`.

The next check, `if (q->q_ptr)...`, determines if this printer is already open. In this case, `EBUSY` is returned to avoid merging printouts from multiple users. `q_ptr` is a driver/module private data pointer. It can be used by the driver for any purpose and is initialized to zero by STREAMS. In this example, the driver sets the value of `q_ptr`, in both the read and write `queue_t` structures, to point to a private data structure for the minor device, `lp_lp[dev]`.

WR is one of three QUEUE pointer macros. As discussed in the *Stream Construction* section, there are no physical pointers between QUEUES, and these macros (see *Utilities* in the *Supplementary STREAMS Material* section) generate the pointer. WR(q) generates the write pointer from the read pointer, RD(q) generates the read pointer from the write pointer and OTHER(q) generates the mate pointer from either.

Driver Processing Procedures

This example only has a write put procedure:

```
static int lpwput(q, mp)
queue_t *q;      /* write queue */
register mblk_t *mp; /* message pointer */
{
    register struct lp *lp;
    int s;

    lp = (struct lp *)q->q_ptr;

    switch (mp->b_datap->db_type) {
    default:
        freemsg(mp);
        break;
    case M_FLUSH:
        /* Canonical flush handling */
        if (*mp->b_rptr & FLUSHW) {
            flushq(q, FLUSHDATA);
            s = spl5();
            /* also flush lp->msg since it is logically
             * at the head of the write queue */
            if (lp->msg) {
                freemsg(lp->msg);
                lp->msg = NULL;
            }
            (void) splx(s);
        }

        if (*mp->b_rptr & FLUSHR) {
            flushq(RD(q), FLUSHDATA);
            *mp->b_rptr &= ~FLUSHW;
            qreply(q, mp);
        } else
            freemsg(mp);
        break;
    case M_IOCTL:
    case M_DATA:
        putq(q, mp);
        s = spl5();
        if (!(lp->flags & BUSY))
            lpout(lp);
        (void) splx(s);
    }
}
```

}

Driver Flush Handling

The write put procedure, `lpwput()`, illustrates driver `M_FLUSH` handling: Note that all drivers are expected to incorporate this flush handling. If `FLUSHW` is set, the write message queue is flushed, and also (for this example) the leading message (`lp->msg`). `sp15` is used to protect the critical code, assuming the device interrupts at level 5. If `FLUSHR` is set, the read queue is flushed, the `FLUSHW` bit is cleared, and the message is sent upstream using `qreply()`. If `FLUSHR` is not set, the message is discarded.

The Stream Head always performs the following actions on flush requests received on the read side from downstream. If `FLUSHR` is set, messages waiting to be sent to user space are flushed. If `FLUSHW` is set, the Stream Head clears the `FLUSHR` bit and sends the `M_FLUSH` message downstream. In this manner, a single `M_FLUSH` message sent from the driver can reach all QUEUES in a stream. A module must send two `M_FLUSH` messages to have the same effect.

`lpwput()` enqueues `M_DATA` and `M_IOCTL` (see the *Driver and Module Ioctls* section, below) messages and, if the device is not busy, starts output by calling `lpout()`. Message types that are not recognized are discarded.

Driver Interrupt

`lpintr()` is the driver interrupt routine:

```

/* Device interrupt routine. */

lpintr(dev)
int dev;    /* minor device number of lp */
{
    register struct lp *lp;

    lp = &lp_lp[dev];
    if (!(lp->flags & BUSY)) {
        printf("lp: unexpected interrupt\n");
        return;
    }
    lp->flags &= ~BUSY;
    lpout(lp);
}

/* Start output to device -
   called at interrupt level, used by put procedure and driver */

lpout(lp)
register struct lp *lp;
{
    register mblk_t *bp;
    queue_t *q;

    q = lp->qptr;
loop:
    if ((bp = lp->msg) == NULL) {

```



```

    if ((bp = getq(q)) == NULL)
        return;
    if (bp->b_datap->db_type == M_IOCTL) {
        lpdoioctl(lp, bp);
        goto loop;
    }
    lp->msg = bp;
}
if (bp->b_rptr >= bp->b_wptr) {
    bp = lp->msg->b_cont;
    lp->msg->b_cont = NULL;
    freeb(lp->msg);
    lp->msg = bp;
    goto loop;
}

lpoutchar(lp, *bp->b_rptr++);
lp->flags |= BUSY;
}

```

lpout () simply takes a character from the queue and sends it to the printer. The processing is logically similar to the service procedure in the *Message Queues and Service Procedures* section. For convenience, the message currently being output is stored in lp->msg.

Two mythical routines need to be supplied:

lpoutchar

send a character to the printer and interrupt when complete

lpsetopt

set the printer interface options

Driver and Module Ioctls

Drivers and modules interface with ioctl(2) system calls through messages. Almost all STREAMS generic ioctl()s (see the streamio(4) man page) go no further than the Stream Head. The capability to send an ioctl() downstream, is provided by the I_STR ioctl. The Stream Head processes an I_STR by constructing an M_IOCTL message (see *Message Types* in the *Supplementary STREAMS Material* chapter) from data provided in the call, and sends that message downstream.

NOTE *In addition, since ioctl() codes in SunOS include the size of the parameter used for the ioctl() as well as an indication of whether this parameter is to be copied to or from the user process, the I_STR ioctl need not be used if the parameter contains 255 or fewer bytes and is of a fixed size.*

The user process that issued the ioctl() is blocked until a module or driver responds with either an M_IOCACK (ack) or M_IOCNAK (nak) message, or until the request “times out” after a user-specified interval. The STREAMS module or driver that generates an ack can also return information to the process. If the Stream Head does not receive one of these messages in the specified time, the ioctl() call fails.

A module that receives an unrecognized `M_IOCTL` message should pass it on unchanged. A driver that receives an unrecognized `M_IOCTL` should nak it.

`lpout()` traps `M_IOCTL` messages and calls `lpdoioctl()` to process them:

```

lpdoioctl(lp, mp)
struct lp *lp;
mblk_t *mp;
{
    struct iocblk *iocp;
    queue_t *q;

    q = lp->qptr;

    /* 1st block contains iocblk structure */
    iocp = (struct iocblk *)mp->b_rptr;

    switch (iocp->ioc_cmd) {
    case SET_OPTIONS:
        /* Count should be exactly one short's worth */
        if (iocp->ioc_count != sizeof(short))
            goto iocnak;
        /* Actual data is in 2nd message block */
        lpsetopt(lp, *(short *)mp->b_cont->b_rptr);

        /* ACK the ioctl */
        mp->b_datap->db_type = M_IOCACK;
        iocp->ioc_count = 0;
        qreply(q, mp);
        break;
    default:
    iocnak:
        /* NAK the ioctl */
        mp->b_datap->db_type = M_IOCNAK;
        qreply(q, mp);
    }
}

```

`lpdoioctl()` illustrates `M_IOCTL` processing: The first part also applies to modules. An `M_IOCTL` message contains a `struct iocblk` in its first block. The first block is followed by zero or more `M_DATA` blocks. The optional `M_DATA` blocks typically contain any user-supplied data.

The form of an `iocblk` is as follows:

```

struct iocblk {
    int      ioc_cmd;      /* ioctl command type */
    ushort   ioc_uid;     /* effective uid of user */
    ushort   ioc_gid;     /* effective gid of user */
    uint     ioc_id;      /* ioctl id */
    uint     ioc_count;   /* count of bytes in data field */
    int      ioc_error;   /* error code */
    int      ioc_rval;    /* return value */
};

```

`ioc_cmd` contains the command supplied by the user. In this example, only one command is recognized, `SET_OPTIONS`. `ioc_count` contains the number of user supplied data bytes. For this example, it must equal the size of a short (2 bytes). The user data is sent directly to the printer interface using `lpsetopt()`. Next, the `M_IOCTL` message is changed to type `M_IOCACK` and the `ioc_count` field is set to zero to indicate that no data is to be returned to the user. Finally, the message is sent upstream using `qreply()`. If `ioc_count` was left non-zero, the Stream Head would copy that many bytes from the 2nd - Nth message blocks into the user buffer.

If the `M_IOCTL` message is not understood or in error for any reason, the driver must set the type to `M_IOCNAK` and send the message upstream. No data can be sent to a user in this case. The Stream Head will cause the `ioctl()` call to fail with the error number `EINVAL`. The driver has the option of setting `ioc_error` to an alternate error number if desired.

NOTE `ioc_error` can be set to a non-zero value by both `M_IOCACK` and `M_IOCNAK`. This will cause that value to be returned as an error number to the process that sent the `ioctl()`.

Driver Close

The driver close clears any message being output. Any messages left on the message queue will be automatically removed by STREAMS.

```

static int lpclose(q)
queue_t *q; /* read queue */
{
    struct lp *lp;
    int s;

    lp = (struct lp *) q->q_ptr;
    /* Free message, queue is automatically flushed by STREAMS */
    s = spl5();
    if (lp->msg) {
        freemsg(lp->msg);
        lp->msg = NULL;
    }
    (void) splx(s);
}

```

3.7. Example Driver

Cloning

The clone mechanism has been developed as a convenience. It allows a user to open a driver without specifying the minor device. When a stream is opened, a flag indicating a clone open is tested by the driver open routine. If the flag is set, the driver returns an unused minor device number. The clone driver (see the `clone(4)` man page) is a system dependent pseudo driver.

Loop-Around Driver

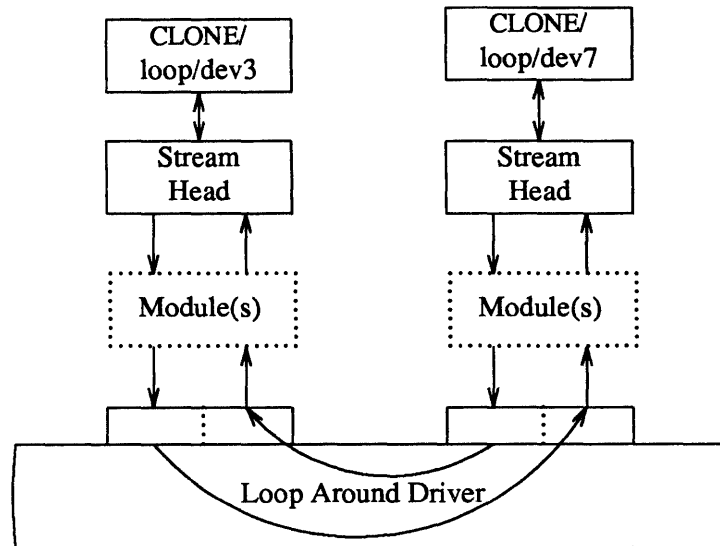
The loop-around driver is not part of the SunOS 4.1 release. It is documented here simply as another example of a STREAMS driver.

The loop-around driver is a pseudo-driver that loops data from one open stream to another open stream. The user processes see the associated files as a full duplex pipe. The streams are not physically linked. The driver is a simple upper multiplexor (see the next section), which the driver simply passes messages from one stream's write QUEUE to the other stream's read QUEUE.

To create a pipe, a process opens two streams, obtains the minor device number associated with one of the returned file descriptors, and sends the device number in a `LOOP_SET ioctl(2)` to the other stream. For each `open()`, the driver open places the passed `queue_t` pointer in a driver interconnection table, indexed by the device number. When the driver later receives the `LOOP_SET` request as an `M_IOCTL` message, it uses the device number to locate the other stream's interconnection table entry, and stores the appropriate `queue_t` pointers in both of the streams' interconnection table entries.

Subsequently, when messages other than `M_IOCTL` or `M_FLUSH` are received by the driver on either stream's write side, the messages are switched to the read QUEUE following the driver on the other stream's read side. The resultant logical connection is shown in figure 3-6. Flow control between the two streams must be handled by special code since STREAMS will not automatically propagate flow control information between two streams that are not physically interconnected.

Figure 3-6 Loop Around Streams



The declarations for the driver are:

```

/*
 * Loop around driver
 */

#include "loop.h"
#include <sys/types.h>
#include <sys/param.h>
#include <sys/sysmacros.h>
#include <sys/stream.h>
#include <sys/stropts.h>
#include <sys/user.h>
#include <sys/errno.h>

static struct module_info minfo = {
    0, "loop", 0, INFPSZ, 512, 128
};

static int loopopen(), loopclose(), loopwput();
static int loopwsrv(), looprsrv();

static struct qinit rinit = {
    NULL, looprsrv, loopopen, loopclose, NULL, &minfo, NULL
};

static struct qinit winit = {
    loopwput, loopwsrv, NULL, NULL, NULL, &minfo, NULL
};

struct streamtab loopinfo = {&rinit, &winit, NULL, NULL};

```

```

struct loop {
    queue_t *qp_ptr; /* back pointer to write queue */
    queue_t *oqp_ptr; /* pointer to connected read queue */
};

#define LOOP_SET    _IOW(1, 1, int) /* should be in a .h file,
                                     see _IOW(x,y,t) in sysiocom.h */

struct loop loop_loop[NLOOP];
int loop_cnt = NLOOP;

```

The `loop` structure contains the interconnection information for a pair of streams. `loop_loop` is indexed by the minor device number. When a stream is opened to the driver, the address of the corresponding `loop_loop` element is placed in `q_ptr` (private data structure pointer) of the read and write side `queue_t`s. Since STREAMS clears `q_ptr` when the `queue_t` is allocated, a `NULL` value of `q_ptr` indicates an initial `open()`. `loop_loop` is used to verify that this stream is connected to another open stream.

The `open` procedure includes canonical clone processing which enables a single file system node to yield a new minor device/vnode each time the driver is opened:

```

static int loopopen(q, dev, flag, sflag)
queue_t *q;
{
    struct loop *loop;

    /*
     * If CLONEOPEN, pick a minor device number to use.
     * Otherwise, check the minor device range.
     */
    if (sflag == CLONEOPEN) {
        for (dev = 0; dev < loop_cnt; dev++) {
            if (loop_loop[dev].qp_ptr == NULL)
                break;
        }
    }
    else
        dev = minor(dev); /* device not a clone, but specific */

    if (dev >= loop_cnt)
        return (OPENFAIL); /* default = ENXIO */

    /* Set up data structures */
    if (q->q_ptr) /* already open */
        return (dev);

    loop = &loop_loop[dev];
    WR(q)->q_ptr = (char *) loop;
    q->q_ptr = (char *) loop;
    loop->oqp_ptr = WR(q);
}

```

```

/*
 * The return value is the minor device.
 * For CLONEOPEN case, this will be used for
 * newly allocated vnode
 */
return (dev);
}

```

In `loopopen()`, `sflag` can be `CLONEOPEN`, indicating that the driver should pick a minor device (i.e., the user does not care which minor device is used). In this case, the driver scans its private `loop_loop` data structure to find an unused minor device number. If `sflag` has not been set to `CLONEOPEN`, the passed-in minor device is used.

The return value is the minor device number. In the `CLONEOPEN` case, this value will be used by the STREAMS driver open code for the newly allocated vnode and will then be passed to the user.

Write Put Procedure

Since the messages are switched to the read QUEUE following the other stream's read side, the driver needs a put procedure only on its write side:

```

static int loopwput(q, mp)
queue_t *q;
mblock_t *mp;
{
    register struct loop *loop;

    loop = (struct loop *)q->q_ptr;

    switch (mp->b_datap->db_type) {
    case M_IOCTL: {
        struct iocblk *iocp;
        int error;

        iocp = (struct iocblk *)mp->b_rptr;
        switch (iocp->ioc_cmd) {
        case LOOP_SET: {
            int to; /* other minor device */
            /*
             * Sanity check. ioc_count contains the amount of
             * user supplied data, which must equal the size of an int.
             */

            if (iocp->ioc_count != sizeof (int)) {
                error = EINVAL;
                goto iocnak;
            }

            /* fetch other dev from 2nd message block */

            to = *(int *)mp->b_cont->b_rptr;

```

```

/*
 * More sanity checks. The minor must be in range, open already.
 * Also, this device and the other one must be disconnected.
 */

if (to >= loop_cnt || to < 0 ||
    !loop_loop[to].qptr) {
    error = ENXIO;
    goto iocnak;
}

if (loop->oqptr || loop_loop[to].oqptr) {
    error = EBUSY;
    goto iocnak;
}

/*
 * Cross connect streams via the loop structures
 */

loop->oqptr = RD(loop_loop[to].qptr);
loop_loop[to].oqptr = RD(q);

/*
 * Return successful ioctl. Set ioc_count
 * to zero, since there is no data to return.
 */

mp->b_datap->db_type = M_IOCACK;
iocp->ioc_count = 0;
qreply(q, mp);
break;
}

default:
    error = EINVAL;
iocnak:
    /*
     * Bad ioctl. Setting ioc_error causes the
     * ioctl call to return that particular errno.
     * By default, ioctl will return EINVAL on failure
     */
    mp->b_datap->db_type = M_IOCNAK;
    iocp->ioc_error = error; /* set returned errno */
    qreply(q, mp);

}
break;
}

```


`loopwput()` shows another use of an `ioctl()` call (see *Driver and Module Ioctls* in the *Drivers* section, above). The driver supports a `LOOP_SET` value of `ioc_cmd` in the `iocblk` of the `M_IOCTL` message. `LOOP_SET` instructs the driver to connect the current open stream to the stream indicated in the message. The second block of the `M_IOCTL` message holds an integer that specifies the minor device number of the stream to connect to.

The driver performs several sanity checks: Does the second block have the proper amount of data? Is the "to" device in range? Is the "to" device open? Is the current stream disconnected? Is the "to" stream disconnected?

If everything checks out, the read `queue_t` pointers for the two streams are stored in the respective `oqptr` fields. This cross-connects the two streams indirectly, via `loop_loop`.

Canonical flush handling is incorporated in the put procedure:

```

case M_FLUSH:
    if (*mp->b_rptr & FLUSHW)
        flushq(q, 0);
    if (*mp->b_rptr & FLUSHR) {
        flushq(RD(q), 0);
        *mp->b_rptr &= ~FLUSHW;
        qreply(q, mp);
    } else
        freemsg(mp);
    break;
default:
    /*
     * If this stream isn't connected, send an M_ERROR upstream.
     */
    if (loop->oqptr == NULL) {
        (void) putctl1(RD(q)->q_next, M_ERROR, ENXIO);
        freemsg(mp);
        break;
    }
    putq(q, mp);
}
}

```

Finally, `loopwput()` enqueues all other messages (e.g., `M_DATA` or `M_PROTO`) for processing by its service procedure. A check is made to see if the stream is connected. If not, an `M_ERROR` is sent upstream to the Stream Head (see below).

`putctl1()` and `putctl()` (see below) are utilities that allocate a non-data (i.e., not `M_DATA`, `M_PROTO`, or `M_PCPROTO`) type message, place one byte in the message (for `putctl1()`) and call the put procedure of the specified QUEUE (see *Utilities* in the *Supplementary STREAMS Material* chapter).

Service Procedures

Service procedures are required on both the write and read sides for purposes of flow control:

```

static int loopwsrv(q)
register queue_t *q;
{
    mblk_t *mp;
    register struct loop *loop;

    loop = (struct loop *)q->q_ptr;

    while ((mp = getq(q)) != NULL) {

        /*
         * Check if we can put the message up the other stream read queue
         */

        if (mp->b_datap->db_type <= QPCTL &&
            !canput(loop->oqptr->q_next)) {
            putbq(q, mp); /* read side is blocked */
            break;
        }

        /* send message */
        /* to queue following other stream read queue */

        putnext(loop->oqptr, mp);
    }
}

static int looprsrv(q)
queue_t *q;
{

    /* Entered only when "back enabled" by flow control */

    struct loop *loop;

    loop = (struct loop *)q->q_ptr; /* guards against race condition */
    if (loop->oqptr == NULL)
        return;

    /* manually enable write service procedure */

    qenable(WR(loop->oqptr));
}

```

The write service procedure, `loopwsrv()`, takes on the canonical form (see the *Message Queues and Service Procedures* section) with a difference. The QUEUE being written to is not downstream, but upstream (found via `oqptr`) on the other

stream.

In this case, there is no read side put procedure so the read service procedure, `looprsrv()`, is not scheduled by an associated put procedure, as has been done previously. `looprsrv()` is scheduled only by being back-enabled when its upstream becomes unstuck from flow control blockage. The purpose of the procedure is to re-enable the writer (`loopwsrv()`) by using `oqptr` to find the related `queue_t`. `loopwsrv()` cannot be directly back-enabled by STREAMS because there is no direct `queue_t` linkage between the two streams. Note that no message ever gets queued to the read service procedure. Messages are kept on the write side so that flow control can propagate up to the Stream Head. There is a defensive check to see if the cross-connect has broken. `qenable()` schedules the write side of the other stream.

Close

`loopclose()` breaks the connection between the streams.

```
static int loopclose(q)
queue_t *q;
{
    register struct loop *loop;

    loop = (struct loop *)q->q_ptr;
    loop->oqptr = NULL;

    /*
     * If we are connected to another stream, break the
     * linkage, and send a hangup message.
     * The hangup message causes the stream head to reject writes,
     * allow the queued data to be read completely, and then
     * return EOF on subsequent reads.
     */

    if (loop->oqptr) {
        ((struct loop *)loop->oqptr->q_ptr)->qptr = NULL;
        ((struct loop *)loop->oqptr->q_ptr)->oqptr = NULL;
        putctl(loop->oqptr->q_next, M_HANGUP);
        loop->oqptr = NULL;
    }
}
```

`loopclose()` sends an `M_HANGUP` message (see above) up the connected stream to the Stream Head.

NOTE *This driver can be implemented much more cleanly by actually linking the `q_next` pointers of the `queue_t` pairs of the two streams.*

3.8. Multiplexing

Multiplexing Configurations

This section describes how STREAMS multiplexing configurations are created and discusses multiplexing drivers. A STREAMS multiplexor is a pseudo-driver with multiple streams connected to it. The primary function of the driver is to switch messages among the connected streams. Multiplexor configurations are created from user level by system calls. The *Other Facilities* section of the *Introduction to STREAMS* contains the required introduction to STREAMS multiplexing.

STREAMS-related system calls are used to set up the “plumbing”, or stream interconnections, for multiplexing pseudo-drivers. The subset of these calls that allows a user to connect (and disconnect) streams below a pseudo-driver is referred to as the multiplexing facility. This type of connection will be referred to as a 1-to-M, or lower, multiplexor configuration. This configuration must always contain a multiplexing pseudo-driver, which is recognized by STREAMS as having special characteristics.

Multiple streams can be connected above a driver by use of `open(2)` calls. This was done for the loop-around driver of the previous section and for the driver handling multiple minor devices in the *Drivers* section. There is no difference between the connections to these drivers, only the functions performed by the driver are different. In the multiplexing case, the driver routes data between multiple streams. In the device driver case, the driver routes data between user processes and associated physical ports. Multiplexing with streams connected above will be referred to as an N-to-1, or upper, multiplexor. STREAMS does not provide any facilities beyond `open()` and `close(2)` to connect or disconnect upper streams for multiplexing purposes.

From the driver’s perspective, upper and lower configurations differ only in the way they are initially connected to the driver. The implementation requirements are the same: route the data and handle flow control. All multiplexor drivers require special developer-provided software to perform the multiplexing data routing and to handle flow control. STREAMS does not directly support flow control among multiple streams.

M-to-N multiplexing configurations are implemented by using both of the above mechanisms in a driver. Complex multiplexing trees can be created by cascading multiplexing streams below one another.

As discussed in the *Drivers* section, the multiple streams that represent minor devices are actually distinct streams in which the driver keeps track of each stream attached to it. The streams are not really connected to their common driver. The same is true for STREAMS multiplexors of any configuration. The multiplexed streams are distinct and the driver must be implemented to do most of the work. As stated above, the only difference between configurations is the manner of connecting and disconnecting. Only lower connections have use of the multiplexing facility.

Connecting Lower Streams

A lower multiplexor is connected as follows: The initial `open()` to a multiplexing driver creates a stream, as in any other driver. As usual, `open()` uses the first two `streamtab` structure entries (see *Opening a Stream* in the *Streams Mechanism* section) to create the driver QUEUES. At this point, the only distinguishing characteristic of this stream are non-NULL entries in the `streamtab` `st_mux[rw]init` (mux) fields:

```

struct streamtab {
    struct qinit    *st_rdinit;    /* defines read QUEUE */
    struct qinit    *st_wrinit;    /* defines write QUEUE */
    struct qinit    *st_muxrinit; /* for multiplexing drivers only */
    struct qinit    *st_muxwinit; /* for multiplexing drivers only */
    char            **st_modlist; /* list of modules to be pushed */
};

```

These fields are ignored by the `open()` (see the rightmost stream in figure 3-7). Any other stream subsequently opened to this driver will have the same `streamtab` and thereby the same mux fields.

Next, another file is opened to create a (soon to be) lower stream. The driver for the lower stream is typically a device driver (see the leftmost stream in figure 3-7). This stream has no distinguishing characteristics. It can include any driver compatible with the multiplexor. Any modules required on the lower stream must be pushed onto it now.

Next, this lower stream is connected below the multiplexing driver with an `I_LINK` or `I_PLINK` `ioctl()` call (see the `streamio(4)` man page). These `ioctls` differ in that the multiplexing configuration created by `I_PLINK` persists after the upper stream is closed whereas that created by `I_LINK` does not. As shown in figure 3-1, all stream components are constructed in a similar manner. The Stream Head points to the stream-head-routines as its procedures (known via its `queue_t`). An `I_LINK` to the upper stream, referencing the lower stream, causes STREAMS to modify the contents of the Stream Head in the lower stream. The pointers to the stream-head-routines, and other values, in the Stream Head are replaced with those contained in the mux fields of the multiplexing driver's `streamtab`. Changing the stream-head-routines on the lower stream means that all subsequent messages sent upstream by the lower stream's driver will, ultimately, be passed to the `put` procedure designated in `st_muxrinit`, the multiplexing driver. The `I_LINK` also establishes this upper stream as the control stream for this lower stream. STREAMS remembers the relationship between these two streams until the upper stream is closed, or the lower stream is unlinked.

Finally, the Stream Head sends to the multiplexing driver an `M_IOCTL` message with `ioc_cmd` set to `I_LINK` or `I_PLINK` (see discussions of the `iocblk` structure in the *Drivers* section, above, and in the *Kernel Structures* section of *Supplementary STREAMS Material* chapter). The `M_DATA` part of the `M_IOCTL` contains a `linkblk` structure:

```

struct linkblk {
    queue_t *l_qtop;      /* lowest level write queue of upper stream */
    queue_t *l_qbot;     /* highest level write queue of lower stream */
    int l_index;         /* system-unique index for lower stream. */
};

```

The multiplexing driver stores information from the `linkblk` in private storage and returns an `M_IOCACK` message (ack). `l_index` is returned to the process requesting the `I_LINK` or `I_PLINK`. This value can be used later by the process to disconnect this stream, as described below. `linkblk` contents are further discussed below.

An `I_LINK` or `I_PLINK` is required for each lower stream connected to the driver. Additional upper streams can be connected to the multiplexing driver by `open()` calls. Any message type can be sent from a lower stream to user process(es) along any of the upper streams. The upper stream(s) provides the only interface between the user process(es) and the multiplexor.

Note that no direct data structure linkage is established for the linked streams. The `q_next` pointers of the lower stream still appear to connect with a Stream Head. Messages flowing upstream from a lower driver (a device driver or another multiplexor) will enter the multiplexing driver (i.e., Stream Head) put procedure with `l_qbot` as the `queue_t` value. The multiplexing driver has to route the messages to the appropriate upper (or lower) stream. Similarly, a message coming downstream from user space on the control, or any other, upper stream has to be processed and routed, if required, by the driver.

Also note that the lower stream (see the headers and file descriptors in figure 3-8) is no longer accessible from user space. This causes all system calls to the lower stream to return `EINVAL`, with the exception of `close()`. This is why all modules have to be in place before the lower stream is linked to the multiplexing driver. As a general rule, the lower stream file should be closed after it is linked (see following section). This does not disturb the multiplexing configuration.

Note that the absence of direct linkage between the upper and lower streams means that STREAMS flow control has to be handled by special code in the multiplexing driver. The flow control mechanism cannot see across the driver.

In general, multiplexing drivers should be implemented so that new streams can be dynamically connected to, and existing streams disconnected from, the driver without interfering with its ongoing operation. The number of streams that can be connected to a multiplexor is developer dependent. However, there is a system limit, `NMUXLINK`, to the number of streams that can be linked in the system (please refer to *SunOS STREAMS Topics Tunable Parameters* section).

Disconnecting Lower Streams

Dismantling a lower multiplexor is accomplished by disconnecting (unlinking) the lower streams. Unlinking can be initiated in three ways: an `I_UNLINK ioctl()` referencing a specific stream, an `I_UNLINK` indicating all lower streams, or the last `close()` (i.e., causes the associated file to be closed) of the control stream. Note that if the connection was established with `I_PLINK`, then it needs to be undone with `I_PUNLINK`. As in the link, an unlink sends a

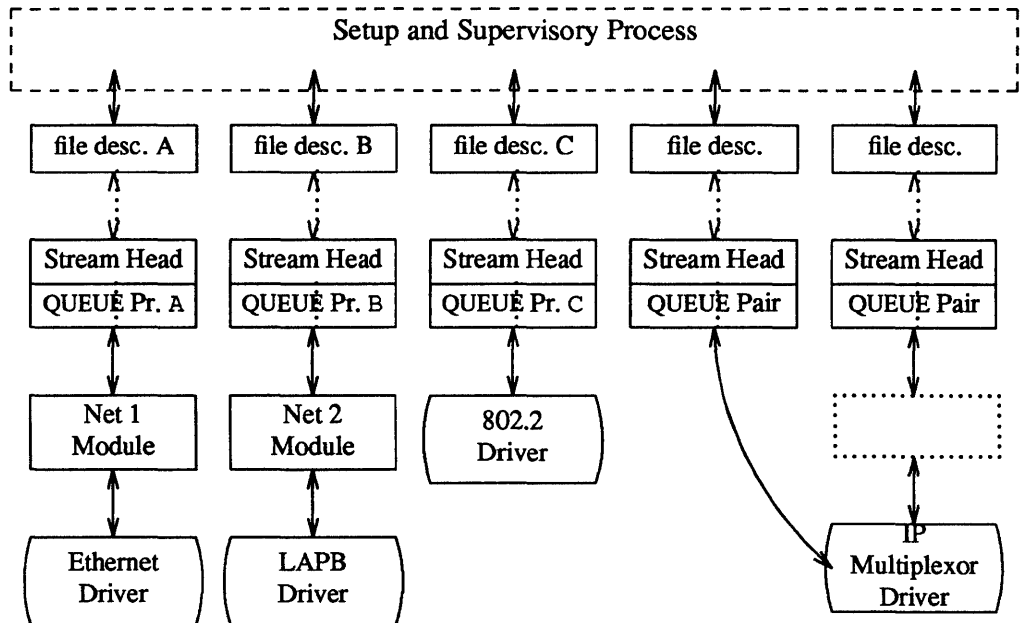
linkblk structure to the driver in an M_IOCTL message. The I_UNLINK call, which unlinks a single stream, uses the l_index value returned in the I_LINK to specify the lower stream to be unlinked. The latter two calls must designate a file corresponding to a control stream. This causes all the lower streams that were previously linked (by this control stream) to be unlinked. However, the driver sees a series of individual unlinks.

If the file descriptor for a lower stream was previously closed, a subsequent unlink will automatically close the stream. Otherwise, the lower stream must be closed by close() following the unlink. STREAMS will automatically dismantle all cascaded multiplexors (below other multiplexing streams) if their controlling stream is closed. This holds true for connections established with I_LINK but not true for I_PLINK, which is the whole point for I_PLINK. An I_UNLINK will leave lower, cascaded multiplexing streams intact unless the stream file descriptor was previously closed.

Multiplexor Construction Example

This section describes an example of multiplexor construction and usage. A multiplexing configuration similar to the Internet figure in the *Other Facilities* section of the *Introduction to STREAMS* is discussed. Figure 3-7 shows the streams before their connection to create the multiplexing configuration of figure 3-8. Multiple upper and lower streams interface to the multiplexor driver. The user processes of figure 3-8 are not shown in figure 3-7.

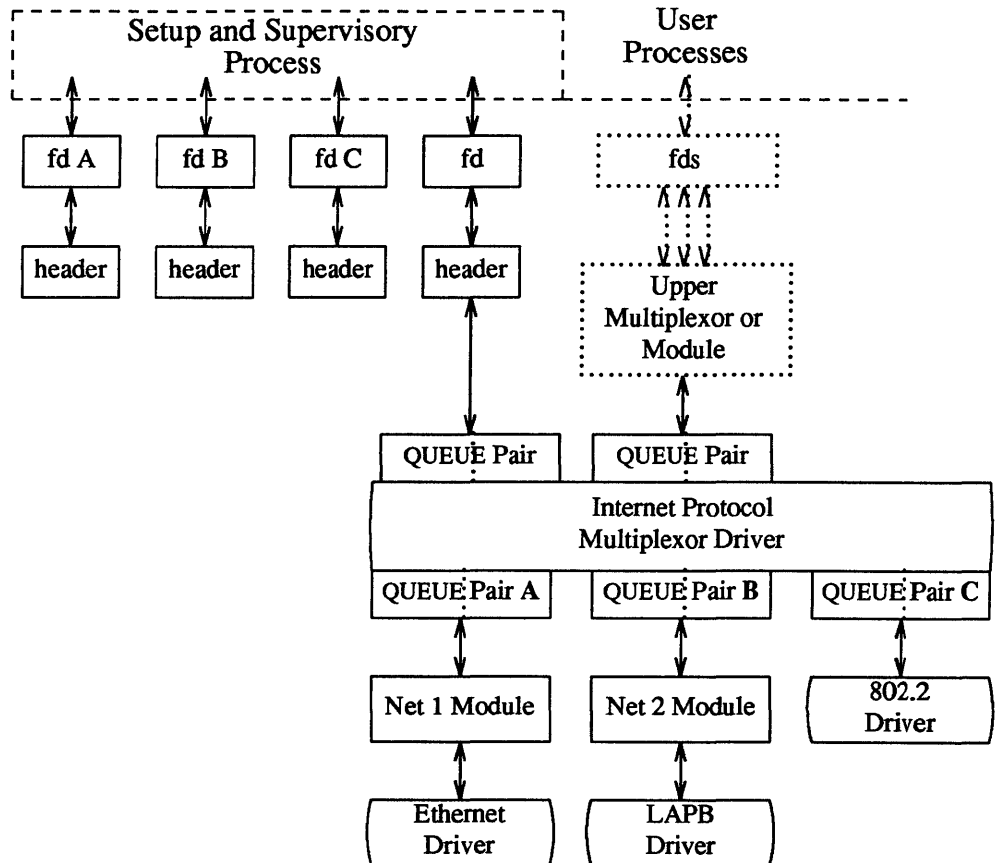
Figure 3-7 *Internet Multiplexor Before Connecting*



The Ethernet, LAPB and IEEE 802.2 device drivers terminate links to other nodes. IP (Internet Protocol) is a multiplexor driver. IP switches datagrams among the various nodes or sends them upstream to a user(s) in the system. The Net modules would typically provide a convergence module that matches the IP and device driver interface.

Figure 3-7 depicts only a portion of the full, larger stream. As shown in the dotted rectangle above the IP multiplexor, there generally would be an upper TCP multiplexor, additional modules and, possibly, additional multiplexors in the stream. Multiplexors could also be cascaded below the IP driver if the device drivers were replaced by multiplexor drivers.

Figure 3-8 *Internet Multiplexor After Connecting*



Streams A, B, and C are opened by the process, and modules are pushed as needed. Two upper streams are opened to the IP multiplexor. The rightmost stream represents multiple streams, each connected to a process using the network. The stream second from the right provides a direct path to the multiplexor for supervisory functions. It is the control stream, leading to a process that sets up and supervises this configuration. It is always directly connected to the IP driver. Although not shown, modules can be pushed on the control stream.

After the streams are opened, the supervisory process typically transfers routing information to the IP drivers (and any other multiplexors above the IP), and initializes the links. As each link becomes operational, its Stream is connected below the IP driver. If a more complex multiplexing configuration is required, the IP multiplexor stream with all its connected links can be connected below another multiplexor driver.

As shown in figure 3-8, the file descriptors for the lower device driver streams are left dangling. The primary purpose in creating these streams was to provide parts for the multiplexor. Those not used for control and not required for error recovery (by reconnecting them through an `I_UNLINK ioctl()`) have no further function. As stated above, these lower streams can be closed to free the file descriptor without any effect on the multiplexor. A setup process installing a configuration containing a large number of drivers should do this to avoid running out of file descriptors.

Multiplexing Driver

This section contains an example of a multiplexing driver that implements an N-to-1 configuration. This configuration might be used for terminal windows, where each transmission to or from the terminal identifies the window. This resembles a typical device driver, with two differences: the device handling functions are performed by a separate driver, connected as a lower stream, and the device information (i.e., relevant user process) is contained in the input data rather than in an interrupt call.

Each upper stream is connected by an `open(2)`. A single lower stream is opened and then it is linked by use of the multiplexing facility. This lower stream might connect to the tty driver. The implementation of this example is a foundation for an M to N multiplexor.

As in the loop-around driver, flow control requires the use of standard and special code, since physical connectivity among the streams is broken at the driver. Different approaches are used for flow control on the lower stream, for messages coming upstream from the device driver, and on the upper streams, for messages coming downstream from the user processes.

The multiplexor declarations are:

```
#include "mux.h"
#include <sys/types.h>
#include <sys/param.h>
#include <sys/sysmacros.h>
#include <sys/stream.h>
#include <sys/stropts.h>
#include <sys/errno.h>

static int muxopen(), muxclose(), muxuwput();
static int muxlwsrv(), muxlrput();

static struct module_info info = {
    0, "mux", 0, INFPSZ, 512, 128
};

static struct qinit urinit = { /* upper read */
    NULL, NULL, muxopen, muxclose, NULL, &info, NULL
};

static struct qinit uwinit = { /* upper write */
    muxuwput, NULL, NULL, NULL, NULL, &info, NULL
};

static struct qinit lrinit = { /* lower read */
    muxlrput, NULL, NULL, NULL, NULL, &info, NULL
};
```

```

};
static struct qinit lwinit = { /* lower write */
    NULL, muxlwsrv, NULL, NULL, NULL, &info, NULL
};

struct streamtab muxinfo =
    {&urinit, &uwinit, &lrinit, &lwinit};

struct mux {
    queue_t *qptr; /* back pointer to read queue */
};

struct mux mux_mux[NMUX];
int mux_cnt = NMUX;

queue_t *muxbot; /* linked lower queue */
int muxerr; /* set if error of hangup on lower stream */

static queue_t *get_next_q();

```

The four `streamtab` entries correspond to the upper read, upper write, lower read, and lower write `qinit` structures. The multiplexing `qinit` structures replace those in each (in this case there is only one) lower Stream Head after the `I_LINK` has completed successfully. In a multiplexing configuration, the processing performed by the multiplexing driver can be partitioned between the upper and lower QUEUES. There must be upper stream write and lower stream read put procedures. Typically, only upper write side and lower read side procedures are used. Application-specific flow control requirements might modify this. If the QUEUE procedures of the opposite upper/lower QUEUE are not needed, the QUEUE can be skipped over, and the message put to the following QUEUE.

In the example, the upper read side procedures are not used. The lower stream read QUEUE put procedure transfers the message directly to the read QUEUE upstream from the multiplexor. There is no lower write put procedure because the upper write put procedure directly feeds the lower write service procedure, as described below.

The driver uses a private data structure, `mux_mux`. `mux_mux[dev]` points back to the opened upper read QUEUE. This is used to route messages coming upstream from the driver to the appropriate upper QUEUE. It is also used to find a free minor device for a `CLONEOPEN` driver open case.

The upper QUEUE open contains the canonical driver open code:

```
static int muxopen(q, dev, flag, sflag)
queue_t *q;
{
    struct mux *mux;

    if (sflag == CLONEOPEN) {
        for (dev = 0; dev < mux_cnt; dev++)
            if (mux_mux[dev].qptr == NULL)
                break;
    }
    else
        dev = minor(dev);

    if (dev >= mux_cnt)
        return (OPENFAIL);

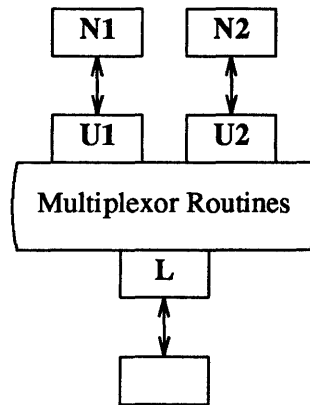
    mux = &mux_mux[dev];
    mux->qptr = q;
    q->q_ptr = (char *) mux;
    WR(q)->q_ptr = (char *) mux;
    return (dev);
}
```

`muxopen` checks for a clone or ordinary open call. It loads `q_ptr` to point at the `mux_mux[]` structure.

The core multiplexor processing is the following: downstream data written to an upper stream is queued on the corresponding upper write message queue. This allows flow control to propagate towards the Stream Head for each upper stream. However, there is no service procedure on the upper write side. All `M_DATA` messages from all the upper message queues are ultimately dequeued by the service procedure on the lower (linked) write side. The upper write streams are serviced in a round-robin fashion by the lower write service procedure. A lower write service procedure, rather than a write put procedure, is used so that flow control, coming up from the driver below, may be handled.

On the lower read side, data coming up the lower stream is passed to the lower read put procedure. The procedure routes the data to an upper stream based on the first byte of the message. This byte holds the minor device number of an upper stream. The put procedure handles flow control by testing the upper stream at the first upper read QUEUE beyond the driver. That is, the put procedure treats the stream component above the driver as the next QUEUE.

Figure 3-9 Example Multiplexor Configuration



This is shown (sort of) in figure 3-9. Multiplexor Routines are all the above procedures. U1 and U2 are `queue_t` pairs, each including a write `queue_t` pointed at by an `l_qtop` in a `linkblk` (see the beginning of this section). L is the `queue_t` pair that contains the write `queue_t` pointed at by `l_qbot`. N1 and N2 are the modules (or Stream Head or another multiplexing driver) seen by L when read side messages are sent upstream.

Upper Write Put Procedure

`muxuwput`, the upper QUEUE write put procedure, traps `ioctl`s, in particular `I_LINK` and `I_UNLINK`:

```

static int muxuwput(q, mp)
queue_t *q;
mbblk_t *mp;

{

    int s;
    struct mux *mux;

    mux = (struct mux *)q->q_ptr;
    switch (mp->b_datap->db_type) {
    case M_IOCTL: {
        struct iocblk *iocp;
        struct linkblk *linkp;

        /*
         * Ioctl. Only channel 0 can do ioctls. Two
         * calls are recognized: I_LINK, and I_UNLINK
         */

        if (mux != mux_mux)
            goto iocnak;

        iocp = (struct iocblk *) mp->b_rptr;
        switch (iocp->ioc_cmd) {

```

```

case I_LINK:
    /*
     * Link. The data contains a linkblk structure
     * Remember the bottom queue in muxbot.
     */

    if (muxbot != NULL)
        goto iocnak;
    linkp = (struct linkblk *) mp->b_cont->b_rptr;
    muxbot = linkp->l_qbot;
    muxerr = 0;
    mp->b_datap->db_type = M_IOCACK;
    iocp->ioc_count = 0;
    qreply(q, mp);
    break;
case I_UNLINK:
    /*
     * Unlink. The data contains a linkblk structure.
     * Should not fail an unlink. Null out muxbot.
     */

    linkp = (struct linkblk *) mp->b_cont->b_rptr;
    muxbot = NULL;
    mp->b_datap->db_type = M_IOCACK;
    iocp->ioc_count = 0;
    qreply(q, mp);
    break;
default:
iocnak:

    /* fail ioctl */

    mp->b_datap->db_type = M_IOCNAK;
    qreply(q, mp);
}

break;
}

```

First, there is a check to enforce that the stream associated with minor device 0 will be the single controlling stream. Iocltls are only accepted on this stream. As described previously, a controlling stream is the one that issues the `I_LINK`. Having a single control stream is a recommended practice. `I_LINK` and `I_UNLINK` include a `linkblk` structure, described previously, containing:

`l_qtop`

The upper write QUEUE from which the ioctl is coming. It should always equal `q`.

l_qbot

The new lower write QUEUE. It is the former Stream Head write QUEUE. It is of most interest since that is where the multiplexor gets and puts its data.

l_index

A unique (system wide) identifier for the link. It can be used for routing, or during selective unlinks, as described above. Since the example only supports a single link, `l_index` is not used.

For `I_LINK`, `l_qbot` is saved in `muxbot` and an ack is generated. From this point on, until an `I_UNLINK` occurs, data from upper queues will be routed through `muxbot`. Note that when an `I_LINK` is received, the lower stream has already been connected. This allows the driver to send messages downstream to perform any initialization functions. Returning an `M_IOCNAK` message (nak) in response to an `I_LINK` will cause the lower stream to be disconnected.

The `I_UNLINK` handling code nulls out `muxbot` and generates an ack. A nak should not be returned to an `I_UNLINK`. The Stream Head assures that the lower stream is connected to a multiplexor before sending an `I_UNLINK M_IOCTL`.

`muxwput` handles `M_FLUSH` messages as a normal driver would:

```

case M_FLUSH:
    if (*mp->b_rptr & FLUSHW)
        flushq(q, FLUSHDATA);
    if (*mp->b_rptr & FLUSHR) {
        flushq(RD(q), FLUSHDATA);
        *mp->b_rptr &= ~FLUSHW;
        qreply(q, mp);
    } else
        freemsg(mp);
    break;
case M_DATA:
    /*
     * Data. If we have no bottom queue --> fail
     * Otherwise, queue the data, and invoke the lower
     * service procedure.
     */
    if (muxerr || muxbot == NULL)
        goto bad;
    putq(q, mp); /* place message on upper write message queue */
    qenable(muxbot); /* lower service write procedure */
    break;
default:
bad:
    /*
     * Send an error message upstream.
     */
    mp->b_datap->db_type = M_ERROR;
    mp->b_rptr = mp->b_wptr = mp->b_datap->db_base;
    *mp->b_wptr++ = EINVAL;
    qreply(q, mp);
}

```

```
}

```

M_DATA messages are not placed on the lower write message queue. They are queued on the upper write message queue. `putq()` recognizes the absence of the upper service procedure and does not schedule the QUEUE. Then, the lower service procedure, `muxlwsrv` is scheduled with `qenable()` (see *Utilities* in the *Supplementary STREAMS Material* chapter) to start output. This is similar to starting output on a device driver. Note that `muxwput` can not access `muxlwsrv` (the lower QUEUE write service procedure, contained in `muxbot`) by the conventional STREAMS calls, `putq()` or `putnext()` (to a `muxlwput`). Both calls require that a message be passed, but the messages remain on the upper stream.

Lower QUEUE Write Service Procedure

`muxlwsrv`, the lower (linked) queue write service procedure is scheduled directly from the upper service procedures. It is also scheduled from the lower stream, by being back-enabled when the lower stream becomes unblocked from downstream flow control.

```
static int muxlwsrv(q)
register queue_t *q;
{
    register mblk_t *mp, *bp;
    register queue_t *nq;

    /*
     * While lower stream is not blocked, find an upper queue to
     * service (get_next_q) and send one message from it downstream.
     */
    while (canput(q->q_next)) {
        nq = get_next_q();
        if (nq == NULL)
            break;
        mp = getq(nq);
        /*
         * Prepend the outgoing message with a single byte header
         * that indicates the minor device number it came from.
         */
        if ((bp = allocb(1, BPRI_MED)) == NULL) {
            printf("mux: allocb failed (size 1)\n");
            freemsg(mp);
            continue;
        }
        *bp->b_wptr++ = (struct mux *)nq->q_ptr - mux_mux;
        bp->b_cont = linkb(bp, mp);
        putnext(q, bp);
    }
}
```

`muxlwsrv` takes data from the upper queues and puts it out through `muxbot`. The algorithm used is simple round robin. While we can put to

`muxbot->q_next`, we select an upper QUEUE (via `get_next_q`) and move a message from it to `muxbot`. Each message is prepended with a one byte header that indicates which upper stream it came from.

Finding messages on upper write queues is handled by `get_next_q()`:

```

/*
 * Round-robin scheduling.
 * Return next upper queue that needs servicing.
 * Returns NULL when no more work needs to be done.
 */

static queue_t *
get_next_q()
{
    static int next;
    int i, start;
    register queue_t *q;

    start = next;
    for (i = next; i < mux_cnt; i++)
        if (q = mux_mux[i].qp_ptr) {
            q = WR(q);
            if (q->q_first) {
                next = i+1;
                return (q);
            }
        }

    for (i = 0; i < start; i++)
        if (q = mux_mux[i].qp_ptr) {
            q = WR(q);
            if (q->q_first) {
                next = i+1;
                return (q);
            }
        }

    return (NULL);
}

```

`get_next_queue()` searches the upper queues in a round robin fashion looking for the first one containing a message. It returns the `queue_t` pointer or

NULL if there is no work to do.

Lower Read Put Procedure

The lower (linked) queue read put procedure is:

```
static int muxlrput(q, mp)
queue_t *q;
mblk_t *mp;
{
    queue_t *uq;
    mblk_t *b_cont;
    int dev;

    switch(mp->b_datap->db_type) {
    case M_FLUSH:

        /*
         * Flush queues. NOTE: sense of tests is reversed
         * since we are acting like a "stream head"
         */

        if (*mp->b_rptr & FLUSHR)
            flushq(q, 0);
        if (*mp->b_rptr & FLUSHW) {
            *mp->b_rptr &= ~FLUSHR;
            qreply(q, mp);
        } else
            freemsg(mp);
        break;

    case M_ERROR:
    case M_HANGUP:
        muxerr = 1;
        freemsg(mp);
        break;

    case M_DATA:
        /*
         * Route message. First byte indicates
         * device to send to. No flow control.
         *
         * Extract and delete device number. If the leading block is
         * now empty and more blocks follow, strip the leading block.
         * The stream head interprets a leading zero length block
         * as an EOF regardless of what follows (sigh).
         */

        dev = *mp->b_rptr++;
        if (mp->b_rptr == mp->b_wptr &&
            (b_cont = mp->b_cont)) {
            freeb(mp);
            mp = b_cont;
        }
    }
}
```

```

    /* Sanity check. Device must be in range */

    if (dev < 0 || dev >= mux_cnt) {
        freemsg(mp);
        break;
    }

    /*
     * If upper stream is open and not backed up,
     * send the message there, otherwise discard it.
     */

    uq = mux_mux[dev].qp_ptr;
    if (uq != NULL && canput(uq->q_next))
        putnext(uq, mp);
    else
        freemsg(mp);
    break;
default:
    freemsg(mp);
}
}

```

`muxlrput` receives messages from the linked stream. In this case, it is acting as a Stream Head. It handles `M_FLUSH` messages. Note the code is reversed from that of a driver, handling `M_FLUSH` messages from upstream.

`muxlrput` also handles `M_ERROR` and `M_HANGUP` messages. If one is received, it locks up the upper streams.

`M_DATA` messages are routed by looking at the first data byte of the message. This byte contains the minor device of the upper stream. If removing this byte causes the leading block to be empty, and more blocks follow, the block is discarded. This is done because the Stream Head interprets a leading zero length block as an EOF [see `read(2)`]. Several sanity checks are made: Does the message have at least one byte? Is the device in range? Is the upper stream open? Is the upper stream not full?

This mux does not do end-to-end flow control. It is merely a router (like the Department of Defense's IP protocol). If everything checks out, the message is put to the proper upper QUEUE. Otherwise, the message is silently discarded.

The upper stream close routine simply clears the mux entry so this queue will no

longer be found by `get_next_queue()`:

```

/*
 * Upper queue close
 */
static int muxclose(q)
queue_t *q;
{
    ((struct mux *)q->q_ptr)->qptr = NULL;
}

```

3.9. Service Interface

Definition

STREAMS provides the means to implement a service interface between any two components in a stream, and between a user process and the topmost module in the stream. A service interface is defined at the boundary between a service user and a service provider. A service interface is a set of primitives and the rules for the allowable sequences of primitives across the boundary. These rules are typically represented by a state machine. In STREAMS, the service user and provider are implemented in a module, driver, or user process. The primitives are carried bidirectionally between a service user and provider in `M_PROTO` and `M_PCPROTO` (generically, `PROTO`) messages. `M_PCPROTO` is the high priority version of `M_PROTO`.

Message Usage

As described in the *Message Types* section of the *Supplementary STREAMS Material* chapter), `PROTO` messages can be multi-block, with the second through last blocks of type `M_DATA`. The first block in a `PROTO` message contains the control part of the primitive in a form agreed upon by the user and provider and the block is not intended to carry protocol headers. (Although its use is not recommended, upstream `PROTO` messages can have multiple `PROTO` blocks at the start of the message. `getmsg()` will compact the blocks into a single control part when sending to a user process.) The `M_DATA` block(s) contains any data part associated with the primitive. The data part may be processed in a module that receives it, or it may be sent to the next stream component, along with any data generated by the module. The contents of `PROTO` messages and their allowable sequences are determined by the service interface specification.

`PROTO` messages can be sent bidirectionally (up and downstream) on a stream and bidirectionally between a stream and a user process. `putmsg(2)` and `getmsg(2)` system calls are analogous, respectively, to `write(2)` and `read(2)` except that the former allow both data and control parts to be (separately) passed, and they observe message boundary alignment across the user-stream boundary. `putmsg()` and `getmsg()` separately copy the control part (`M_PROTO` or `M_PCPROTO` block) and data part (`M_DATA` blocks) between the stream and user process.

An `M_PCPROTO` message is normally used to acknowledge `M_PROTO` messages and not to carry protocol expedited data. `M_PCPROTO` insures that the

acknowledgement reaches the service user before any other message. If the service user is a user process, the Stream Head will only store a single `M_PCPROTO` message, and discard subsequent `M_PCPROTO` messages until the first one is read with `getmsg(2)`.

The following rules pertain to service interfaces:

- Modules and drivers that support a service interface must act upon all `PROTO` messages and not pass them through.
- Modules may be inserted between a service user and a service provider to manipulate the data part as it passes between them. However, these modules may not alter the contents of the control part (`PROTO` block, first message block) nor alter the boundaries of the control or data parts. That is, the message blocks comprising the data part may be changed, but the message may not be split into separate messages nor combined with other messages. In addition, modules and drivers must observe the rule that high priority messages are not subject to flow control and forward them accordingly (e.g., see the beginning of `modwsrv()` in the *Message Queues and Service Procedures* section). High priority messages also bypass flow control at the user-stream boundary (e.g., see `putmsg(2)`).

Example

The example below is part of a module that illustrates the concept of a service interface. The module implements a simple datagram interface.

Declarations

The service interface primitives are defined in the declarations:

```
#include <sys/types.h>
#include <sys/param.h>
#include <sys/stream.h>
#include <sys/errno.h>

/*
 * Primitives initiated by the service user:
 */
#define BIND_REQ          1    /* bind request */
#define UNITDATA_REQ     2    /* unit data request */
/*
 * Primitives initiated by the service provider:
 */
#define OK_ACK           3    /* bind acknowledgment */
#define ERROR_ACK       4    /* error acknowledgment */
#define UNITDATA_IND     5    /* unit data indication */
/*
 * The following structures define the format of the
 * stream message block of the above primitives.
 */
struct bind_req {           /* bind request */
    long    PRIM_type;     /* always BIND_REQ */
    long    BIND_addr;    /* addr to bind */
};
```

```

struct unitdata_req {      /* unitdata request */
    long    PRIM_type;    /* always UNITDATA_REQ */
    long    DEST_addr;    /* dest addr      */
};
struct ok_ack {           /* ok acknowledgment */
    long    PRIM_type;    /* always OK_ACK  */
};
struct error_ack {       /* error acknowledgment */
    long    PRIM_type;    /* always ERROR_ACK */
    long    UNIX_error;   /* SunOS error code */
};
struct unitdata_ind {    /* unitdata indication */
    long    PRIM_type;    /* always UNITDATA_IND */
    long    SRC_addr;     /* source addr      */
};
union primitives {       /* union of all primitives */
    long    type;
    struct bind_req bind_req;
    struct unitdata_req unitdata_req;
    struct ok_ack ok_ack;
    struct error_ack error_ack;
    struct unitdata_ind unitdata_ind;
};
struct dgproto {         /* structure per minor device */
    short state;         /* current provider state */
    long  addr;          /* net address */
};
/* Provider states */

#define IDLE    0
#define BOUND  1

```

In general, the M_PROTO or M_PCPROTO block is described by a data structure containing the service interface information. In this example, union primitives is that structure.

Two commands are recognized by the module:

BIND_REQ

Give this stream a protocol address, i.e. give it a name on the network. After a BIND_REQ is completed, datagrams from other senders will find their way through the network to this particular stream.

UNITDATA_REQ

Send a datagram to the specified address.

Three messages are generated:

OK_ACK

A positive acknowledgement (ack) of BIND_REQ.

ERROR_ACK

A negative acknowledgement of BIND_REQ.

UNITDATA_IND

A datagram from the network has been received (this code is not shown).

The ack of a BIND_REQ informs the user that the request was syntactically correct (or incorrect if ERROR_ACK). The receipt of a BIND_REQ is acknowledged with an M_PCPROTO to insure that the acknowledgement reaches the user before any other message. For example, a UNITDATA_IND could come through before the bind has completed, and the user would get confused.

The driver uses a per-minor device data structure, dgproto, which contains the following:

state

current state of the stream (endpoint) IDLE or BOUND

addr

network address that has been bound to this stream

It is assumed (though not shown) that the module open procedure sets the write queue `q_ptr` to point at one of these structures.

Service Interface Procedure

The write put procedure is:

```
static int protowput(q, mp)
queue_t *q;
mblock_t *mp;
{
    union primitives *proto;
    struct dgproto *dgproto;
    int err;

    dgproto = (struct dgproto *) q->q_ptr;

    switch (mp->b_datap->db_type) {
    default:
        /* don't understand it */
        mp->b_datap->db_type = M_ERROR;
        mp->b_rptr = mp->b_wptr = mp->b_datap->db_base;
        *mp->b_wptr++ = EPROTO;
        qreply(q, mp);
        break;
    case M_FLUSH:
        /* standard flush handling goes here ... */
        break;
    case M_PROTO:
        /* Protocol message -> user request */
        proto = (union primitives *) mp->b_rptr;
        switch (proto->type) {
        default:
            mp->b_datap->db_type = M_ERROR;
            mp->b_rptr = mp->b_wptr = mp->b_datap->db_base;
            *mp->b_wptr++ = EPROTO;
            qreply(q, mp);
        }
    }
}
```

```

        return;

    case BIND_REQ:
        if (dgproto->state != IDLE) {
            err = EINVAL;
            goto error_ack;
        }
        if (mp->b_wptr - mp->b_rptr
            != sizeof(struct bind_req)) {
            err = EINVAL;
            goto error_ack;
        }
        if (err = chkaddr(proto->bind_req.BIND_addr))
            goto error_ack;

        dgproto->state = BOUND;
        dgproto->addr = proto->bind_req.BIND_addr;
        mp->b_datap->db_type = M_PCPROTO;
        proto->type = OK_ACK;
        mp->b_wptr =
            mp->b_rptr + sizeof(struct ok_ack);
        qreply(q, mp);
        break;

    error_ack:
        mp->b_datap->db_type = M_PCPROTO;
        proto->type = ERROR_ACK;
        proto->error_ack.UNIX_error = err;
        mp->b_wptr =
            mp->b_rptr + sizeof(struct error_ack);
        qreply(q, mp);
        break;

    case UNITDATA_REQ:
        if (dgproto->state != BOUND)
            goto bad;
        if (mp->b_wptr - mp->b_rptr
            != sizeof(struct unitdata_req))
            goto bad;
        if (err=chkaddr(proto->unitdata_req.DEST_addr))
            goto bad;
        if (mp->b_cont) {
            putq(q, mp->b_cont);

            /* start device or mux output ... */
        }
        break;

    bad:
        freemsg(mp);
        break;
}
}
}

```

The write put procedure switches on the message type. The only types accepted are `M_FLUSH` and `M_PROTO`. For `M_FLUSH` messages, the driver will perform the canonical flush handling (not shown). For `M_PROTO` messages, the driver assumes the message block contains a union primitive and switches on the `type` field. Two types are understood: `BIND_REQ`, and `UNITDATA_REQ`.

For a `BIND_REQ`, the current state is checked; it must be `IDLE`. Next, the message size is checked. If it is the correct size, the passed-in address is verified for legality by calling `chkaddr`. If everything checks, the incoming message is converted into an `OK_ACK` and sent upstream. If there was any error, the incoming message is converted into an `ERROR_ACK` and sent upstream.

For `UNITDATA_REQ`, the state is also checked; it must be `BOUND`. As above, the message size and destination address are checked. If there is any error, the message is simply discarded. (This action may seem rash, but it is in accordance with the interface specification, which is not shown. Another specification might call for the generation of a `UNITDATA_ERROR` indication.) If all is well, the data part of the message, if it exists, is put on the queue, and the lower half of the driver is started.

If the write put procedure receives a message type that it does not understand, either a bad `b_datap->db_type` or a bad `proto->type`, the message is converted into an `M_ERROR` message and sent upstream.

Another piece of code not shown is the generation of `UNITDATA_IND` messages. This would normally occur in the device interrupt if this is a hardware driver (like Ethernet) or in the lower read put procedure if this is a multiplexor. The algorithm is simple: The data part of the message is appended to an `M_PROTO` message block that contains a `unitdata_ind` structure and sent upstream.

3.10. Advanced Topics

Recovering From No Buffers

The `bufcall()` utility (see *Utilities* in the *Supplementary STREAMS Material* chapter) is used to recover from an `allocb()` failure. The call syntax is as follows:

```
int bufcall(size, pri, func, arg);
int size, pri, (*func)();
long arg;
```

`bufcall()` will call `(*func)(arg)` when a buffer of `size` bytes at `pri` priority is available. When `func` is called, it has no user context and must return without sleeping. Also, because of interrupt processing, there is no guarantee that when `func` is called, a buffer will actually be available (someone else may steal it). `bufcall()` returns a nonzero integer on success, indicating that the request has been successfully recorded, or 0 on failure. On a failure return, the requested function will never be called. The caller should save the return value for possible subsequent use as the argument to `unbufcall()`, which is used to

Care must be taken to avoid deadlock when holding resources while waiting for `bufcall()` to call `(*func)(arg).bufcall()` should be used sparingly.

cancel an outstanding `bufcall()` request, for example within a module's close routine.

Two examples are provided. Example one is a device receive interrupt handler:

```
#include <sys/types.h>
#include <sys/param.h>
#include <sys/stream.h>

dev_rintr(dev)
{
    /* process incoming message ... */

    /* allocate new buffer for device */
    dev_re_load(dev);
}
/*
 * Reload device with a new receive buffer
 */
dev_re_load(dev)
{
    mblk_t *bp;

    if ((bp = allocb(DEVBLKSZ, BPRI_MED)) == NULL) {
        log(LOG_ERR("dev: allocb failure (size %d)\n",
            DEVBLKSZ);
        /*
         * Allocation failed. Use bufcall to
         * schedule a call to ourself.
         */
        (void) bufcall(DEVBLKSZ, BPRI_MED, dev_re_load,
            dev);
        return;
    }

    /* pass buffer to device ... */
}
```

`dev_rintr` is called when the device has posted a receive interrupt. The code retrieves the data from the device (not shown). `dev_rintr` must then give the device another buffer to fill by a call to `dev_re_load`, which calls `allocb()` with the appropriate buffer size (`DEVBLKSZ`, definition not shown) and priority. If `allocb()` fails, `dev_re_load` uses `bufcall()` to call itself when STREAMS determines a buffer of the appropriate size and priority is available.

NOTE *Since `bufcall()` may fail, there is still a chance that the device may hang. A better strategy, in the event `bufcall()` fails, would be to discard the current input message and resubmit that buffer to the device. Losing input data is generally better than hanging.*

The second example is a write service procedure, `mod_wsrv()`, which needs to prepend each output message with a header (similar to the multiplexor example of the *Multiplexing* section). `mod_wsrv()` illustrates a case for potential deadlock:

```
static int mod_wsrv(q)
queue_t *q;
{
    int qenable();
    mblk_t *mp, *bp;
    while (mp = getq(q)) {
        /* check for priority messages and canput ... */
        /*
         * Allocate a header to prepend to the message. If
         * the allocb fails, use bufcall to reschedule ourself.
         */
        if ((bp = allocb(HDRSZ, BPRI_MED)) == NULL) {
            if (!bufcall(HDRSZ, BPRI_MED, qenable, q)) {
                /*
                 * The bufcall request has failed. Discard
                 * the message and keep running to avoid hanging.
                 */
                freemsg(mp);
                continue;
            }
            /*
             * Put the message back and exit, we will be re-enabled later
             */
            putbq(q, mp);
            return;
        }
        /* process message .... */
    }
}
```

However, if `allocb()` fails, `mod_wsrv()` wants to recover without loss of data and calls `bufcall()`. In this case, the routine passed to `bufcall()` is `qenable()` (see below and in the *Utilities* section of the *Supplementary STREAMS Material* chapter). When a buffer is available (of size `HDRSZ`, definition not shown), the service procedure will be automatically re-enabled. Before exiting, the current message is put back on the queue. This example deals with `bufcall()` failure by discarding the current message and continuing in the service procedure loop.

A module that issues a `bufcall()` request should be prepared to cancel that request (with `unbufcall()`) when its close routine is called. If it doesn't, the kernel is likely to crash later on when STREAMS calls the function given as `bufcall()`'s argument at a time when the data structures it manipulates are

no longer valid.

Advanced Flow Control

STREAMS provides mechanisms to alter the normal queue scheduling process. `putq()` will not schedule a QUEUE if `noenable(q)` had been previously called for this QUEUE. `noenable()` instructs `putq()` to queue the message when called by this QUEUE, but not to schedule the service procedure. `noenable()` does not prevent the QUEUE from being scheduled by a flow control back-enable. The inverse of `noenable()` is `enableok(q)`.

An example of this is driver upstream flow control. Although device drivers typically discard input when unable to send it to a user process, STREAMS allows driver read side flow control, possibly for handling temporary upstream blocks. This is done through a driver read service procedure which is disabled during the driver open with `noenable()`. If the driver input interrupt routine determines messages can be sent upstream (from `canput()`), it sends the message with `putnext()`. Otherwise, it calls `putq()` to queue the message. The message waits on the message queue (possibly with queue length checked when new messages are enqueued by the interrupt routine) until the upstream QUEUE becomes unblocked. When the blockage abates, STREAMS back-enables the driver read service procedure. The service procedure sends the messages upstream using `getq()` and `canput()`, as in *Message Queues and Service Procedures*. This is similar to `looprsvr()` in the *Complete Driver* where the service procedure is present only for flow control.

`qenable()`, another flow control utility, allows a module or driver to cause one of its QUEUES, or another module's QUEUES, to be scheduled. In addition to the usage shown in the *Complete Driver* and *Multiplexing* sections, `qenable()` might be used when a module or driver wants to delay message processing for some reason. An example of this is a buffer module that gathers messages in its message queue and forwards them as a single, larger message. This module uses `noenable()` to inhibit its service procedure and queues messages with its `put` procedure until a certain byte count or "in queue" time has been reached. When either of these conditions is met, the `put` procedure calls `qenable()` to cause its service procedure to run.

Another example is a communication line discipline module that implements end-to-end (i.e., to a remote system) flow control. Outbound data is held on the write side message queue until the read side receives a transmit window from the remote end of the network. Then, the read side schedules the write side service procedure to run.

Signals

STREAMS allows modules and drivers to cause a signal to be sent to user process(es) through an `M_SIG` or `M_PCISIG` message (see *Message Types* in the *Supplementary STREAMS Material* chapter) sent upstream. `M_PCISIG` is a high priority version of `M_SIG`. For both messages, the first byte of the message specifies the signal for the Stream Head to generate. If the signal is not `SIGPOLL` (see `signal(2)` and `sigset(2)`), then the signal is sent to the process group associated with the stream (see below). If the signal is `SIGPOLL`, the signal is only sent to processes that have registered for the signal by using the `I_SETSIG` `ioctl(2)` (see also the `streamio(4)` call).

A process group is associated with a stream during the open of the driver or module. If the `NEWCTTY` flag is ORed into the value returned by the `open()` procedure, the process on whose behalf the module or driver is being opened becomes a “session process group leader” by executing the `setspgldr()` call (which is executed by the `setpgrp()` call in the System V environment, but not in the 4BSD environment). If that process does not already have a controlling tty, and the stream does not already have a process group, then the stream is assigned to the process group that the process is the leader of and becomes that process’s controlling tty.

If the driver or module wants to have a process group associated with the stream, it should OR the `NEWCTTY` flag into its return value.

`M_SIG` can be used by modules or drivers that wish to insert an explicit in band signal into a message stream. For example, an `M_SIG` message can be sent to the user process immediately before a particular service interface message to gain the immediate attention of the user process. When the `M_SIG` reaches the head of the Stream Head read message queue, a signal will be generated and the `M_SIG` message will be removed. This leaves the service interface message as the next message to be processed by the user. Use of `M_SIG` would typically be defined as part of the service interface of the driver or module.

Control of Stream Head Processing

The `M_SETOPTS` message (see *Message Types* in the *Supplementary STREAMS Material* chapter) allows a driver or module to exercise control over certain Stream Head processing. An `M_SETOPTS` can be sent upstream at any time. The Stream Head responds to the message by altering the processing associated with certain system calls. The options to be modified are specified by the contents of the `stroptions` structure (see *Message Types*) contained in the message.

Six Stream Head characteristics can be modified. As described in *Message Types*, four correspond to fields contained in `queue_t` (min/max packet sizes and high/low water marks). The other two are discussed here.

Read Options

The value for read options (`so_readopt`) corresponds to the three modes a user can set via the `I_SRDOPT` `ioctl()` (see `streamio`) call:

byte-stream (**RNORM**)

The `read(2)` call completes when the byte count is satisfied, the Stream Head read queue becomes empty, or a zero length message is encountered. In the last case, the zero length message is put back on the queue. A subsequent `read()` will return 0 bytes.

message non-discard (**RMSGN**)

The `read()` call completes when the byte count is satisfied or at a message boundary, whichever comes first. Any data remaining in the message is put back on the Stream Head read queue.

message discard (**RMSGD**)

The `read()` call completes when the byte count is satisfied or at a message boundary. Any data remaining in the message is discarded.

Byte-stream mode approximately models pipe data transfer. Message non-discard mode approximately models a tty in canonical mode.

Write Offset

The value for write offset (`so_wroff`) is a hook to allow more efficient data handling. It works as follows: In every data message generated by a `write(2)` system call and in the first `M_DATA` block of the data portion of every message generated by a `putmsg(2)` call, the Stream Head will leave `so_wroff` bytes of space at the beginning of the message block. Expressed as a C language construct:

```
bp->b_rptr = bp->b_datap->db_base + write offset.
```

The write offset value must be smaller than the maximum STREAMS message size, `STRMSGSZ` (see *Tunable Parameters* in the *Supplementary STREAMS Material*) Chapter. In certain cases (e.g., if a buffer large enough to hold the `offset+data` is not currently available), the write offset might not be included in the block. To be general, modules and drivers should not assume that the offset exists in a message, but should always check the message.

The intended use of write offset is to leave room for a module or a driver to place a protocol header before user data in the message rather than by allocating and prepending a separate message. This feature is not general, and its use is discouraged. A more general technique is to put protocol header information in a separate message block and link the user data to it.

SunOS STREAMS Topics

4.1. Configuring STREAMS Drivers

The configuration of STREAMS device drivers is not fundamentally different from the configuration of regular device drivers. This section, therefore, presumes familiarity with the *Configuring the Kernel* section of *Writing Device Drivers* manual, which explains in some detail how new drivers are configured into the kernel.

Note that, while STREAMS give programmers a good deal of flexibility in regard to configuration issues, STREAMS drivers and protocol modules must still be precompiled into the kernel. STREAMS drivers are *not* dynamically loadable.

SunOS STREAMS drivers use exactly the same *autoconfiguration* interface as do regular SunOS drivers. This interface is designed to allow drivers (and modules) to easily define their per-instance data structures, using the information supplied by `config`. However, if a given driver or module chooses to use some other scheme for allocating its resources (such as using `kmem_alloc()` when a previously unopened device is opened), it is free to do so. This differs significantly from the System V driver/kernel interface, which arranges for such storage to be allocated elsewhere.

Each character device that is configured into the Sun kernel results in an entry being placed in the kernel `cdevsw` table. Entries for STREAMS drivers are no exception — they too are placed in `cdevsw`. However, since system calls to STREAMS drivers must be processed by the STREAMS routines, their `cdevsw` interface differs from that of non-STREAMS drivers. `config`, it should be noted, knows nothing about STREAMS drivers. It handles them correctly because, as far as it's concerned, they are just regular character drivers. There is nothing in the format of entries in a `config` file that distinguishes STREAMS devices/modules from other character devices.

There is, however, a difference between STREAMS and non-STREAMS `cdevsw` entries, in that STREAMS entries have only the `d_str` field set while other entries *never* have this field set. `d_str` provides the appropriate single entry point for all system calls on STREAMS files, as shown below:

```
extern struct cdevsw {
    .
    .
    .
    struct streamtab *d_str;
} cdevsw[];
```

The `d_str` entry name is formed by appending the string “info” to the STREAMS driver prefix. The “info” entry is a pointer to the driver/module declared `streamtab` structure (see *Kernel Structures*) in Appendix A. The `streamtab` structure contains pointers to the `qinit` structures for the driver/module’s read and write queues. Its declaration must be externally visible:

```
struct streamtab xxinfo = { ...
```

If the driver declares a `streamtab` named `xxinfo`, the `d_str` entry will contain a non-NULL pointer and the kernel will recognize the driver as a STREAMS driver and will call it by way of the appropriate STREAMS routines. If the `d_str` entry is NULL, the normal character I/O `cdevsw` interface will be used. Note that only `streamtab` must be externally visible in STREAMS drivers and modules, since it is used to uniquely identify the appropriate open, close, put, service, and administration routines. These driver/module routines should generally be declared `static`.

Module Configuration

When adding a new STREAMS module to a kernel, one must add an entry to the `fmodsw` array in `/sys/sun/str_conf.c`. This file is analogous to `/sys/sun/conf.c` (see the *Configuring the Kernel* chapter of *Writing Device Drivers*) and its entries should be similarly conditional on the number of module instances being positive. For example, for the `xx` device:

```
#if NXX > 0
extern struct streamtab xx_info;
#endif
    .
    .
    .
struct fmodsw fmodsw[] =
{
    .
    .
    .
    #if NXX > 0
        { "xx", &xx_info),
    #endif
    .
    .
    .
}
```


The first of the two fields in each `fmodsw` entry is the name of the module, which will be used in all STREAMS-related `ioctl()` calls upon this module. The second is a pointer to the module's `streamtab` structure.

Tunable Parameters

Certain system parameters referenced by STREAMS are configurable when building a new kernel. These parameters can be changed only by editing the file `param.c`. In this discussion, the term “queues” refers to `queue_t` structures. The tunable parameters are:

NMUXLINK

Total number of streams in system that can be linked as lower streams to multiplexor drivers (by an `I_LINK ioctl(2)`, see `streamio(4)`).

NSTREVENT

Initial number of internal event cells available in system to support `bufcall()` and `poll(2)` calls.

MAXSEPGCNT

The number of additional pages of memory that can be dynamically allocated for event cells. If this value is 0, only the allocation defined by `NSTREVENT` is available for use. If the value is not 0 and if the kernel runs out of event cells, it will under some circumstances attempt to allocate an extra page of memory from which new event cells can be created. `MAXSEPGCNT` places a limit on the number of pages that can be allocated for this purpose. Once a page has been allocated for event cells, however, it cannot be recovered later for use elsewhere.

NSTRPUSH

Maximum number of modules that may be pushed onto a single stream.

STRMSGSZ

Maximum bytes of information that a single system call can pass to a stream to be placed into the data part of a message (in `M_DATA` blocks). Any `write(2)` exceeding this size will be broken into multiple messages. A `putmsg(2)` with a data part exceeding this size will fail.

STRCTLSZ

Maximum bytes of information that a single system call can pass to a stream to be placed into the control part of a message (in an `M_PROTO` or `M_PCPROTO` block). A `putmsg(2)` with a control part exceeding this size will fail.

System Error Messages

Messages are reported to the console as a result of various error conditions detected by STREAMS. These messages and the action to be taken on their occurrence are described below. In certain cases, a tunable parameter (see previous section) may have to be changed.

allocq: out of queues

A pair of queues could not be allocated for the Stream Head during the `open()` of a driver, or a pair of queues could not be allocated for a pushable module (`I_PUSH ioctl`). This error message should never be seen, as additional space for queues is allocated dynamically when needed.

allocstr: out of streams

Cannot allocate any more streams.

bufcall: could not allocate stream event

A call to `bufcall()` has failed because all stream event cells have been allocated. If this occurs repeatedly, increase `NSTREVENT`.

esbcall: could not allocate stream event

Could not allocate and stream events.

munlink: could not perform ioctl, closing anyway

A linked multiplexor could not be unlinked when the controlling stream for that link was closed. The linked stream will be unlinked and the controlling stream will be closed anyway.

xdballoc: out of dblks***xmballoc: out of mblks***

The indicated resource could not be allocated during an `allocb()` call. These messages should never be seen as additional `dblks` and `mblks` space is allocated dynamically when needed.

4.2. STREAMS in SunOS

SunOS 4.1 includes implementations of two fundamental system mechanisms in terms of STREAMS. These are:

1. The system terminal driver, which controls serial-line I/O, and
2. The Network Interface Tap (NIT) mechanism, which permits a process to talk to the “raw” Ethernet. NIT is the only networking facility which is thus far implemented in terms of STREAMS, though a TCP/IP implementation that can be accessed via STREAMS is planned.

STREAM Modules

The following STREAMS modules, necessary to support the `tty` driver and the Network Interface Tap, are included in SunOS 4.1.

- The “standard `tty` driver” module, which implements most of the standard `tty` driver behavior; it’s a replacement for the current standard `tty` line discipline. (See `ldterm(4M)`).
- The “`ioctl` mapping” module, which maps old V7 and 4BSD `ioctl()` calls into new-style `ioctl()` calls. This gets pushed on top of the standard `tty` driver module, giving a stream that responds either to the old-style or new-style `ioctl()` calls. (See `ttcompat(4M)`).
- The keyboard and mouse modules, which replace the old keyboard and mouse line disciplines. (See `kb(4M)` and `ms(4M)`).
- The NIT “packet filter” module, which is given a set of criteria for selecting Ethernet packets, and passes only the selected packets upstream, discarding the others. For example, the Reverse ARP daemon requests that it receive only Reverse ARP packets; filtering can be done more efficiently in this fashion than if all packets were handed to the program and it had to do the filtering itself. This also makes it easier to handle a high rate of arrival of packets, since the program doesn’t have to handle the ones it’s not interested

in. (See `nit_pf(4M)`).

- The NIT “buffering” module, which buffers up received Ethernet packets and delivers them to the user program in a single chunk. Such buffering reduces the number of `read()` calls done while monitoring the Ethernet, as is necessary when the rate at which packets arrive is very high. (See `nit_buf(4M)`).

SunOS STREAMS Extension

In order to support STREAMS terminal and pseudo-terminal drivers, SunOS has extended the AT&T STREAMS mechanism. This extension is called “Autopush” and is totally transparent to customer written STREAMS drivers, with the exception of the `streamtab/modlist` initialization. Future versions of STREAMS will replace this feature with the System V implementation.

The Autopush feature includes a mechanism by which STREAMS drivers can specify a list of STREAMS modules to be automatically pushed onto the stream at device open time. This (or a similar) feature is necessary to allow tty drivers to present an interface compatible with that which existed in previous system releases.

Drivers wishing to use this facility must specify the address of a NULL-terminated array of module names in the `st_modlist` field of the `streamtab` structure. The modules named in this array will be pushed at open time in the order of their appearance in the array. For example:

```
char *mypushlist[] = {
    "ldterm",
    "ttcompat",
    NULL
};

struct streamtab mydrivertab = {
    .
    .
    .
    char **mypushlist; /* st_modlist field */
};
```

These definitions arrange to push the standard modules implementing tty semantics when `mydriver` is opened.

STREAMS Portability

The set of internal interfaces and utility routines defined by the SunOS kernel differs considerably from that defined by the System V kernel. The STREAMS/kernel interface is well specified, however, and System V STREAMS modules and drivers that use only the interfaces it defines (see *Accessible Symbols and Functions* in the *Supplementary STREAMS Material* chapter of this manual) should be adaptable to the SunOS kernel without many problems. However, it's easy to use kernel facilities (data structures and routines) other than those defined in the STREAMS interface. Any such use is likely to be non-portable between System V and SunOS.

Similarly, STREAMS modules and drivers written for SunOS will only be portable to System V systems if their kernel interfaces are confined to the explicitly listed *Accessible Symbols and Functions*. If System V compatibility is not an issue, then STREAMS modules and drivers can use any of the driver-support routines listed in the *Kernel Support Routines* appendix of *Writing Device Drivers*.

Note that STREAMS drivers, as opposed to modules, will always require a certain degree of rewriting for use on System V machines, since the SunOS autoconfiguration interface differs significantly from that used in System V. See *The Bus-Resource Interface* section of *Writing Device Drivers* manual for the details of the Sun interface.

User Line Disciplines

NOTE *User-built line discipline code must be re-implemented within the STREAMS framework before it will work with release 4.1. This is not a simple matter of porting, for there is nothing within the STREAMS framework that really corresponds to line disciplines. In fact, STREAMS modules themselves can be seen as generalizations of line disciplines, in the sense that the character translation that line disciplines used to manage is now the job of STREAMS modules themselves. In a word, pre-STREAMS driver code is "ported", because the lower-level code that actually interacts with the device is typically preserved within the new STREAMS framework. Line discipline code, however, is best "rewritten" within the STREAMS framework. See Sun Consulting.*

There is another reason why pre-4.0 line discipline code will no longer work — such code probably accesses tty-specific internal structures, such as `clist` buffers. These structures no longer exist, having been replaced by STREAMS structures. For information on how to proceed with the conversion of a line discipline, contact the Sun consulting department.

Character device drivers that do not implement line disciplines can also be converted to STREAMS form, though in this case the conversion is entirely optional. This is because the SunOS STREAMS implementation preserves the *external* interfaces to the character devices and drivers (e.g. through the standard tty compatibility module, `ttcompat (4M)`, that implements most of the 4BSD tty interfaces under STREAMS). Thus, drivers that do not directly access underlying system data structures will continue to work without changes.

Drivers that have fancy read and write routines (routines that do anything more than just import parameters and perhaps start another routine) are probably not good candidates for conversion into STREAMS form, since STREAMS read/write modules should just set up data for the STREAMS queues.

A line-printer driver is an example of a character driver that could be written in terms of STREAMS, but doesn't need to be, and doesn't need to be converted to STREAMS if it already exists. After all, while a line-printer driver does transform a stream of characters (this transformation could certainly be built into a STREAMS module), its transformation is unlikely to be of interest to other programs. Thus, there's little to be gained by encapsulating it in a module. And, since line-printer drivers implement no line discipline, they will continue to work

with SunOS 4.1.

Supplementary STREAMS Material

A.1. Kernel Structures

This appendix summarizes previously described kernel structures commonly encountered in STREAMS module and driver development.

STREAMS kernel structures are contained in `<sys/stream.h>`.

Many of the fields in the structures described below are intended for the private use of the STREAMS framework code. Modules and drivers should not access such fields in any way, as their meaning and even existence is subject to change from release to release. **These fields are omitted from the descriptions below.**

Since these private fields can change from release to release, modules and drivers should not assume that the structures below have the same size from one release to the next.

Certain of the public fields listed below are read-only for modules and drivers; that is, only the STREAMS framework code is allowed to change them. These fields are marked with the legend "RO:" at the beginning of the comments describing them.

streamtab

As discussed in the *STREAMS Mechanism* section of the *STREAMS Module and Driver Programming* chapter, this structure defines a module or driver:

```

struct streamtab {
    struct qinit    *st_rdinit;    /* defines read QUEUE */
    struct qinit    *st_wrinit;    /* defines write QUEUE */
    struct qinit    *st_muxrinit; /* for multiplexing drivers only */
    struct qinit    *st_muxwinit; /* for multiplexing drivers only */
    char            **st_modlist; /* NULL-terminated list of
                                   modules to be pushed */
};

```

QUEUE Structures

Two sets of QUEUE structures form a module. The structures, discussed in the *STREAMS Mechanism* and *Message Queues and Service Procedures* sections of the *STREAMS Module and Driver Programming* chapter, are `queue_t`, `qinit`, `module_info` and, optionally, `module_stat`:

```

struct queue {
    struct qinit *q_qinfo; /* RO: procedures and limits for queue */
    struct msgb *q_first; /* RO: head of msg queue for this QUEUE */
    struct msgb *q_last; /* RO: tail of msg queue for this QUEUE */
    struct queue *q_next; /* RO: next QUEUE in stream */
    caddr_t q_ptr; /* to private data structure */
    ushort q_count; /* RO: count of characters on message queue */
    ushort q_flag; /* RO: QUEUE state */
    short q_minpsz /* min packet size accepted by this QUEUE */
    short q_maxpsz; /* max packet size accepted by this QUEUE */
    ushort q_hiwat; /* msg queue high water mark, for flow control */
    ushort q_lowat; /* msg queue low water mark, for flow control */
};
typedef struct queue queue_t;

```

When a `queue_t` pair is allocated, their contents are zero unless specifically initialized. The following fields are initialized:

- `q_qinfo` - from `streamtab.st_[rd/wr]init` (or `st_mux[rw]init`)
- `q_minpsz`, `q_maxpsz`, `q_hiwat`, `q_lowat` - from `module_info`
- `q_ptr` - optionally, by the driver/module open routine

```

struct qinit {
    int (*qi_putp) (); /* put procedure */
    int (*qi_srvp) (); /* service procedure */
    int (*qi_qopen) (); /* called on each open or a push */
    int (*qi_qclose) (); /* called on last close or a pop */
    int (*qi_qadmin) (); /* reserved for future use */
    struct module_info *qi_minfo; /* information structure */
    struct module_stat *qi_mstat; /* optional stats structure */
};

```

```

struct module_info {
    ushort mi_idnum; /* module ID number */
    char *mi_idname; /* module name */
    short mi_minpsz; /* min packet size accepted, for developer use */
    short mi_maxpsz; /* max packet size accepted, for developer use */
    short mi_hiwat; /* hi-water mark, for flow control */
    ushort mi_lowat; /* lo-water mark, for flow control */
};

```



```

struct module_stat {
    long    ms_pcnt;    /* count of calls to put proc */
    long    ms_scnt;    /* count of calls to service proc */
    long    ms_ocnt;    /* count of calls to open proc */
    long    ms_ccnt;    /* count of calls to close proc */
    long    ms_acnt;    /* count of calls to admin proc */
    char    *ms_xptr;    /* pointer to private statistics */
    short   ms_xsize;    /* length of private statistics buffer */
};

```

Note that in the event these counts are calculated by modules or drivers, the counts will be cumulative over all instantiations of modules with the same `fmodsw` entry and drivers with the same `cdevsw` entry.

A.2. Message Structures

As described in the *Messages* section of *STREAMS Module and Driver Programming*, a message is composed of a linked list of triples, consisting of two structures and a data buffer:

```

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 in buffer */
    struct datab *b_datap;    /* RO: data block */
};
typedef struct msgb mblk_t;

```

```

struct datab {
    unsigned char *db_base;    /* RO: first byte of buffer */
    unsigned char *db_lim;    /* RO: last byte+1 of buffer */
    unsigned char db_ref;    /* RO: cnt of msgs pointing to this block */
    unsigned char db_type;    /* message type */
};
typedef struct datab dblk_t;

```

`iocblk`

As described in the *Drivers* section of the *STREAMS Module and Driver Programming* chapter and in *Message Types*, below, this is contained in an `M_IOCTL` message block:

```

struct iocblk {
    int      ioc_cmd;      /* ioctl command type */
    ushort   ioc_uid;     /* effective uid of user */
    ushort   ioc_gid;     /* effective gid of user */
    uint     ioc_id;      /* ioctl id */
    uint     ioc_count;   /* count of bytes in data field */
    int      ioc_error;   /* error code */
    int      ioc_rval;    /* return value */
};

```

linkblk

As described in the *Multiplexing* section of *STREAMS Module and Driver Programming*, this is used in lower multiplexor drivers:

```

struct linkblk {
    queue_t *l_qtop;     /* lowest level write queue of upper stream */
    queue_t *l_qbot;     /* highest level write queue of lower stream */
    int      l_index;    /* system-unique index for lower stream. */
};

```

A.3. Message Types

Here the STREAMS message types are defined. The message types differ in their intended purposes, their treatment at the stream head, and in their message queueing priority (see the *Message Queues and Service Procedures* section of the *STREAMS Module and Driver Programming* chapter).

STREAMS does not prevent a module or driver from generating any message type and sending it in any direction on the stream. However, established processing and direction rules should be observed. Stream head processing according to message type is fixed, although certain parameters can be altered.

The message types are described below, classified according to their message queueing priority. Ordinary messages are described first, with high priority messages following. In certain cases, two message types may perform similar functions, differing in priority. Message construction is described in the *Messages* section of the *STREAMS Module and Driver Programming* chapter. The use of the word module will generally imply “module or driver.”

Ordinary Messages

These message types are subject to flow control. These are referred to as non-priority messages when received at user level.

M_DATA

Intended to contain ordinary data. Messages allocated by the `allocb()` routine (see *Message Types*, below) are type `M_DATA` by default. `M_DATA` messages are generally sent bidirectionally on a stream and their contents can be passed between a process and the stream head. In the `getmsg(2)` and `putmsg(2)` system calls, the contents of `M_DATA` message blocks are referred to as the data part. Messages composed of multiple message blocks will typically have `M_DATA` as the message type for all message blocks following the first.

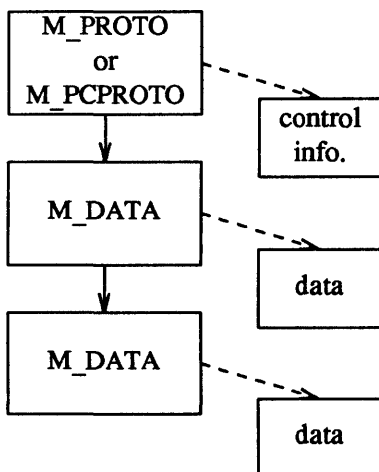
M_PROTO

Intended to contain internal control information and associated data. The message format is one M_PROTO message block followed by zero or more M_DATA message blocks as shown below: The semantics of the M_DATA and M_PROTO message block are determined by the STREAMS module that receives the message.

The M_PROTO message block will typically contain implementation dependent control information. M_PROTO messages are generally sent bidirectionally on a stream, and their contents can be passed between a process and the stream head. The contents of the first message block of an M_PROTO message is generally referred to as the control part, and the contents of any following M_DATA message blocks are referred to as the data part. In the getmsg(2) and putmsg(2) system calls, the control and data parts are passed separately. These calls refer to M_PROTO messages as non-priority messages.

Note that, although its use is not recommended, the format of M_PROTO and M_PCPROTO (generically PROTO) messages sent upstream to the stream head allows multiple PROTO blocks at the beginning of the message. getmsg() will compact the blocks into a single control part when passing them to the user process.

Figure A-1 M_PROTO and M_PCPROTO Message Structure

**M_IOCTL**

Generated by the stream head in response to an I_STR and certain other ioctl(2) system calls (see the streamio(4) man page.) When one of these ioctl()s is received from a user process, the stream head uses values from the process and supplied in the call to create an M_IOCTL message containing them, and sends the message downstream. M_IOCTL messages are intended to perform the general ioctl functions of character device drivers.

The user values are supplied in a structure of the following form, provided as an argument to the ioctl() call (see I_STR in the streamio(4) man

page).

```

struct strioctl {
    int ic_cmd;      /* downstream request */
    int ic_timeout; /* ACK/NAK timeout */
    int ic_len;     /* length of data arg */
    char *ic_dp;    /* ptr to data arg */
};

```

where `ic_cmd` is the request (or command) defined by a downstream module or driver, `ic_timeout` is the time the stream head will wait for acknowledgement to the `M_IOCTL` message before timing out, `ic_dp` is a pointer to an optional data argument. On input, `ic_len` contains the length of the data argument passed in and, on return from the call, it contains the length of the data, if any, being returned to the user.

The form of an `M_IOCTL` message is one `M_IOCTL` message block linked to zero or more `M_DATA` message blocks. STREAMS constructs an `M_IOCTL` message block by placing an `iocblk` structure in its data buffer:

```

struct iocblk {
    int ioc_cmd;      /* ioctl command type */
    ushort ioc_uid;   /* effective user id number */
    ushort ioc_gid;   /* effective group id number */
    uint ioc_id;      /* ioctl identifier */
    uint ioc_count;   /* byte count for ioctl data */
    int ioc_error;    /* error code */
    int ioc_rval;     /* return value */
};

```

The `iocblk` structure is defined in `<sys/stream.h>`. `ioc_cmd` corresponds to `ic_cmd`. `ioc_uid` and `ioc_gid` are the effective user and group IDs for the user sending the `ioctl()`, and can be tested to determine if the user issuing the `ioctl()` call is authorized to do so. `ioc_count` is the number of data bytes, if any, contained in the message and corresponds to `ic_len`.

`ioc_id` is an identifier generated internally, and is used to match each `M_IOCTL` message sent downstream with a response that must be sent upstream to the stream head. The response is contained in an `M_IOCACK` (positive acknowledgement) or an `M_IOCNAK` (negative acknowledgement) message. Both these message types have the same format as an `M_IOCTL` message and contain an `iocblk` structure in the first block with optional data blocks following. If one of these messages reaches the stream head with an identifier that does not match that of the currently-outstanding `M_IOCTL` message, the response message is discarded. A common means of assuring that the correct identifier is returned, is for the replying module to convert the `M_IOCTL` message type into the appropriate response type and set `ioc_count` to 0, if no data is returned. Then, the `qreply()`

utility (see *Utilities*, below) is used to send the response to the stream head.

`ioc_error` holds any return error condition set by a downstream module. If this value is non-zero, it is returned to the user in `errno`. Note that both an `M_IOCNAK` and an `M_IOCACK` may return an error. `ioc_rval` holds any `M_IOCACK` return value set by a responding module.

If a user supplies data to be sent downstream, the stream head copies the data, pointed to by `ic_dp` in the `striocctl` structure, into `M_DATA` message blocks and links the blocks to the initial `M_IOCTL` message block. `ioc_count` is copied from `ic_len`. If there is no data, `ioc_count` is zero.

If a module wants to send data to a user process as part of its response, it must construct an `M_IOCACK` message that contains the data. The first message block of this message contains the `iocblk` data structure, with any data stored in one or more `M_DATA` message blocks linked to the first message block. The module must set `ioc_count` to the number of data bytes sent. On completion of the call, this number is passed to the user in `ic_len`. Data associated with an `M_IOCNAK` message is not returned to the user process, and is discarded by the stream head.

The first module or driver that understands the request contained in the `M_IOCTL` acts on it, and generally returns an `M_IOCACK` message. Intermediate modules that do not recognize a particular request must pass it on. If a driver does not recognize the request, or the receiving module can not acknowledge it, an `M_IOCNAK` message must be returned.

The stream head waits for the response message and returns any information contained in an `M_IOCACK` to the user. The stream head will “time out” if no response is received in `ic_timeout` interval.

M_CTL

Generated by modules that wish to send information to a particular module or type of module. `M_CTL` messages are typically used for inter-module communication, as when adjacent STREAMS protocol modules negotiate the terms of their interface. An `M_CTL` message cannot be generated by a user-level process and is always discarded if passed to the stream head.

M_BREAK

Sent to a driver to request that `BREAK` be transmitted on whatever media the driver is controlling.

The message format is not defined by STREAMS and its use is developer dependent. This message may be considered a special case of an `M_CTL` message. An `M_BREAK` message cannot be generated by a user-level process and is always discarded if passed to the stream head.

M_DELAY

Sent to a media driver to request a real-time delay on output. The data buffer associated with this message type is expected to contain an integer to indicate the number of machine ticks of delay desired. `M_DELAY` messages are typically used to prevent transmitted data from exceeding the buffering

capacity of slower terminals.

The message format is not defined by STREAMS and its use is developer dependent. Not all media drivers may understand this message. This message may be considered a special case of an `M_CTL` message. An `M_DELAY` message cannot be generated by a user-level process and is always discarded if passed to the stream head.

M_PASSFP

This is used by STREAMS to pass a file pointer from the stream head at one end of a stream pipe to the stream head at the other end of the same stream pipe. (A stream pipe is a stream that is terminated at both ends by a stream head; one end of the stream can always find the other by following the `q_next` pointers in the stream. The means by which such a structure is created is not described in this document.)

The message is generated as a result of an `I_SENDFD ioctl()` (see the `streamio(4)` man page) issued by a process to the sending stream head. STREAMS places the `M_PASSFP` message directly on the destination stream head's read queue to be retrieved by an `I_RECVFD ioctl()` (see the `streamio(4)` man page). The message is placed without passing it through the stream (i.e., it is not seen by any modules or drivers in the stream). This message type should never be present on any queue except the read queue of a stream head. Consequently, modules and drivers do not need to recognize this message type, and it can be ignored by module and driver developers.

M_SETOPTS

Alters some characteristics of the stream head. It is generated by any downstream module, and is interpreted by the stream head. The data buffer of the message has the following structure:

```
struct stroptions {
    short so_flags;      /* options to set */
    short so_readopt;   /* read option */
    ushort so_wroff;    /* write offset */
    short so_minpsz;    /* minimum read packet size */
    short so_maxpsz;    /* maximum read packet size */
    ushort so_hiwat;    /* read queue high-water mark */
    ushort so_lowat;    /* read queue low-water mark */
};
```

where `so_flags` specifies which options are to be altered, and can be any combination of the following:

SO_ALL

Update all options according to the values specified in the remaining fields of the `stroptions` structure.

SO_READOPT

Set the read mode (see the `read(2)` man page) to `RNORM` (byte stream), `RMSGD` (message discard), or `RMSGN` (message non-discard) as

specified by the value of `so_readopt`.

SO_WROFF

Direct the stream head to insert an offset specified by `so_wroff` into the first message block of all `M_DATA` messages created as a result of a `write()` system call. The same offset is inserted into the first `M_DATA` message block, if any, of all messages created by a `putmsg()` system call. The default offset is zero.

The offset must be less than the maximum message buffer size (system dependent). Under certain circumstances, a write offset may not be inserted. A module or driver must test that `b_rptr` in the `mblk_t` structure is greater than `db_base` in the `dbl_t` structure to determine that an offset has been inserted in the first message block.

SO_MINPSZ

Change the minimum packet size value associated with the stream head read queue to `so_minpsz` (see `q_minpsz` in the `queue_t` structure, in the *Kernel Structures* section, above) This value is advisory for the module immediately below the stream head. It is intended to limit the size of `M_DATA` messages that the module should put to the stream head. There is no intended minimum size for other message types. The default value in the stream head is 0.

SO_MAXPSZ

Change the maximum packet size value associated with the stream head read queue to `so_maxpsz` (see `q_maxpsz` in the `queue_t` structure, in the *Kernel Structures* section, above). This value is advisory for the module immediately below the stream head. It is intended to limit the size of `M_DATA` messages that the module should put to the stream head. There is no intended maximum size for other message types. The default value in the stream head is `INFPSZ`, the maximum STREAMS allows.

SO_HIWAT

Change the flow control high water mark on the stream head read queue to the value specified in `so_hiwat`.

SO_LOWAT

Change the flow control low water mark (see `q_minpsz` in the `queue_t` structure, in the *Kernel Structures* section, above) on the stream head read queue to the value specified in `so_lowat`.

M_SIG

Sent upstream by modules or drivers to post a signal to a process. When the message reaches the stream head, the first data byte of the message is transformed into a signal, as defined in `<sys/signal.h>`, to the process(es) according to the following.

If the signal is not `SIGPOLL` and the stream containing the sending module or driver is a controlling TTY, the signal is sent to the associated process group. If the stream does not have a process group, then `M_SIG` is discarded and the signal is ignored.

If the signal is SIGPOLL, it will be sent only to those processes that have explicitly registered to receive the signal (see `I_SETSIG` in the `streamio(4)` man page).

High Priority Messages

High priority messages are not subject to flow control.

M_PCPROTO

This message type has the same format and characteristics as the `M_PROTO` message type, except for priority and the following additional attributes.

When an `M_PCPROTO` message is placed on a queue, its service procedure is always enabled. The stream head will allow only one `M_PCPROTO` message to be placed in its read queue at a time. If an `M_PCPROTO` message is already in the queue when another arrives, the second message is silently discarded and its message blocks freed.

This message type is intended to allow data and control information to be sent outside the normal flow control constraints.

The `getmsg(2)` and `putmsg(2)` system calls refer to `M_PCPROTO` messages as priority messages.

M_ERROR

This message type is sent upstream by modules or drivers to report some downstream error condition. When the message reaches the stream head, the stream is marked so that all subsequent system calls issued to the stream, excluding `close(2)` and `poll(2)`, will fail with `errno` set to the first data byte of the message. `POLLERR` is set if the stream is being `poll()`ed (see the `poll(2)` man page. All processes sleeping on a system call to the stream are awakened. An `M_FLUSH` message with a `FLUSHRW` argument is sent downstream.

M_HANGUP

This message type is sent upstream by a driver to report that it can no longer send data upstream. For example, this might be due to an error, or to a remote line connection being dropped. When the message reaches the stream head, the stream is marked so that all subsequent `write(2)` and `putmsg(2)` system calls issued to the stream will fail and return an `ENXIO` error. Those `ioctl()`s that cause messages to be sent downstream are also rejected. `POLLHUP` is set if the stream is being `poll()`ed (see the `poll(2)` man page).

However, subsequent `read(2)` or `getmsg(2)` calls to the stream will not generate an error. These calls will return any messages (according to their function) that were on, or in transit to, the stream head read queue before the `M_HANGUP` message was received. When all such messages have been read, `read()` will return 0, and `getmsg()` will set each of its two length fields to 0.

This message also causes a `SIGHUP` signal to be sent to the process group, if the device is a controlling TTY (see `M_SIG`).

M_IOCACK

This message type signals the positive acknowledgement of a previous `M_IOCTL` message. The message may contain information sent by the receiving module or driver. The stream head returns the information to the user if there is a corresponding outstanding `M_IOCTL` request. The format and use of this message type is described further under `M_IOCTL`.

M_IOCNAK

This message type signals the negative acknowledgement (failure) of a previous `M_IOCTL` message. When the stream head receives an `M_IOCNAK`, the outstanding `ioctl()` request, if any, will fail. The format and usage of this message type is described further under `M_IOCTL`.

M_FLUSH

This message type requests all modules and drivers that receive it to flush their message queues (discard all messages in those queues) as indicated in the message. An `M_FLUSH` can originate at the stream head, or in any module or driver. The first byte of the message contains flags that specify one of the following actions:

FLUSHR:

Flush the read queue of the module.

FLUSHW:

Flush the write queue of the module.

FLUSHRW:

Flush both the read and the write queue of the module.

Each module passes this message to its neighbor after flushing its appropriate queue(s), until the message reaches one of the ends of the stream.

Drivers are expected to include the following processing for `M_FLUSH` messages. When an `M_FLUSH` message is sent downstream through the write queues in a stream, the driver at the stream end discards it if the message action indicates that the read queues in the stream are not to be flushed (only `FLUSHW` set). If the message indicates that the read queues are to be flushed, the driver sets the `M_FLUSH` message flag to `FLUSHR`, and sends the message up the stream's read queues. When a flush message is sent up a stream's read side, the stream head checks to see if the write side of the stream is to be flushed. If only `FLUSHR` is set, the stream head discards the message. However, if the write side of the stream is to be flushed, the stream head sets the `M_FLUSH` flag to `FLUSHW` and sends the message down the stream's write side. *All modules that enqueue messages must recognize and process this message type.*

M_PCSIG

This message type has the same format and characteristics as the `M_SIG` message type except for priority.

M_START and M_STOP

These messages request devices to start or stop their output. They are intended to produce momentary pauses in a device's output, not to turn

devices on or off.

The message format is not defined by STREAMS and its use is developer dependent. These messages may be considered special cases of an `M_CTL` message. These messages cannot be generated by a user-level process and each is always discarded if passed to the stream head.

A.4. Utilities

The utilities contained in this appendix represent an interface that will be maintained in subsequent versions of SunOS. Other than these utilities (see also the *Accessible Symbols and Functions* section, below) functions contained in the STREAMS kernel code may change in future releases.

This appendix specifies the set of utilities that STREAMS provides to assist development of modules and drivers. There are many utility routines and macros.

The general purpose of the utilities is to perform functions that are commonly used in modules and drivers. A utility must always be used when operating on a message queue and when accessing the buffer pool.

The utilities are contained in either the system source file `os/str_buf.c` or, if they are macros, in `<sys/stream.h>`.

All structure definitions are contained in the *Kernel Structures* section, above, unless otherwise indicated. All routine references are found in this section unless otherwise indicated. The following definitions are used.

Blocked

A queue that can not be enabled due to flow control (see the *Flow Control* section in the *Introduction to STREAMS* chapter of the *System Services Overview*).

Enable

To schedule a queue.

Free

De-allocate STREAMS storage.

Message block (bp)

A triplet consisting of an `mblk_t` structure, a `dbl_t` structure, and a data buffer. It is referenced by its `mblk_t` structure (see the *Messages* section of the *STREAMS Module and Driver Programming* chapter).

Message (mp)

One or more linked message blocks. A message is referenced by its first message block.

Message queue

Zero or more linked messages associated with a queue (`queue_t` structure).

Queue (q)

A `queue_t` structure. This is generally the same as `QUEUE` in the rest of this document (e.g., see the definitions for `enable` and `schedule`). When it appears with “message” in certain utility description lines, it means “message queue.”

Schedule

Place a queue on the internal linked list of queues which will subsequently have their service procedure called by the STREAMS scheduler.

The word *module* will generally mean “module and/or driver.” The phrase “next/following module” will generally refer to a module, driver, or stream head. Message queuing priority (see the *Message Queues and Service Procedures* section of the *STREAMS Module and Driver Programming* chapters and the *Message Types* section, above) can be ordinary or high priority (to avoid “priority priority”).

Buffer Allocation Priority

STREAMS buffers are normally allocated with `allocb()`, described above. An associated set of allocation priorities has been established, which are also used in other utility routines:

BPRI_LO

Low priority. At this priority, `allocb()` may fail even though the requested buffer size is available. This priority is used by the stream head write routine to hold data associated with user calls.

BPRI_MED

Medium priority. This priority is typically used for normal data and control block allocation. As above, `allocb()` may fail at this priority even though a buffer of the requested size is available. However, for a given block size, a `BPRI_LO allocb()` call will fail before a `BPRI_MED allocb()` call.

BPRI_HI

High priority. This priority is typically used only for critical control message allocations. Calls to `allocb()` will succeed if a buffer of the appropriate size is available. Developers should exercise restraint in use of `BPRI_HI` allocation requests.

The values `BPRI_LO`, `BPRI_MED`, and `BPRI_HI` are defined in `<sys/stream.h>`.

STREAMS does not guarantee successful buffer allocation—any set of resources can be exhausted under the right or wrong conditions. The `bufcall()` function will help modules recover from buffer allocation failures, but it does not guarantee that the resources will ever be available. Developers should be aware of this when implementing modules.

adjmsg() — Trim Bytes in a Message

```
int adjmsg(mp, len)
    mblk_t *mp;
    int len;
```

`adjmsg()` trims bytes from either the head or tail of the message specified by *mp*. If *len* is greater than zero, it removes *len* bytes from the beginning of *mp*. If *len* is less than zero, it removes *len* bytes from the end of *mp*. If *len* is zero, `adjmsg()` does nothing. `adjmsg()` only trims bytes across message blocks of the same type. It will fail if *mp* points to a message containing fewer than *len* bytes of similar type at the message position indicated. `adjmsg()` returns 1 on success, and 0 on failure.

allocb () — Allocate a Message Block

```
mblk_t *allocb(size, pri)
      int size, pri;
```

`allocb ()` returns a pointer to a message block of type `M_DATA`, in which the data buffer contains at least `size` bytes. `pri` indicates the priority of the allocation request, and can have the values `BPRI_LO`, `BPRI_MED`, or `BPRI_HI` (see *Buffer Allocation Priority*, above). If a block can not be allocated as requested, `allocb ()` returns a `NULL` pointer.

backq () — Get Pointer to Queue Behind a Given Queue

```
queue_t *backq(q)
      queue_t *q;
```

`backq ()` returns a pointer to the queue behind a given queue. That is, it returns a pointer to the queue whose `q_next` (see `queue_t` structure) pointer is `q`. If no such queue exists (as when `q` is at a stream end), `backq ()` returns `NULL`.

bufcall () — Recover from Failure of allocb ()

```
int bufcall(size, pri, func, arg)
      uint size;
      int pri;
      int (*func) ();
      long arg;
```

`bufcall ()` is provided to assist in the event of a block allocation failure. If `allocb ()` returns `NULL`, indicating a message block is not currently available, `bufcall ()` may be invoked.

`bufcall ()` arranges for `(*func)(arg)` to be called when a buffer of `size` bytes at `pri` priority (see *Buffer Allocation Priority*, below) is available. When `func` is called, it has no user context. It cannot reference the *user* structure and must return without sleeping. `bufcall ()` does not guarantee that the desired buffer will be available when `func` is called since interrupt processing may acquire it.

`bufcall ()` returns a non-zero value on success, indicating that the request has been successfully recorded, or 0 on failure. On a failure return, `func` will never be called. A failure indicates a (temporary) inability to allocate required internal data structures or a `size` value that is so large as to be never satisfiable.

The success return value should be saved for possible future use as described below. When a module or driver is closed, its close routine should cancel all pending `bufcall ()` requests by calling `unbufcall ()` with the saved return value from `bufcall ()` as argument.

canput () — Test for Room in a Queue

```
int canput(q)
      queue_t *q;
```

`canput ()` determines if there is room left in a message queue. If `q` does not have a service procedure, `canput ()` will search further in the same direction in

the stream until it finds a queue containing a service procedure (this is the first queue on which the passed message can actually be enqueued). If such a queue cannot be found, the search terminates on the queue at the end of the stream. `canput()` tests the queue found by the search. If the message queue in this queue is not full (see the *Flow Control* section in the *Introduction to STREAMS* chapter of the *System Services Overview*) `canput()` returns 1. This return indicates that a message can be put to queue *q*. If the message queue is full, `canput()` returns 0. In this case, the caller is generally referred to as blocked.

`copyb()` — Copy a Message Block

```
mblk_t *copyb(bp)
mblk_t *bp;
```

`copyb()` copies the contents of the message block pointed at by *bp* into a newly-allocated message block of at least the same size. `copyb()` allocates a new block by calling `allocb()` with *pri* set to `BPRI_MED` (see *Buffer Allocation Priority*, above). All data between the *b_rptr* and *b_wptr* pointers of a message block are copied to the new block, and these pointers in the new block are given the same offset values they had in the original message block. On successful completion, `copyb()` returns a pointer to the new message block containing the copied data. Otherwise, it returns a NULL pointer.

`copymsg()` — Copy a Message

```
mblk_t *copymsg(mp)
mblk_t *mp;
```

`copymsg()` uses `copyb()` to copy the message blocks contained in the message pointed at by *mp* to newly-allocated message blocks, and links the new message blocks to form the new message. On successful completion, `copymsg()` returns a pointer to the new message. Otherwise, it returns a NULL pointer.

`datamsq()` — Test Whether Message is a Data Message

```
#define datamsq(mp) ...
```

The `datamsq` macro returns TRUE if *mp* (declared as `mblk_t *mp`) points to a data type message. Types `M_DATA`, `M_PROTO`, `M_PCPROTO`, and `M_DELAY` count as data (see the *Message Types* section, above). If *mp* points to any other message type, `datamsq` returns FALSE.

`dupb()` — Duplicate a Message Block Descriptor

```
mblk_t *dupb(bp)
mblk_t *bp;
```

`dupb()` duplicates the message block descriptor (`mblk_t` structure) pointed at by *bp* by copying it into a newly allocated message block descriptor. A message block is formed with the new message block descriptor pointing to the same data block as the original descriptor. The reference count in the data block descriptor (`dbl_t` structure) is incremented. `dupb()` does not copy the data buffer, only the message block descriptor.

On successful completion, `dupb()` returns a pointer to the new message block. If `dupb()` cannot allocate a new message block descriptor, it returns `NULL`.

This routine allows message blocks that exist on different queues to reference the same data block. In general, if the contents of a message block with a reference count greater than 1 are to be modified, `copyb()` should be used to create a new message block and only the new message block should be modified. This insures that other references to the original message block are not invalidated by unwanted changes.

dupmsg() — Duplicate a Message

```
mbblk_t *dupmsg(mp)
    mblk_t *mp;
```

`dupmsg()` calls `dupb()` to duplicate the message pointed at by `mp`, by copying all individual message block descriptors, and then linking the new message blocks to form the new message. `dupmsg()` does *not* copy data buffers, only message block descriptors. On successful completion, `dupmsg()` returns a pointer to the new message. Otherwise, it returns `NULL`.

enableok() — Re-allow Queue to be Scheduled

```
#define enableok(q) ...
```

The `enableok()` macro cancels the effect of an earlier `noenable()` on the same queue `q` (declared as `queue_t *q`). It allows a queue to be scheduled for service that had previously been excluded from queue service by a call to `noenable()`.

esballoc() — Extended STREAMS Buffer Allocation

```
mbblk_t *
esballoc(base, size, pri, fr_rtn)
    unsigned char *base;
    int size, pri;
    frtn_t *fr_rtn;
```

`esballoc()` allocates message and data blocks that point directly to a client-supplied buffer. Message and data blocks allocated this way are indistinguishable from the normal data blocks. The client-supplied buffers are processed as if they were normal STREAMS buffers. `fr_rtn` points to a `frtn_t` structure (defined in `<sys/stream.h>`) containing the address of a function and an argument to that function. This function is called with that argument when the message block that `esballoc` returns is freed.

esbcall() — bufcall for extended buffers

```
int
esbcall(pri, func, arg)
    int pri;
    int (*func)();
    long arg;
```

`esbbcall()` is an analog of `bufcall()` that can be called when `esbal-loc()` fails. See the description of `bufcall()` for more information.

flushq() — Flush a Queue

```
flushq(q, flag)
      queue_t *q;
      int flag;
```

`flushq()` removes messages from the message queue in queue *q* and frees them, using `freemsg()`. If *flag* is set to `FLUSHDATA`, then `flushq()` discards all `M_DATA`, `M_PROTO`, `M_PCPROTO`, and `M_DELAY` messages (see `datamsg`), but leaves all other messages on the queue. If *flag* is set to `FLUSHALL`, all messages are removed from the message queue and freed. `FLUSHALL` and `FLUSHDATA` are defined in `<sys/stream.h>`.

If a queue behind *q* is blocked, `flushq()` may enable the blocked queue, as described in `putq()`.

freeb() — Free a Message Block

```
freeb(bp)
      mblk_t *bp;
```

`freeb()` will free (de-allocate) the message block descriptor pointed at by *bp*, and free the corresponding data block if the reference count (see `dupb()`) in the data block descriptor (`dblk_t` structure) is equal to 1. If the reference count is greater than 1, `freeb()` will not free the data block, but will decrement the reference count.

freemsg() — Free All Message Blocks in a Message

```
freemsg(mp)
      mblk_t *mp;
```

`freemsg()` uses `freeb()` to free all message blocks and their corresponding data blocks for the message pointed at by *mp*.

getq() — Get a Message from a Queue

```
mblk_t *getq(q)
      queue_t *q;
```

`getq()` gets the next available message from the queue pointed at by *q*. `getq()` returns a pointer to the message and removes that message from the queue. If no message is queued, `getq()` returns `NULL`.

`getq()`, and certain other utility routines, affect flow control in the stream as follows: If `getq()` returns `NULL`, the queue is internally marked so that the next time a message is placed on it, it will be scheduled for service (enabled, see `genable()`). Also, if the data in the enqueued messages drops below the low-water mark, `q_lowat`, and a queue behind the current queue had previously attempted to place a message in the queue and failed (i.e., was blocked, see `canput()`), then the queue behind the current queue is scheduled for service (see the *Flow*

Control section in the *Introduction to STREAMS* chapter).

insq() — Put a Message at a Specific Place in a Queue

```
insq(q, emp, nmp)
    queue_t *q;
    mblk_t *emp, *nmp;
```

`insq()` places the message pointed at by `nmp` in the message queue contained in the queue pointed at by `q` immediately before the already-enqueued message pointed at by `emp`. If `emp` is `NULL`, the message is placed at the end of the queue. If `emp` is non-`NULL`, it must point to a message that exists on the queue `q`, or a system panic could result.

Note that the message is placed where indicated, without consideration of message queueing priority. The queue will be scheduled in accordance with the rules described in `putq()` for ordinary priority messages.

linkb() — Concatenate Two Messages into One

```
linkb(mp1, mp2)
    mblk_t *mp1;
    blk_t *mp2;
```

`linkb()` puts the message pointed at by `mp2` at the tail of the message pointed at by `mp1`.

msgdsize() — Get Number of Data Bytes in a Message

```
int msgdsize(mp)
    mblk_t *mp;
```

`msgdsize()` returns the number of bytes of data in the message pointed at by `mp`. Only bytes included in data blocks of type `M_DATA` are included in the total.

noenable() — Prevent a Queue from Being Scheduled

```
#define noenable(q) ...
```

The `noenable()` macro prevents the queue `q` (declared as `queue_t *q`) from being scheduled for service by `putq()` or `putbq()` when these routines enqueue an ordinary priority message, or by `insq()` when it enqueues any message. `noenable()` does not prevent the scheduling of queues when a high priority message is enqueued, unless it is enqueued by `insq()`.

OTHERQ() — Get Pointer to the Mate Queue

```
#define OTHERQ(q) ...
```

The `OTHERQ()` macro returns a pointer to the mate queue of `q` (declared as `queue_t *q`). If `q` is the read queue for the module, it returns a pointer to the module's write queue. If `q` is the write queue for the module, it returns a pointer to the read queue.

**pullupmsg () —
Concatenate Bytes in a
Message**

```
int pullupmsg(mp, len)
    mblk_t *mp;
    int len;
```

`pullupmsg()` concatenates and aligns the first *len* data bytes of the passed message into a single, contiguous message block. Proper alignment is hardware-dependent. To perform its function, `pullupmsg()` allocates a new message block by calling `allocb()` with *pri* set to `BPRI_MED` (see *Buffer Allocation Priority*, above). `pullupmsg()` only concatenates across message blocks of similar type. It will fail if *mp* points to a message of less than *len* bytes of similar type. A *len* value of -1 requests a pull-up of all the like-type blocks in the beginning of the message pointed at by *mp*.

At completion of concatenation, `pullupmsg()` replaces *mp* with a pointer to the new message block, so that *mp* still points to the same message block at the end of the operation. However, the contents of the message block may have been altered. On success, `pullupmsg()` returns 1. On failure, it returns 0.

**putbq () — Return a
Message to the Beginning of a
Queue**

```
putbq(q, bp)
    queue_t *q;
    mblk_t *bp;
```

`putbq()` puts the message pointed at by *bp* at the beginning of the queue pointed at by *q*, in a position in accordance with the message's type. High priority messages are placed at the head of the queue, and ordinary messages are placed after all high priority messages, but before all other ordinary messages. The queue will be scheduled in accordance with the same rules described in `putq()`. This utility is typically used to replace a message on a queue from which it was just removed.

**putctl () — Put a Control
Message**

```
int putctl(q, type)
    queue_t *q;
    int type;
```

`putctl()` creates a control (not data, see `datamsg`, above) message of type *type*, and calls the *put* procedure in the queue pointed at by *q*, with a pointer to the created message as an argument. `putctl()` allocates new blocks by calling `allocb()` with *pri* set to `BPRI_HI` (see the *Buffer Allocation Priority* section, above). On successful completion, `putctl()` returns 1. It returns 0 if it cannot allocate a message block, or if *type* `M_DATA`, `M_PROTO`, or `M_PCPROTO` was specified.

**putctl1 () — Put One-byte
Parameter Control Message**

```
int putctl1(q, type, p)
    queue_t *q;
    int type;
    int p;
```

`putcctl1()` creates a control (not data, see `datamsg`, above) message of type *type* with a one-byte parameter *p*, and calls the *put* procedure in the queue pointed at by *q*, with a pointer to the created message as an argument. `putcctl1()` allocates new blocks by calling `allocb()` with *pri* set to `BPRI_HI` (see the *Buffer Allocation Priority* section, above). On successful completion, `putcctl1()` returns 1. It returns 0 if it cannot allocate a message block, or if *type* `M_DATA`, `M_PROTO`, or `M_PCPROTO` was specified.

`putnext()` — Put a Message to the Next Queue

```
#define putnext(q, mp) ...
```

The `putnext()` macro calls the *put* procedure of the next queue in a STREAM, and passes it a message pointer as an argument. The parameters must be declared as `queue_t *q` and `mblk_t *mp`. *q* is the calling queue (not the next queue) and *mp* is the message to be passed. `putnext()` is the typical means of passing messages to the next queue in a STREAM.

`putq()` — Put a Message on a Queue

```
putq(q, bp)
    queue_t *q;
    mblk_t *bp;
```

`putq()` puts the message pointed at by *bp* on the message queue contained in the queue pointed at by *q* and enables that queue. `putq()` queues messages appropriately by type (i.e., message queueing priority, see the *Message Queues and Service Procedures* section of the *STREAMS Module and Driver Programming* chapter).

`putq()` will always enable the queue when a high priority message is queued. `putq()` will enable the queue when an ordinary message is queued if the following condition is set, and enabling is not inhibited by `noenable()`. The condition is set if the module has just been pushed (see `I_PUSH` in `streamio(4)`), or if no message was queued on the last `getq()` call and no message has been queued since.

`putq()` is intended to be used from the *put* procedure in the same queue in which the message will be queued. A module should not call `putq()` directly to pass messages to a neighboring module. `putq()` may be used as the `qi_putp` *put* procedure value in either or both of a module's `qinit` structures. This effectively bypasses any *put* procedure processing and uses only the module's service procedure(s).

`qenable()` — Enable a Queue

```
qenable(q)
    queue_t *q;
```

`qenable()` places the queue pointed at by *q* on the linked list of queues that are ready to be called by the STREAMS scheduler (see the definition for "Schedule" above, and the *Put and Service Procedures* section in the *Introduction to STREAMS* chapter).

qreply () — Send Reverse-Direction Message

```
qreply(q, bp)
queue_t *q;
mblk_t *bp;
```

`qreply ()` sends the message pointed at by `bp` up (or down) the stream in the reverse direction from the queue pointed at by `q`. This is done by locating the partner of `q` (see `OTHERQ ()`, below), and then calling the `put` procedure of that queue's neighbor (as in `putnext ()`). `qreply ()` is typically used to send back a response (`M_IOCACK` or `M_IOCNAK` message) to an `M_IOCTL` message (see *Message Types*, above).

qsize () — Find the Number of Messages on a Queue

```
int qsize(q)
queue_t *q;
```

`qsize ()` returns the number of messages present in queue `q`. If there are no messages on the queue, `qsize ()` returns 0.

RD () — Get Pointer to the Read Queue

```
#define RD(q) ...
```

The `RD ()` macro accepts a write queue pointer, `q` (declared as `queue_t *q`), as an argument and returns a pointer to the read queue for the same module.

rmvb () — Remove a Message Block from a Message

```
mblk_t *rmvb(mp, bp)
mblk_t *mp;
mblk_t *bp;
```

`rmvb ()` removes the message block pointed at by `bp` from the message pointed at by `mp`, and then restores the linkage of the message blocks remaining in the message. `rmvb ()` does not free the removed message block. `rmvb ()` returns a pointer to the head of the resulting message. If `bp` is not contained in `mp`, `rmvb ()` returns a -1. If there are no message blocks in the resulting message, `rmvb ()` returns a NULL pointer.

rmvq () — Remove a Message from a Queue

```
rmvq(q, mp)
queue_t *q;
mblk_t *mp;
```

`rmvq ()` removes the message pointed at by `mp` from the message queue in the queue pointed at by `q`, and then restores the linkage of the messages remaining on the queue. If `mp` does not point to a message that is present on the queue `q`, a system panic could result.

splstr() — Set Processor Level

```
int splstr()
```

`splstr()` increases the system processor level to block interrupts at a level appropriate for STREAMS modules when those modules are executing critical portions of their code. `splstr()` returns the processor level at the time of its invocation. Module developers are expected to use the standard kernel function `splx(s)`, where *s* is the integer value returned by `splstr()`, to restore the processor level to its previous value after the critical portions of code are passed.

testb() — Check for an Available Buffer

```
int testb(size, pri)
    register size;
    uint pri;
```

`testb()` checks for the availability of a message buffer of size *size* at priority *pri* (see *Buffer Allocation Priority*, below) without actually retrieving the buffer. `testb()` returns 1 if the buffer is available, and 0 if no buffer is available. A successful return value from `testb()` does not guarantee that a subsequent `allocb()` call will succeed (e.g., in the case of an interrupt routine taking buffers).

unbufcall() — Cancel an outstanding `bufcall` request

```
void unbufcall(id)
    int id;
```

`unbufcall()` cancels the pending `bufcall()` request denoted by *id*, where *id* is the value returned previously by `bufcall()`. This routine is typically called in a module's close routine to clean up before returning. (Clean up is often necessary to prevent executing the function supplied to `bufcall()` with a no-longer-valid argument.)

unlinkb() — Remove Message Block from Message Head

```
mblk_t *unlinkb(mp)
    mblk_t *mp;
```

`unlinkb()` removes the first message block pointed at by *mp* and returns a pointer to the head of the resulting message. `unlinkb()` returns a NULL pointer if there are no more message blocks in the message.

WR() — Get Pointer to the Write Queue

```
#define WR(q) ...
```

The `WR` macro accepts a read queue pointer, *q* (declared as `queue_t *q`), as an argument and returns a pointer to the write queue for the same module.

A.5. Design Guidelines

This appendix summarizes STREAMS module and driver design guidelines and rules presented in previous chapters. Additional rules that developers must observe are included. Where appropriate, the section of this document containing detailed information is named.

Unless otherwise noted, “module” implies “modules and drivers”.

General Rules

The following are general rules that developers should follow when writing modules.

1. Modules cannot access information in the `user` structure associated with a process. Modules are not associated with any process, and therefore have no concept of process or user context.

The capability to pass `user` structure information upstream using messages has been provided where required. This can be done in `M_IOCTL` handling (see the *Drivers* section of the *STREAMS Module and Driver Programming* chapter and also *Message Types*, above). A module can send error codes upstream in an `M_IOCACK` or `M_IOCNAK` message, where they will be placed in `u_error` by the stream head. Return values may also be sent upstream in a `M_IOCACK` message, and will be placed in `u_rval1`. Information can also be passed to the `user` structure via a `M_ERROR` message (see the *Complete Driver* section of the *STREAMS Module and Driver Programming* chapter and also *Message Types*, above). The stream head will recognize this message type and inform the next system call that an error has occurred downstream by setting `u_error`. Note that in both instances, the downstream module cannot access the `user` structure, but it informs the stream head to do so.

2. In general, modules should not require the data in an `M_DATA` message to follow a particular format, such as a specific alignment. This makes it easier to arbitrarily push modules on top of each other in a sensible fashion. Not following this rule may limit module re-usability (the ability to use the module in multiple applications).
3. Every module must process an `M_FLUSH` message according to the value of the argument passed in the message. (See the *Message Queues and Service Procedures* and *Drivers* sections of *STREAMS Module and Driver Programming*, and also *Message Types*, above).
4. A module should not change the contents of a data block whose reference count is greater than 1 (see `dupmsg()` in the *Utilities* section, above) because other modules that have references to the block may not want the data changed. To avoid problems, it is recommended that the module copy the data to a new block and then change the new one.
5. Modules should only manipulate message queues and manage buffers with the routines provided for those purposes, (see the *Utilities* section, above).
6. Filter modules pushed between a service user and a service provider (see the *Service Interface* section of the *STREAMS Module and Driver Programming* chapter) may not alter the contents of the `M_PROTO` or `M_PCPROTO` block in messages. The contents of the data blocks may be manipulated, but the

message boundaries must be preserved.

System Calls

These rules pertain to modules and drivers as noted.

1. *open* and *close* routines may sleep, but the sleep must return to the routine in the event of a signal. That is, if they sleep, they must be at priority \leq PZERO, or with PCATCH set in the sleep priority.
2. The *open* routine must return ≥ 0 on success or OPENFAIL if it fails. This ensures that a failure will be reported to the user process. `errno` may be set on failure. However, if the *open* routine returns OPENFAIL and `errno` is not set, STREAMS will automatically set `errno` to ENXIO.
3. If a module or driver recognizes and acts on an M_IOCTL message, it must reply by sending a M_IOCACK message upstream. A unique id is associated with each M_IOCTL, and the M_IOCACK or M_IOCNAK message must contain the id of the M_IOCTL it is acknowledging.
4. A module (not a driver) must pass on any M_IOCTL message it does not recognize (see *Message Types*, above). If an unrecognized M_IOCTL reaches a driver, the driver must reply by sending an M_IOCNAK message upstream.

Data Structures

Only the contents of `q_ptr`, `q_minpsz`, `q_maxpsz`, `q_hiwat`, and `q_lowat` in a `queue_t` structure may be altered. The latter four quantities are set when the module or driver is opened, but may be modified subsequently.

As described in the *SunOS STREAMS Topics* chapter, every module and driver is configured in with the address of a `streamtab` structure (see also the *STREAMS Mechanism* section of the *STREAMS Module and Driver Programming* chapter). For a driver, a pointer to its `streamtab` is included in `cdevsw`. For a module, a pointer to its `streamtab` is included in `fmodsw`.

Header Files

The following header files are generally required in modules and drivers:

types.h

contains type definitions used in the STREAMS header files

stream.h

contains required structure and constant definitions

stropts.h

primarily for users, but contains definitions of the arguments to the M_FLUSH message type also required by modules

One or more of the header files described below may also be included (also see the following section). No standard SunOS system header files should be included except as described in the following section. The intent is to prevent attempts to access data that cannot or should not be accessed.

errno.h

defines various system error conditions, and is needed if errors are to be returned upstream to the user

sysmacros.h

contains miscellaneous system macro definitions

param.h

defines various system parameters, particularly the value of the PCATCH sleep flag

signal.h

defines the system signal values, and should be used if signals are to be processed or sent upstream

file.h

defines the file open flags, and is needed if O_NDELAY is interpreted

Accessible Symbols and Functions

The following lists the only symbols and functions that modules or drivers may refer to (in addition to those defined by STREAMS), if hardware and UNIX-system release independence is to be maintained. Drivers and modules that use symbols not listed here will not be compatible with System V systems.

user.h (from open/close procedures only)

```
struct proc *u_procp      /* process structure pointer */
char u_error              /* system call error number */
ushort u_uid              /* effective user ID */
ushort u_gid              /* effective group ID */
ushort u_ruid             /* real user ID */
ushort u_rgid            /* real group ID */
```

proc.h (from open/close procedures only)

```
short p_pid              /* process ID */
short p_pgrp             /* process group ID */
```

Functions accessible from open/close procedures only

```
flg = sleep(chan, pri)  /* sleep until wakeup */
```

Universally accessible functions

```
bcopy(from, to, nbytes) /* copy data quickly */
bzero(buffer, nbytes)  /* zero data quickly */
t = max(a, b)           /* return max of args */
t = min(a, b)           /* return min of args */
mem=rm_alloc(map, size) /* allocate resource */
rmfree(map, size, addr) /* de-allocate resource */
rminit(mp, size, addr, name, mapsize) /* initialize resource map */
printf(format, ...)    /* print message */
s = spln()              /* set priority level */
timeout(func, arg, ticks) /* schedule event */
untimeout(func, arg)   /* cancel event */
wakeup(chan)           /* wake up sleeper */
```

sysmacros.h

```
t = major(dev)          /* return major device */
t = minor(dev)          /* return minor device */
```

kernel.h

```

struct timeval boottime /* time since system came up */
struct timeval time     /* current time */

```

param.h

```

PZERO           /* zero sleep priority */
PCATCH         /* catch signal sleep flag */
hz             /* clock ticks per second */
NULL           /* 0 */

```

types.h

```

dev_t          /* combined major/minor device */
time_t        /* time counter */

```

All data elements are software read-only except:

```

u_error       /* may be set on a failure return of open */

```

Rules for Put and Service Procedures

To ensure proper data flow between modules, the following rules should be observed in put and service procedures. The following rules pertain to put procedures.

1. A put procedure must not sleep.
2. Each QUEUE must define a put procedure in its `qinit` (see *Kernel Structures*, above) structure for passing messages between modules.
3. A put procedure must use the `putq()` (see *Utilities*, above) utility to enqueue a message on its own message queue. This is necessary to ensure that the various fields of the `queue_t` structure are maintained consistently.
4. When passing messages to a neighbor module, a module may not call `putq()` directly, but must call its neighbor's put procedure (see `putnext()` in *Utilities*). Note that this rule is distinct from the one above it. The previous rule states that a module must call `putq()` to place messages on its own message queue, whereas this rule states that a module must not call `putq()` directly to place messages on a neighbor's queue.

However, the `q_qinfo` structure that points to a module's put procedure may point to `putq()` (i.e. `putq()` is used as the put procedure for that module). When a module calls a neighbor's put procedure that is defined in this manner, it will be calling `putq()` indirectly. If any module uses `putq()` as its put procedure in this manner, the module must define a service procedure. Otherwise, no messages will ever be sent to the next module. Also, because `putq()` does not process `M_FLUSH` messages, any module that uses `putq()` as its put procedure must define a service procedure to process `M_FLUSH` messages.

5. The put procedure of a QUEUE with no service procedure must call the put procedure of the next QUEUE directly, if a message is to be passed to that QUEUE. If flow control is desired, a service procedure must be provided.

Service procedures must observe the following rules:

1. A service procedure must not sleep.
2. The service procedure must use `getq()` to remove a message from its message queue, so that the flow control mechanism is maintained.
3. The service procedure should process all messages on its message queue. The only exception is if the stream ahead is blocked (i.e., `canput()` fails, see *Utilities*, above). Adherence to this rule is the only guarantee that STREAMS will enable (schedule for execution) the service procedure when necessary, and that the flow control mechanism will not fail.

If a service procedure exits for any other reason (e.g., buffer allocation failure), it must take explicit steps to assure it will be re-enabled.

4. The service procedure must follow the steps below for each message that it processes. STREAMS flow control relies on strict adherence to these steps.

Step 1:

Remove the next message from the message queue using `getq()`. It is possible that the service procedure could be called when no messages exist on the queue, so the service procedure should never assume that there is a message on its message queue. If there is no message, return.

Step 2:

If all the following conditions are met:

- `canput()` fails and
 - the message type is not a priority type (see *Message Types*) and
 - the message is to be put on the next QUEUE,
- then, continue at Step 3. Otherwise, continue at Step 4.

Step 3:

The message must be replaced on the head of the message queue from which it was removed using `putbq()` (see *Utilities*). Following this, the service procedure is exited. The service procedure should not be re-enabled at this point. It will be automatically back-enabled by flow control.

Step 4:

If all the conditions of Step 2 are not met, the message should not be returned to the queue. It should be processed as necessary. Then, return to Step 1.

A.6. STREAMS Glossary

<i>Back Enable</i>	To enable (by STREAMS) a preceding blocked QUEUE when STREAMS determines that a succeeding QUEUE has reached its low <i>water mark</i> .
<i>Blocked</i>	A QUEUE that cannot be enabled due to <i>flow control</i> .
<i>Clone Device</i>	A STREAMS device that returns an unused minor device when initially opened, rather than requiring the minor device to be specified in the <code>open(2)</code> call.

<i>Close Procedure</i>	The module routine that is called when a module is popped from a stream and the driver routine that is called when a <i>driver</i> is closed.
<i>Control Stream</i>	In a <i>multiplexor</i> , the <i>upper stream</i> on which a previous <code>I_LINK ioctl</code> (to the associated file, see <code>streamio(4)</code>) caused a <i>lower stream</i> to be connected to the multiplexor driver at the end of the <i>upper stream</i> .
<i>Downstream</i>	The direction from <i>stream head</i> towards <i>driver</i> .
<i>Device Driver</i>	In the STREAMS context, the term “device driver” refers to the end of the <i>stream</i> closest to an external interface. The principal functions of a <i>device driver</i> are handling an associated physical device, and transforming data and information between the external interface and <i>stream</i> .
<i>Driver</i>	A module that forms the <i>stream end</i> . It can be a <i>device driver</i> or a <i>pseudo-device driver</i> . In STREAMS, a <i>driver</i> is physically identical to a module (i.e., composed of two QUEUES), but it has additional attributes.
<i>Enable</i>	Schedule a QUEUE.
<i>Flow Control</i>	The STREAMS mechanism that regulates the flow of messages within a stream and the flow from user space into a stream.
<i>Lower Stream</i>	A stream connected below a multiplexor <i>pseudo-device driver</i> , by means of an <code>I_LINK ioctl</code> . The far end of a lower stream terminates at a <i>device driver</i> or another multiplexor driver.
<i>Message</i>	One or more linked <i>message blocks</i> . A message is referenced by its first message block and its type is defined by the <i>message type</i> of that block.
<i>Message block</i>	Carries data or information, as identified by its <i>message type</i> , in a stream. A <i>message block</i> is a triplet consisting of a data buffer and associated control structures, an <code>mblk_t</code> structure, and a <code>dblkc_t</code> structure.
<i>Message Queue</i>	A linked list of zero or more <i>messages</i> connected to a QUEUE.
<i>Message type</i>	A defined set of values identifying the contents of a <i>message block</i> and <i>message</i> .
<i>Module</i>	A pair of QUEUES. In general, module implies a <i>pushable</i> module.
<i>Multiplexor</i>	A STREAMS mechanism that allows messages to be routed among multiple STREAMS in the kernel. A multiplexor includes at least one multiplexing <i>pseudo-device driver</i> connected to one or more <i>upper</i> STREAMS and one or more <i>lower</i> STREAMS.
<i>Open Procedure</i>	The routine in each STREAMS <i>driver</i> and <i>module</i> called by STREAMS on each <code>open(2)</code> system call made on the stream. A <i>module</i> ’s open procedure is also called when the <i>module</i> is pushed.
<i>Pop</i>	A STREAMS <code>ioctl()</code> (see <code>streamio(4)</code>) that causes the <i>pushable module</i> immediately below the <i>stream head</i> to be removed (popped) from a stream (modules can also be popped as the result of a <code>close(2)</code>).
<i>Pseudo-device Driver</i>	A software <i>driver</i> , not directly associated with a physical device, that performs functions internal to a <i>stream</i> such as a <i>multiplexor</i> or <i>log driver</i> .

<i>Push</i>	A STREAMS <code>ioctl()</code> (see <code>streamio(4)</code>) that causes a <i>pushable module</i> to be inserted (pushed) in a stream immediately below the <i>stream head</i> .
<i>Pushable Module</i>	A module interposed (pushed) between the <i>stream head</i> and <i>driver</i> . Pushable modules perform intermediate transformations on messages flowing between the <i>stream head</i> and <i>driver</i> . A <i>driver</i> is a non-pushable module and a <i>stream head</i> includes a non-pushable module.
<i>Put Procedure</i>	The routine in a QUEUE that receives messages from the preceding QUEUE. It is the single entry point into a QUEUE from a preceding QUEUE. The procedure may perform processing on the message and will then generally either queue the message for subsequent processing by this QUEUE's <i>service procedure</i> , or will pass the message to the <i>put procedure</i> of the following QUEUE.
<i>QUEUE</i>	A STREAMS defined set of C structures. A module is composed of a read (<i>upstream</i>) QUEUE and a write (<i>downstream</i>) QUEUE. A QUEUE will typically contain a put and service procedure, a <i>message queue</i> , and private data. The read QUEUE (cf. <i>read queue</i>) in a <i>module</i> will also contain the open procedure and close procedure for the <i>module</i> . The primary structure is the <code>queue_t</code> structure, occasionally used as a synonym for a QUEUE.
<i>Read Queue</i>	The <i>message queue</i> in a <i>module</i> or <i>driver</i> containing <i>messages</i> moving <i>upstream</i> . Associated with a <code>read(2)</code> system call and input from a <i>driver</i> .
<i>Schedule</i>	Place a QUEUE on the internal list of QUEUES which will subsequently have their service procedure called by the STREAMS scheduler.
<i>Service Interface</i>	A set of primitives that define a service at the boundary between a <i>service user</i> and a <i>service provider</i> and the rules (typically represented by a state machine) for allowable sequences of the primitives across the boundary. At a stream/user boundary, the primitives are typically contained in the control part of a message; within a stream, in <code>M_PROTO</code> or <code>M_PCPROTO</code> message blocks.
<i>Service Procedure</i>	The routine in a QUEUE that receives messages queued for it by the <i>put procedure</i> of the QUEUE. The procedure is called by the STREAMS <i>scheduler</i> . It may perform processing on the message and will generally pass the message to the <i>put procedure</i> of the following QUEUE.
<i>Service Provider</i>	In a <i>service interface</i> , the entity (typically a <i>module</i> or <i>driver</i>) that responds to request primitives from the <i>service user</i> with response and event primitives.
<i>Service User</i>	In a <i>service interface</i> , the entity that generates request primitives for the <i>service provider</i> and consumes response and event primitives.
<i>Stream</i>	The kernel aggregate created by connecting STREAMS components, resulting from an application of the STREAMS mechanism. The primary components are the <i>stream head</i> , the <i>driver</i> , and zero or more <i>pushable modules</i> between the <i>stream head</i> and <i>driver</i> .
<i>Stream End</i>	The end of the <i>stream</i> furthest from the user process, containing a <i>driver</i> .
<i>Stream Head</i>	The end of the <i>stream</i> closest to the user process. It provides the interface between the <i>stream</i> and the user process.

- STREAMS** A kernel mechanism that supports development of network services and data communication *drivers*. It defines interface standards for character input/output within the kernel, and between the kernel and user level. The STREAMS mechanism comprises integral functions, utility routines, kernel facilities, and a set of structures.
- Upper Stream** A stream terminating above a multiplexor *pseudo-device driver*. The far end of an upper stream originates at the *stream head* or another multiplexor driver.
- Upstream** The direction from *driver* towards *stream head*.
- Water Marks** Limit values used in *flow control*. Each QUEUE has a high water mark and a low water mark. The high water mark value indicates the upper limit related to the number of characters contained on the *message queue* of a QUEUE. When the enqueued characters in a QUEUE reach its high water mark, STREAMS causes another QUEUE that attempts to send a message to this QUEUE to become *blocked*. When the characters in this QUEUE are reduced to the low water mark value, the other QUEUE will be unblocked by STREAMS.
- Write queue** The *message queue* in a *module* or *driver* containing *messages* moving *downstream*. Associated with a `write(2)` system call and output from a user process.

Index

A

accessing the datagram provider, STREAMS, 64
adding modules, STREAMS, 76
advanced operations, STREAMS, 46
advanced topics, STREAMS, 132
application programming, STREAMS, 37
autopush, SunOS extension to STREAMS, 143

B

bappend(), 84
basic operations, STREAMS, 40
blocking, process, 49
buffer allocation priority, STREAMS, 159
building a multiplexor, STREAMS, 52

C

canonical processing, 90
CANONPROC, STREAMS, 25
cdevsw, upon open, 75, 140
clone open, STREAMS, 50
cloning, STREAMS, 104
closing STREAMS, 111
closing a service, STREAMS, 67
closing a stream, STREAMS, 76

D

data structures
 STREAMS, 170
data, routing multiplexed, STREAMS, 58
datagram
 receiving a STREAMS datagram, 68
 sending a STREAMS datagram, 67
datagram service interface, STREAMS, 62
definition of “stream” in STREAMS context, 38
development facilities, STREAMS, 39
dismantling a multiplexor, STREAMS, 57
driver
 STREAMS close, 103
 STREAMS declarations, 96
 STREAMS development facilities, 72
 STREAMS environment, 80
 STREAMS flow control, 95
 STREAMS flush handling, 100
 STREAMS interrupt, 100
 STREAMS ioctls, 101

driver, *continued*

 STREAMS module programming, 71
 STREAMS open, 98
 STREAMS processing procedures, 99
 STREAMS programming, 96

E

error logging, STREAMS, 32
error messages, STREAMS, 141
example driver, STREAMS, 104
example of STREAMS multiplexing, 115
example of STREAMS multiplexor configuration, 120
external variables, STREAMS, 77

F

filter module declarations, STREAMS, 83
flow control, STREAMS, 27, 89, 173
fmodsw, STREAMS, 76
freeing messages, STREAMS, 77

H

high priority messages, STREAMS, 156

I

inserting modules, STREAMS, 42
ioctl
 STREAMS I_LINK, 113
 STREAMS I_PLINK, 113

K

kernel
 STREAMS functions, 19
kernel structures, STREAMS, 147

L

line disciplines, STREAMS, 144
log command, 32
loop-around driver, STREAMS, 104

M

manual overview, 1
message allocation, STREAMS, 21, 85
message blocks, STREAMS, 20
message form and linkage, STREAMS, 82
message format, STREAMS, 81
message generation, STREAMS, 83

message handling, STREAMS, 59 *thru* 70
 message interface, STREAMS, 60
 message priority, STREAMS, 89
 message queues, STREAMS, 19, 87, 89
 message reception, STREAMS, 83
 message structures, STREAMS, 149
 message types, STREAMS, 15, 150
 modularity STREAMS, 2
 module and driver control, STREAMS, 43
 module configuration, STREAMS, 140
 module declarations, STREAMS, 77
 module environment, STREAMS, 80
 module ioctls, STREAMS, 101
 module procedures, STREAMS, 79
 module reusability, STREAMS, 7
 multiplexed data, routing, STREAMS, 58
 multiplexed Streams, STREAMS, 50 *thru* 58
 multiplexing driver, STREAMS, 117
 multiplexing,
 STREAMS", 111

O

O_NDELAY, 76
 opening a stream, STREAMS, 75
 overview of STREAMS drivers, 94

P

protocol portability, STREAMS, 6
 protocol substitution, STREAMS, 7
 put procedure rules, STREAMS, 172
 put procedure, STREAMS, 86
 put procedures, STREAMS, 22

Q

QUEUE data structures
 QUEUE data structures, STREAMS , 74
 queue priority, STREAMS, 89
 queue_t, 87

R

receiving a STREAMS datagram, 68
 receiving messages, STREAMS, 16
 removing modules, STREAMS, 76
 routing multiplexed data, STREAMS, 58

S

sending a datagram, STREAMS, 67
 sending messages, STREAMS, 16
 service interface messages, STREAMS, 59
 service interface, STREAMS, 127
 service procedure rules, STREAMS, 172
 service procedures, STREAMS, 22, 87, 88, 110
 strace command, 32
 stream construction, STREAMS, 73
 stream end, 10
 stream head, STREAMS, 9
 STREAMS
 accessible functions, 171

STREAMS, *continued*

accessible symbols, 171
 accessing the datagram provider, 64
 adding modules, 76
 advanced I/O, 46
 advanced operations, 46
 advanced topics, 132
 advanced topics, flow control, 135
 advanced topics, read options, 136
 advanced topics, recovering from no buffers, 132
 advanced topics, signals, 135
 advanced topics, Stream Head processing, 136
 advanced topics, write offset, 137
 advanced view, 8
 application programming, 37 *thru* 70
 asynchronous I/O, 49
 asynchronous protocol example, 14
 bappend(), 84
 basic operations, 40
 basic view, 3
 benefits, 5
 big example, 104
 buffer allocation priority, 159
 building a multiplexor, 52
 building a stream, 11
 CANONPROC, 25
 clone open, 50
 cloning, 104
 close, 111
 closing a service, 67
 closing a stream, 76
 configuring drivers, 139
 context, 1
 creating service interfaces, 5
 data structures, 170
 datagram provider access, 64
 datagram receiving, 68
 datagram service interface, 62
 definition of "stream", 38
 design and system calls, 170
 design guidelines, 169
 development facilities, 39
 device driver streams, 95
 dismantling a multiplexor, 57
 driver cdevsw interface, 139
 driver close, 103
 driver declarations, 96
 driver development facilities, 72
 driver environment, 80
 driver flow control, 95
 driver flush handling, 100
 driver interrupt, 100
 driver ioctls, 101
 driver open, 98
 driver processing, 24
 driver processing procedures, 99
 driver programming, 96
 drivers, 34
 environment, 34
 error and trace logging, 32
 error messages, 141
 example of asynchronous protocol, 14
 example of message use, 16
 expanded streams, 11

STREAMS, *continued*

- external variables, 77
- filter module declarations, 83
- flow control, 27, 89, 173
- fmodsw, 76
- freeing messages, 77
- functional parts, 3
- functions, accessible, 171
- general design rules, 169
- glossary, 35, 173
- header files, 170
- high priority messages, 156
- I/O, advanced, 46
- I/O, asynchronous, 49
- I/O, polling, 46
- I_LINK, 113
- I_PLINK, 113
- inserting modules, 42
- Internet multiplexing, 29
- Internet multiplexor after connecting, 116
- Internet multiplexor before connecting, 115
- introduction, 1
- ioctl's, 101
- kernel level functions, 19
- kernel processing, 23
- kernel structures, 147
- kernel structures, *iocblk*, 149
- kernel structures, *linkblk*, 150
- kernel structures, *streamtab*, 147
- kernel structures, QUEUE, 147
- line disciplines, 144
- loop-around driver, 104
- M_PCPROTO messages, 151
- M_PROTO messages, 151
- manipulating modules, 6
- manipulating STREAMS modules, 6
- manual pages, 5
- mechanism, 72
- message allocation, 21, 85
- message blocks, 20
- message form and linkage, 82
- message format, 81
- message generation, 83
- message handling, 59 *thru* 70
- message interface, 60
- message priority, 89
- message queue priority, 26
- message queues, 19, 87, 89
- message reception, 83
- message structures, 149
- message types, 15, 150
- message types, ordinary, 150
- message use example, 16
- modularity, 2
- module and driver control, 43
- module configuration, 140
- module declarations, 77
- module environment, 80
- module ioctl's, 101
- module procedures, 79
- module programming, 71
- module reusability, 7
- modules, 9, 34
- monitoring, 31

STREAMS, *continued*

- multiplexed Streams, 50 *thru* 58
- multiplexing, 29, 111
- multiplexing configurations, 112
- multiplexing driver, 117
- multiplexing, connecting lower streams, 113
- multiplexing, disconnecting lower streams, 114
- multiplexing, example, 115
- multiplexor configuration, example, 120
- multiplexor, lower QUEUE write, 123
- multiplexor, lower read put, 125
- multiplexor, upper write put, 120
- opening a stream, 75
- overview of drivers, 94
- polling I/O, 46
- portability, 143
- primer, 1
- protocol migration, 7
- protocol portability, 6
- protocol substitution, 7
- pushable modules, 12
- put and service procedures, 22
- put procedure, 86
- put procedure rules, 172
- QUEUE data structures, 74
- queue priority, 89
- queue_t*, 87
- receiving a datagram, 68
- removing modules, 76
- routing multiplexed data, 58
- sending a datagram, 67
- sending and receiving messages, 16
- service interface, 127
- service interface messages, 59
- service interface procedure, 130
- service interface, declarations, 128
- service interface, messages, 127
- service interfaces, 5
- service procedure rules, 172
- service procedures, 87, 88, 110
- single I/O pathway, 38
- standard SunOS modules, 142
- stream construction, 73
- stream end, 10
- stream head, 9
- streamtab*, 140
- SunOS STREAMS Topics, 139
- SunOS extension, 143
- SunOS modules, 142
- supplementary material, 147
- symbols, accessible, 171
- system calls, 4, 13
- system error messages, 141
- tunable parameters, 141
- user line disciplines, 144
- utilities, 158
- write put procedure, 107
- write side processing, 25
- X.25 multiplexing, 30

STREAMS utilities

- adjmsg()*, 159
- allocb()*, 160
- backq()*, 160
- bufcall()*, 160

STREAMS utilities, *continued*

canput (), 160
copyb (), 161
copymsg (), 161
dupb (), 161
dupmsg (), 162
enableok (), 162
esballoc (), 162
esbcall (), 162
flushq (), 163
freeb (), 163
freemsg (), 163
insq (), 164
linkb (), 164
msgdsize (), 164
noenable (), 164
OTHERQ (), 164
pullupmsg (), 165
putbq (), 165
putctl (), 165
putctl1 (), 165
putnext (), 166
putq (), 166
qenable (), 166
qreply (), 167
qsize (), 167
RD (), 167
rmvb (), 167
rmvq (), 167
setq (), 163
splstr (), 168
sx (), 161
testb (), 168
unbufcall (), 168
unlinkb (), 168
WB (), 168
x (), 160
strerr command, 32
system calls, STREAMS, 4

T

trace logging, STREAMS, 32
tunable parameters, STREAMS, 141

U

utilities, STREAMS, 158

W

write put procedure, STREAMS, 107

Notes

Notes