#### value?

ter 12

by the

on one Ilticast

of the

nod to

multi-

ice has )K? tion is

ised to

# 13

# IGMP: Internet Group Management Protocol

## 13.1 Introduction

IGMP conveys group membership information between hosts and routers on a local network. Routers periodically multicast IGMP queries to the all-hosts group. Hosts respond to the queries by multicasting IGMP report messages. The IGMP specification appears in RFC 1112. Chapter 13 of Volume 1 describes the specification of IGMP and provides some examples.

From an architecture perspective, IGMP is a transport protocol above IP. It has a protocol number (2) and its messages are carried in IP datagrams (as with ICMP). IGMP usually isn't accessed directly by a process but, as with ICMP, a process can send and receive IGMP messages through an IGMP socket. This feature enables multicast routing daemons to be implemented as user-level processes.

Figure 13.1 shows the overall organization of the IGMP protocol in Net/3.

The key to IGMP processing is the collection of in\_multi structures shown in the center of Figure 13.1. An incoming IGMP query causes igmp\_input to initialize a countdown timer for each in\_multi structure. The timers are updated by igmp\_fasttimo, which calls igmp\_sendreport as each timer expires.

We saw in Chapter 12 that ip\_setmoptions calls igmp\_joingroup when a new in\_multi structure is created. igmp\_joingroup calls igmp\_sendreport to announce the new group and enables the group's timer to schedule a second announcement a short time later. igmp\_sendreport takes care of formatting an IGMP message and passing it to ip\_output.

On the left and right of Figure 13.1 we see that a raw socket can send and receive IGMP messages directly.

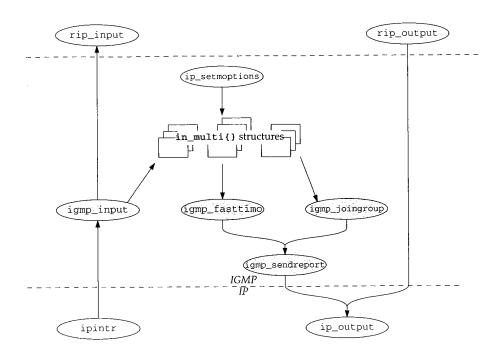


Figure 13.1 Summary of IGMP processing.

## 13.2 Code Introduction

The IGMP protocol is implemented in four files listed in Figure 13.2.

File	Description
netinet/igmp.h	IGMP protocol definitions
netinet/igmp_var.h	IGMP implementation definitions
netinet/in_var.h	IP multicast data structures
netinet/igmp.c	IGMP protocol implementation

Figure 13.2 Files discussed in this chapter.

## **Global Variables**

Three new global variables, shown in Figure 13.3, are introduced in this chapter.

## Statistics

IGMP statistics are maintained in the igmpstat variables shown in Figure 13.4.

Chapter 13

Ģ

**INTEL Ex.1013.408** 

SNMP

Section 13.2

Variable	Datatype	Description
igmp_all_hosts_group	u_long	all-hosts group address in network byte order
igmp_timers_are_running	int	true if any IGMP timer is active, false otherwise
igmpstat	struct igmpstat	IGMP statistics (Figure 13.4).

Figure 13.3 Global variables introduced in this chapter.

igmpstat member	Description
igps_rcv_badqueries	#messages received as invalid queries
igps_rcv_badreports	#messages received as invalid reports
igps_rcv_badsum	#messages received with bad checksum
igps_rcv_ourreports	#messages received as reports for local groups
igps_rcv_queries	#messages received as membership queries
igps_rcv_reports	#messages received as membership reports
igps_rcv_tooshort	#messages received with too few bytes
igps_rcv_total	total #IGMP messages received
igps_snd_reports	#messages sent as membership reports



Figure 13.5 shows some sample output of these statistics, from the netstat -p igmp command on vangogh.cs.berkeley.edu.

netstat -p igmp output	igmpstat member
18774 messages received	igps_rcv_total
0 messages received with too few bytes	igps_rcv_tooshort
0 messages received with bad checksum	igps_rcv_badsum
18774 membership queries received	igps_rcv_queries
0 membership queries received with invalid field(s)	igps_rcv_badqueries
0 membership reports received	igps_rcv_reports
0 membership reports received with invalid field(s)	igps_rcv_badreports
0 membership reports received for groups to which we belong	igps_rcv_ourreports
0 membership reports sent	igps_snd_reports

Figure 13.5 Sample IGMP statistics.

From Figure 13.5 we can tell that vangogh is attached to a network where IGMP is being used, but that vangogh is not joining any multicast groups, since igps\_snd\_reports is 0.

#### **SNMP Variables**

There is no standard SNMP MIB for IGMP, but [McCloghrie and Farinacci 1994a] describes an experimental MIB for IGMP.

## 13.3 igmp Structure

An IGMP message is only 8 bytes long. Figure 13.6 shows the igmp structure used by Net/3.

	····		igmp.h
43 str	uct igmp {		0,
44	u_char igmp_type;	<pre>/* version &amp; type of IGMP message</pre>	*/
45	u_char igmp_code;	<pre>/* unused, should be zero</pre>	*/
46	u_short igmp_cksum;	/* IP-style checksum	*/
47	<pre>struct in_addr igmp_group;</pre>	<pre>/* group address being reported</pre>	*/
48 };		<pre>/* (zero for queries)</pre>	*/
			igmp.h

Figure 13.6 igmp structure.

43-44 A 4-bit version code and a 4-bit type code are contained within igmp\_type. Figure 13.7 shows the standard values.

Version	Туре	igmp_type	Description
1	1	0x11 (IGMP_HOST_MEMBERSHIP_QUERY)	membership query
1	2	0x12 (IGMP_HOST_MEMBERSHIP_REPORT)	membership report
1	3	0x13	DVMRP message (Chapter 14)

#### Figure 13.7 IGMP message types.

Only version 1 messages are used by Net/3. Multicast routers send type 1 (IGMP\_HOST\_MEMBERSHIP\_QUERY) messages to solicit membership reports from hosts on the local network. The response to a type 1 IGMP message is a type 2 (IGMP\_HOST\_MEMBERSHIP\_REPORT) message from the hosts reporting their multicast membership information. Type 3 messages transport multicast routing information between routers (Chapter 14). A host never processes type 3 messages. The remainder of this chapter discusses only type 1 and 2 messages.

<sup>45-46</sup> igmp\_code is unused in IGMP version 1, and igmp\_cksum is the familiar IP checksum computed over all 8 bytes of the IGMP message.

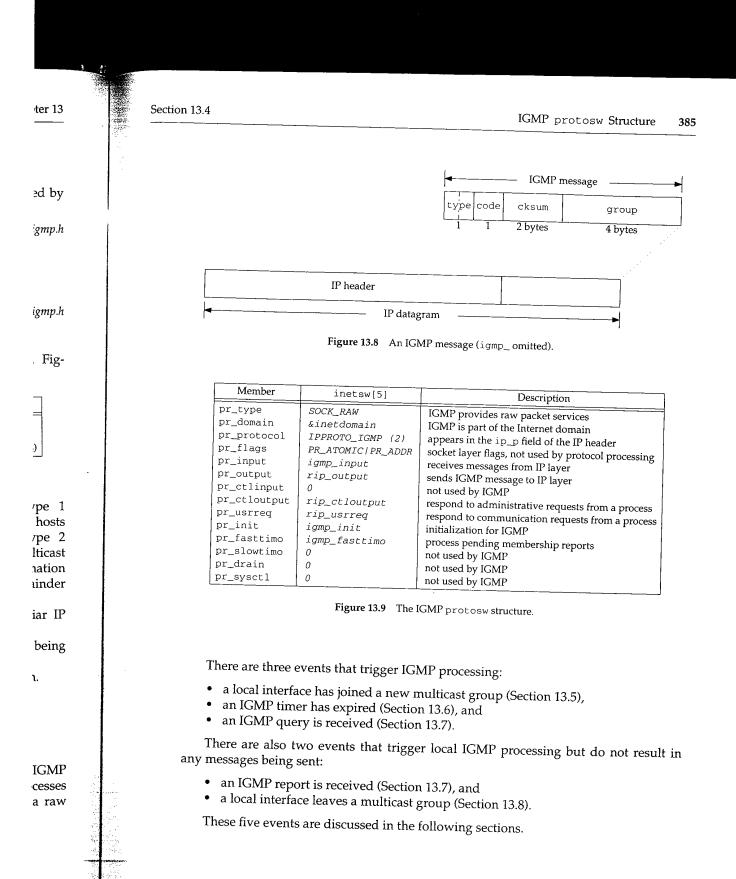
<sup>47-48</sup> igmp\_group is 0 for queries. For replies, it contains the multicast group being reported.

Figure 13.8 shows the structure of an IGMP message relative to an IP datagram.

### 13.4 IGMP protosw Structure

Figure 13.9 describes the protosw structure for IGMP.

Although it is possible for a process to send raw IP packets through the IGMP protosw entry, in this chapter we are concerned only with how the kernel processes IGMP messages. Chapter 32 discusses how a process can access IGMP using a raw socket.



INTEL Ex.1013.411

Chapter 13

## 13.5 Joining a Group: igmp\_joingroup Function

We saw in Chapter 12 that igmp\_joingroup is called by in\_addmulti when a new in\_multi structure is created. Subsequent requests to join the same group only increase the reference count in the in\_multi structure; igmp\_joingroup is not called. igmp\_joingroup is shown in Figure 13.10

164 void	
165 igmp_joingroup(inm)	
166 struct in_multi *inm;	
167 {	
168 int s = splnet();	
169 if (inm->inm_addr.s_addr == igmp_all_hosts_group	
inm->inm_ifp == &loif)	
<pre>inm-&gt;inm_timer = 0;</pre>	
172 else {	
igmp_sendreport(inm); 173 igmp_sendreport(inm);	;
<pre>173 igmp_sendlepole(inm, / RANDOM_DELAY(inm-&gt;inm_addr)) 174 inm-&gt;inm_timer = IGMP_RANDOM_DELAY(inm-&gt;inm_addr))</pre>	
175 igmp_timers_are_running = 1;	
176 }	
177 splx(s);	
178	

## \_\_\_\_ igmp.c

Figure 13.10 igmp\_joingroup function.

<sup>164-178</sup> inm points to the new in\_multi structure for the group. If the new group is the all-hosts group, or the membership request is for the loopback interface, inm\_timer is disabled and igmp\_joingroup returns. Membership in the all-hosts group is never reported, since every multicast host is assumed to be a member of the group. Sending a membership report to the loopback interface is unnecessary, since the local host is the only system on the loopback network and it already knows its membership status.

In the remaining cases, a report is sent immediately for the new group, and the group timer is set to a random value based on the group. The global flag igmp\_timers\_are\_running is set to indicate that at least one timer is enabled. igmp\_fasttimo (Section 13.6) examines this variable to avoid unnecessary processing.

When the timer for the new group expires, a second membership report is issued. The duplicate report is harmless, but it provides insurance in case the first report is lost or damaged. The report delay is computed by IGMP\_RANDOM\_DELAY (Figure 13.11).

According to RFC 1122, report timers should be set to a random time between 0 and 10 (IGMP\_MAX\_HOST\_REPORT\_DELAY) seconds. Since IGMP timers are decremented five (PR\_FASTHZ) times per second, IGMP\_RANDOM\_DELAY must pick a random value between 1 and 50. If r is the random number computed by adding the total number of IP packets received, the host's primary IP address, and the multicast group, then

 $0 \leq (r \bmod 50) \leq 49$ 

and

59-73

$$1 \le (r \mod 50) + 1 \le 50$$

13.6

Sectio

154-16

語語ではあります

Section 13.6

С

.C

ıe

is

er

a ie

ıe

١g

d.

g.

d.

st

ιd

зd

ıe

of

```
59 /*
                                                                       igmp_var.h
60 * Macro to compute a random timer value between 1 and (IGMP_MAX_REPORTING_
61 * DELAY * countdown frequency). We generate a "random" number by adding
62 * the total number of IP packets received, our primary IP address, and the
63 * multicast address being timed-out. The 4.3 random() routine really
   * ought to be available in the kernel!
64
65
   */
66 #define IGMP_RANDOM_DELAY(multiaddr) \
67
      /* struct in_addr multiaddr; */ \
68
       ( (ipstat.ips_total + \
          ntohl(IA_SIN(in_ifaddr)->sin_addr.s_addr) + \
69
70
          ntohl((multiaddr).s_addr) \
71
         ) \
72
          (IGMP_MAX_HOST_REPORT_DELAY * PR_FASTHZ) + 1 \
         8
73
      )
```

Figure 13.11 IGMP\_RANDOM\_DELAY function.

Zero is avoided because it would disable the timer and no report would be sent.

## 13.6 igmp\_fasttimo Function

Before looking at igmp\_fasttimo, we need to describe the mechanism used to traverse the in\_multi structures.

To locate each in\_multi structure, Net/3 must traverse the in\_multi list for each interface. During a traversal, an in\_multistep structure (shown in Figure 13.12) records the position.

123 struct in\_multistep {
124 struct in\_ifaddr \*i\_ia;
125 struct in\_multi \*i\_inm;
126 };

Figure 13.12 in\_multistep function.

— in\_var.h

in\_var.h

igmp\_var.h

<sup>123-126</sup> i\_ia points to the *next* in\_ifaddr interface structure and i\_inm points to the *next* in\_multi structure for the *current* interface.

The IN\_FIRST\_MULTI and IN\_NEXT\_MULTI macros (shown in Figure 13.13) traverse the lists.

<sup>154–169</sup> If the in\_multi list has more entries, i\_inm is advanced to the next entry. When IN\_NEXT\_MULTI reaches the end of a multicast list, i\_ia is advanced to the next interface and i\_inm to the first in\_multi structure associated with the interface. If the interface has no multicast structures, the while loop continues to advance through the interface list until all interfaces have been searched.

<sup>170-177</sup> The in\_multistep array is initialized to point to the first in\_ifaddr structure in the in\_ifaddr list and i\_inm is set to null. IN\_NEXT\_MULTI finds the first in\_multi structure.

#### **INTEL Ex.1013.413**

Chapter 13

-in var.h

-in var.h

147 /\* 148 \* Macro to step through all of the in\_multi records, one at a time. 149 \* The current position is remembered in "step", which the caller must \* provide. IN\_FIRST\_MULTI(), below, must be called to initialize "step" 150 151 \* and get the first record. Both macros return a NULL "inm" when there \* are no remaining records. 152 153 \*/ 154 #define IN\_NEXT\_MULTI(step, inm) \ /\* struct in\_multistep step; \*/ \ 155 /\* struct in\_multi \*inm; \*/ \ 156 157 { \ 158 if (((inm) = (step).i\_inm) != NULL) \ (step).i\_inm = (inm)->inm\_next; \ 159 else \ 160 161 while ((step).i\_ia != NULL) { \ (inm) = (step).i\_ia->ia\_multiaddrs; \ 162 163 (step).i\_ia = (step).i\_ia->ia\_next; \ if ((inm) != NULL) { \ 164 (step).i\_inm = (inm)->inm\_next; \ 165 166 break;  $\setminus$ 167 } \ 168 } 169 } 170 #define IN\_FIRST\_MULTI(step, inm) \ 171 /\* struct in\_multistep step; \*/ \ /\* struct in\_multi \*inm; \*/ \ 172 173 { \ 174 (step).i\_ia = in\_ifaddr; \ (step).i\_inm = NULL;  $\$ 175 176 IN\_NEXT\_MULTI((step), (inm)); \ 177 }

Figure 13.13 IN\_FIRST\_MULTI and IN\_NEXT\_MULTI structures.

We know from Figure 13.9 that igmp\_fasttimo is the fast timeout function for IGMP and is called five times per second. igmp\_fasttimo (shown in Figure 13.14) decrements multicast report timers and sends a report when the timer expires.

If igmp\_timers\_are\_running is false, igmp\_fasttimo returns immediately 187-198 instead of wasting time examining each timer.

199–213

igmp\_fasttimo resets the running flag and then initializes step and inm with IN\_FIRST\_MULTI. The igmp\_fasttimo function locates each in\_multi structure with the while loop and the IN\_NEXT\_MULTI macro. For each structure:

- If the timer is 0, there is nothing to be done.
- If the timer is nonzero, it is decremented. If it reaches 0, an IGMP membership ٠ report is sent for the group.
- If the timer is still nonzero, then at least one timer is still running, so igmp\_timers\_are\_running is set to 1.

ię

21

Section 13.6

ter 13

var.h

з'

Э

ı var.h

on for

13.14)

liately

1 with

acture

ership

1g, so

igmp.c 187 void 188 igmp\_fasttimo() 189 ( 190 struct in\_multi \*inm; 191 int s; 192 struct in\_multistep step; 193 /\* \* Quick check to see if any work needs to be done, in order 194 195 \* to minimize the overhead of fasttimo processing. 196 \*/ if (!igmp\_timers\_are\_running) 197 198 return; 199 s = splnet(); 200 igmp\_timers\_are\_running = 0; 201 IN\_FIRST\_MULTI(step, inm); 202 while (inm != NULL) { if  $(inm->inm_timer == 0)$  ( 203 204 /\* do nothing \*/ 205 } else if (--inm->inm\_timer == 0) { 206 igmp\_sendreport(inm); 207 } else { 208 igmp\_timers\_are\_running = 1; 209 } 210 IN\_NEXT\_MULTI(step, inm); 211 } 212 splx(s); 213 } igmp.c Figure 13.14 igmp\_fasttimo function. igmp\_sendreport Function The igmp\_sendreport function (shown in Figure 13.15) constructs and sends an IGMP report message for a single multicast group. The single argument inm points to the in\_multi structure for the group being 214-232 reported. igmp\_sendreport allocates a new mbuf and prepares it for an IGMP message. igmp\_sendreport leaves room for a link-layer header and sets the length of the mbuf and packet to the length of an IGMP message.

- <sup>233-245</sup> The IP header and IGMP message is constructed one field at a time. The source address for the datagram is set to INADDR\_ANY, and the destination address is the multicast group being reported. ip\_output replaces INADDR\_ANY with the unicast address of the outgoing interface. Every member of the group receives the report as does every multicast router (since multicast routers receive *all* IP multicasts).
  - Finally, igmp\_sendreport constructs an ip\_moptions structure to go along with the message sent to ip\_output. The interface associated with the in\_multi structure is selected as the outgoing interface; the TTL is set to 1 to keep the report on the local network; and, if the local system is configured as a router, multicast loopback is enabled for this request.

390 IGMP: Internet Group Management Protocol

「「「「「「「」」」」

igmp.c 214 static void 215 igmp\_sendreport(inm) 216 struct in\_multi \*inm; 217 { struct mbuf \*m; 218 219 struct igmp \*igmp; 220 struct ip \*ip; struct ip\_moptions \*imo; 221 struct ip\_moptions simo; 222 MGETHDR(m, M\_DONTWAIT, MT\_HEADER); 223 224 if (m == NULL) 225 return; 226 /\* \* Assume max\_linkhdr + sizeof(struct ip) + IGMP\_MINLEN 227 \* is smaller than mbuf size returned by MGETHDR. 228 \*/ 229 230 m->m\_data += max\_linkhdr; m->m\_len = sizeof(struct ip) + IGMP\_MINLEN; 231 m->m\_pkthdr.len = sizeof(struct ip) + IGMP\_MINLEN; 232 233 ip = mtod(m, struct ip \*); 234  $ip \rightarrow ip_tos = 0;$ ip->ip\_len = sizeof(struct ip) + IGMP\_MINLEN; 235 ip->ip\_off = 0; 236 ip->ip\_p = IPPROTO\_IGMP; 237 ip->ip\_src.s\_addr = INADDR\_ANY; 238 239 ip->ip\_dst = inm->inm\_addr; igmp = (struct igmp \*) (ip + 1);240 igmp->igmp\_type = IGMP\_HOST\_MEMBERSHIP\_REPORT; 241 igmp->igmp\_code = 0; 242 243 igmp->igmp\_group = inm->inm\_addr; igmp->igmp\_cksum = 0; 244 igmp->igmp\_cksum = in\_cksum(m, IGMP\_MINLEN); 245 246 imo = &simo; bzero((caddr\_t) imo, sizeof(\*imo)); 247 imo->imo\_multicast\_ifp = inm->inm\_ifp; 248 249 imo->imo\_multicast\_ttl = 1; 250 /\* \* Request loopback of the report if we are acting as a multicast 251 \* router, so that the process-level routing demon can hear it. 252 253 \*/ 254 { extern struct socket \*ip\_mrouter; 255 imo->imo\_multicast\_loop = (ip\_mrouter != NULL); 256 257 } ip\_output(m, NULL, NULL, 0, imo); 258 259 ++igmpstat.igps\_snd\_reports; 260 }

Figure 13.15 igmp\_sendreport function.

**INTEL Ex.1013.416** 

igmp.c

Section 13.7

391

The process-level multicast router must hear the membership reports. In Section 12.14 we saw that IGMP datagrams are always accepted when the system is configured as a multicast router. Through the normal transport demultiplexing code, the messages are passed to igmp\_input, the pr\_input function for IGMP (Figure 13.9).

#### Input Processing: igmp\_input Function 13.7

In Section 12.14 we described the multicast processing portion of ipintr. We saw that a multicast router accepts any IGMP message, but a multicast host accepts only IGMP messages that arrive on an interface that is a member of the destination multicast group (i.e., queries and membership reports for which the receiving interface is a member).

The accepted messages are passed to igmp\_input by the standard protocol demultiplexing mechanism. The beginning and end of igmp\_input are shown in Figure 13.16. The code for each IGMP message type is described in following sections.

#### Validate IGMP message

The function ipintr passes m, a pointer to the received packet (stored in an mbuf), 52-96 and iphlen, the size of the IP header in the datagram.

The datagram must be large enough to contain an IGMP message (IGMP\_MINLEN), must be contained within a standard mbuf header (m\_pullup), and must have a correct IGMP checksum. If any errors are found, they are counted, the datagram is silently discarded, and igmp\_input returns.

The body of igmp\_input processes the validated messages based on the code in igmp\_type. Remember from Figure 13.6 that igmp\_type includes a version code and a type code. The switch statement is based on the combined value stored in igmp\_type (Figure 13.7). Each case is described separately in the following sections.

## Pass IGMP messages to raw IP

There is no default case for the switch statement. Any valid message (i.e., one 157-163 that is properly formed) is passed to rip\_input where it is delivered to any process listening for IGMP messages. IGMP messages with versions or types that are unrecognized by the kernel can be processed or discarded by the listening processes.

> The mrouted program depends on this call to rip\_input so that it receives membership queries and reports.

## Membership Query: IGMP\_HOST\_MEMBERSHIP\_QUERY

RFC 1075 recommends that multicast routers issue an IGMP membership query at least once every 120 seconds. The query is sent to group 224.0.0.1 (the all-hosts group). Figure 13.17 shows how the message is processed by a host.

Chapter 13

1. Carlos - Carlos -

「「「「「「」」」

「「「「「「「「「」」」

And a state of the second

-

法の法律にはないないない

**Name** 

· igmp.c 52 void 53 igmp\_input(m, iphlen) 54 struct mbuf \*m; 55 int iphlen; 56 { struct igmp \*igmp; 57 struct ip \*ip; 58 59 igmplen; int struct ifnet \*ifp = m->m\_pkthdr.rcvif; 60 61 int minlen; struct in\_multi \*inm; 62 63 struct in\_ifaddr \*ia; 64 struct in\_multistep step; 65 ++igmpstat.igps\_rcv\_total; ip = mtod(m, struct ip \*); 66 67 igmplen = ip->ip\_len; 68 /\* \* Validate lengths 69 \*/ 70 if (igmplen < IGMP\_MINLEN) { 71 ++igmpstat.igps\_rcv\_tooshort; 72 73 m\_freem(m); 74 return; 75 } minlen = iphlen + IGMP\_MINLEN; 76 if ((m->m\_flags & M\_EXT || m->m\_len < minlen) && 77 78  $(m = m_pullup(m, minlen)) == 0) {$ 79 ++igmpstat.igps\_rcv\_tooshort; 80 return; 81 } /\* 82 \* Validate checksum 83 \*/ 84 m->m\_data += iphlen; 85 86 m->m\_len -= iphlen; 87 igmp = mtod(m, struct igmp \*); 88 if (in\_cksum(m, igmplen)) { 89 ++igmpstat.igps\_rcv\_badsum; 90 m\_freem(m); 91 return; 92 } 93 m->m\_data -= iphlen; 94 m~>m\_len += iphlen; 95 ip = mtod(m, struct ip \*); 96 switch (igmp->igmp\_type) { /\* switch cases \*/ 157 }

S

Section 13.7

igmp.c

igmp.c

158 /\*
159 \* Pass all valid IGMP packets up to any process(es) listening
160 \* on a raw IGMP socket.
161 \*/
162 rip\_input(m);
163 }

Figure 13.16 igmp\_input function.

97	case IGMP_HOST_MEMBERSHIP_QUERY:	— igmp.c
98	++igmpstat.igps_rcv_queries;	0 1
	('igmpocae.igps_icv_queries;	
99	if (ifp == &loif)	
100	break;	
101	if (ip->ip_dst.s_addr != igmp_all_hosts_group) {	
102	++igmpstat.igps_rcv_badqueries;	
103	<pre>m_freem(m);</pre>	
104	return;	
105	}	
106	/*	
107	* Start the timers in all of our membership records for	
108	* the interface on which the query arrived, except those	
109	* that are already running and those that belong to the	
110	* "all-hosts" group.	
111	*/	
112	<pre>IN_FIRST_MULTI(step, inm);</pre>	
113	while (inm != NULL) (	
114	if (inm->inm_ifp == ifp && inm->inm_timer == 0 &&	
115	inm~>inm_addr.s_addr != igmp_all_hosts_group) {	
116	inm->inm_timer =	
117	<pre>IGMP_RANDOM_DELAY(inm-&gt;inm_addr);</pre>	
118	igmp_timers_are_running = 1;	
119	• }	
120	<pre>IN_NEXT_MULTI(step, inm);</pre>	
121	}	
122	break;	

Figure 13.17 Input processing of the IGMP query message.

97-122 Queries that arrive on the loopback interface are silently discarded (Exercise 13.1). Queries by definition are sent to the all-hosts group. If a query arrives addressed to a different address, it is counted in igps\_rcv\_badqueries and discarded.

The receipt of a query message does not trigger an immediate flurry of IGMP membership reports. Instead, igmp\_input resets the membership timers for each group associated with the interface on which the query was received to a random value with IGMP\_RANDOM\_DELAY. When the timer for a group expires, igmp\_fasttimo sends a membership report. Meanwhile, the same activity is occurring on all the other hosts that received the IGMP query. As soon as the random timer for a particular group expires on one host, it is multicast to that group. This report cancels the timers on the other hosts so that only one report is multicast to the network. The routers, as well as any other members of the group, receive the report.

The one exception to this scenario is the all-hosts group. A timer is never set for this group and a report is never sent.

## Membership Report: IGMP\_HOST\_MEMBERSHIP\_REPORT

The receipt of an IGMP membership report is one of the two events we mentioned in Section 13.1 that does not result in an IGMP message. The effect of the message is local to the interface on which it was received. Figure 13.18 shows the message processing.

123	case IGMP_HOST_MEMBERSHIP_REPORT:	- igm
124	++igmpstat.igps_rcv_reports;	
125	if (ifp == &loif)	
126	break;	
127	if (!IN_MULTICAST(ntohl(igmp->igmp_group.s_addr))	
128	igmp->igmp_group.s_addr != ip->ip_dst.s_addr) {	
129	++igmpstat.igps_rcv_badreports;	
130	<pre>m_freem(m);</pre>	
131	return;	
132	}	
133	/*	
134	* KLUDGE: if the IP source address of the report has an	
135	* unspecified (i.e., zero) subnet number, as is allowed for	
136	* a booting host, replace it with the correct subnet number	
137	* so that a process-level multicast routing demon can	
L38	* determine which subnet it arrived from. This is necessary	
139	* to compensate for the lack of any way for a process to	
140	* determine the arrival interface of an incoming packet.	
141	*/	
142	if ((ntohl(ip->ip_src.s_addr) & IN_CLASSA_NET) == 0) {	
143	IFP_TO_IA(ifp, ia);	
L44	if (ia)	
L45	ip->ip_src.s_addr = htonl(ia->ia_subnet);	
.46	}	
47	/*	
148	* If we belong to the group being reported, stop	
149	* our timer for that group.	
150	*/	
151	IN_LOOKUP_MULTI(igmp->igmp_group, ifp, inm);	
.52	if (inm != NULL) {	
153	<pre>inm-&gt;inm_timer = 0;</pre>	
154	++igmpstat.igps_rcv_ourreports;	
155	}	
156	break;	- igm

Figure 13.18 Input processing of the IGMP report message.

いたので、「ないない」を見ていたが、「ないない」を見ていた。

Section 13.8

## Leaving a Group: igmp\_leavegroup Function 395

123-156

13

зs

is

n

**i**l

.с

Reports sent to the loopback interface are discarded, as are membership reports sent to the incorrect multicast group. That is, the message must be addressed to the group The course and the message.

The source address of an incompletely initialized host might not include a network or host number (or both). igmp\_report looks at the class A network portion of the address, which can only be 0 when the network and subnet portions of the address are 0. If this is the case, the source address is set to the subnet address, which includes the network ID and subnet ID, of the receiving interface. The only reason for doing this is to inform a process-level daemon of the receiving interface, which is identified by the subnet number.

If the receiving interface belongs to the group being reported, the associated report timer is reset to 0. In this way the first report sent to the group stops any other hosts from issuing a report. It is only necessary for the router to know that at least one interface on the network is a member of the group. The router does not need to maintain an explicit membership list or even a counter.

## 13.8 Leaving a Group: igmp\_leavegroup Function

We saw in Chapter 12 that in\_delmulti calls igmp\_leavegroup when the last reference count in the associated in\_multi structure drops to 0.

```
179 void
180 igmp_leavegroup(inm)
181 struct in_multi *inm;
182 {
183 /*
184 * No action required on leaving a group.
185 */
186 }
```

Figure 13.19 igmp\_leavegroup function.

179–186

As we can see, IGMP takes no action when an interface leaves a group. No explicit notification is sent—the next time a multicast router issues an IGMP query, the interface does not generate an IGMP report for this group. If no report is generated for a group, the multicast router assumes that all the interfaces have left the group and stops forwarding multicast packets for the group to the network.

If the interface leaves the group while a report is pending (i.e., the group's report timer is running), the report is never sent, since the timer is discarded by in\_delmulti (Figure 12.36) along with the in\_multi structure for the group when icmp\_leavegroup returns.

igmp.c

igmp.c

## 13.9 Summary

In this chapter we described IGMP, which communicates IP multicast membership information between hosts and routers on a single network. IGMP membership reports are generated when an interface joins a group, and on demand when multicast routers issue an IGMP report query message.

The design of IGMP minimizes the number of messages required to communicate membership information:

- Hosts announce their membership when they join a group.
- Response to membership queries are delayed for a random interval, and the first response suppresses any others.
- Hosts are silent when they leave a group.
- Membership queries are sent no more than once per minute.

Multicast routers share the IGMP information they collect with each other (Chapter 14) to route multicast datagrams toward remote members of the multicast destination group.

#### Exercises

- 13.1 Why isn't it necessary to respond to an IGMP query on the loopback interface?
- 13.2 Verify the assumption stated on lines 226 to 229 in Figure 13.15.
- **13.3** Is it necessary to set random delays for membership queries that arrive on a point-to-point network interface?

14.1

14

## first

r 13

ship orts lters

cate

## r 14)

ation

-point

#### 14.1 Introduction

The previous two chapters discussed multicasting on a single network. In this chapter we look at multicasting across an entire internet. We describe the operation of the mrouted program, which computes the multicast routing tables, and the kernel functions that forward multicast datagrams between networks.

**IP Multicast Routing** 

Technically, multicast packets are forwarded. In this chapter we assume that every multicast packet contains an entire datagram (i.e., there are no fragments), so we use the term datagram exclusively. Net/3 forwards IP fragments as well as IP datagrams.

Figure 14.1 shows several versions of mrouted and how they correspond to the BSD releases. The mrouted releases include both the user-level daemons and the kernel-level multicast code.

mrouted version	Description
1.2	modifies the 4.3BSD Tahoe release
2.0	included with 4.4BSD and Net/3
3.3	modifies SunOS 4.1.3

Figure 14.1 mrouted and IP multicasting releases.

IP multicast technology is an active area of research and development. This chapter discusses version 2.0 of the multicast software, which is included in Net/3 but is considered an obsolete implementation. Version 3.3 was released too late to be discussed fully in this text, but we will point out various 3.3 features along the way.

397

£

S

1

Because commercial multicast routers are not widely deployed, multicast networks are often constructed using multicast *tunnels*, which connect two multicast routers over a standard IP unicast internet. Multicast tunnels are supported by Net/3 and are constructed with the Loose Source Record Route (LSRR) option (Section 9.6). An improved tunneling technique encapsulates the IP multicast datagram within an IP unicast datagram and is supported by version 3.3 of the multicast code but is not supported by Net/3.

As in Chapter 12, we use the generic term *transport protocols* to refer to the protocols that send and receive multicast datagrams, but UDP is the only Internet protocol that supports multicasting.

## 14.2 Code Introduction

The three files listed in Figure 14.2 are discussed in this chapter.

File	Description
netinet/ip_mroute.h	multicast structure definitions
netinet/ip_mroute.c	multicast routing functions
netinet/raw_ip.c	multicast routing options

Figure 14.2 Files discussed in this chapter.

#### **Global Variables**

The global variables used by the multicast routing code are shown in Figure 14.3.

Variable	Datatype	Description
cached_mrt cached_origin cached_originmask mrtstat mrttable numvifs viftable	<pre>struct mrt u_long struct mrtstat struct mrt *[] vifi_t struct vif[]</pre>	one-behind cache for multicast routing multicast group for one-behind cache mask for multicast group for one-behind cache multicast routing statistics hash table of pointers to multicast routes number of enabled multicast interfaces array of virtual multicast interfaces

Figure 14.3 Global variables introduced in this chapter.

#### Statistics

All the statistics collected by the multicast routing code are found in the mrtstat structure described by Figure 14.4. Figure 14.5 shows some sample output of these statistics, from the netstat -gs command.

Section 14.3

mrtstat member	Description	Used by SNMP
mrts_mrt_lookups	#multicast route lookups	
mrts_mrt_misses	#multicast route cache misses	
mrts_grp_lookups	#group address lookups	
mrts_grp_misses	#group address cache misses	
mrts_no_route	#multicast route lookup failures	
mrts_bad_tunnel	#packets with malformed tunnel options	
mrts_cant_tunnel	#packets with no room for tunnel options	

Figure 14.4 Statistics collected in this chapter.

netstat -gs output	mrtstat members
multicast routing:	
329569328 multicast route lookups	mrts_mrt_lookups
9377023 multicast route cache misses	mrts_mrt_misses
242754062 group address lookups	mrts_grp_lookups
159317788 group address cache misses	mrts_grp_misses
65648 datagrams with no route for origin	mrts_no_route
0 datagrams with malformed tunnel options	mrts_bad_tunnel
0 datagrams with no room for tunnel options	mrts_cant_tunnel

Figure 14.5 Sample IP multicast routing statistics.

These statistics are from a system with two physical interfaces and one tunnel interface. These statistics show that the multicast route is found in the cache 98% of the time. The group address cache is less effective with only a 34% hit rate. The route cache is described with Figure 14.34 and the group address cache with Figure 14.21.

#### SNMP Variables

There is no standard SNMP MIB for multicast routing, but [McCloghrie and Farinacci 1994a] and [McCloghrie and Farinacci 1994b] describe some experimental MIBs for multicast routers.

## 14.3 Multicast Output Processing Revisited

In Section 12.15 we described how an interface is selected for an outgoing multicast datagram. We saw that ip\_output is passed an explicit interface in the ip\_moptions structure, or ip\_output looks up the destination group in the routing tables and uses the interface returned in the route entry.

If, after selecting an outgoing interface, ip\_output loops back the datagram, it is queued for input processing on the interface selected for *output* and is considered for forwarding when it is processed by ipintr. Figure 14.6 illustrates this process.

vorks over conroved dataed by

tocols >l that



e

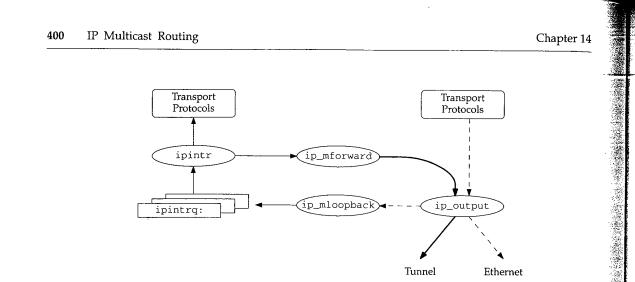


Figure 14.6 Multicast output processing with loopback.

In Figure 14.6 the dashed arrows represent the original outgoing datagram, which in this example is multicast on a local Ethernet. The copy created by ip\_mloopback is represented by the thin arrows; this copy is passed to the transport protocols for input. The third copy is created when ip\_mforward decides to forward the datagram through another interface on the system. The thickest arrows in Figure 14.6 represents the third copy, which in this example is sent on a multicast tunnel.

If the datagram is *not* looped back, ip\_output passes it directly to ip\_mforward, where it is duplicated and also processed as if it were received on the interface that ip\_output selected. This process is shown in Figure 14.7.

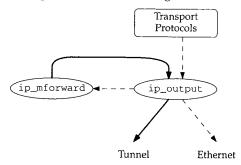


Figure 14.7 Multicast output processing with no loopback.

Whenever ip\_mforward calls ip\_output to send a multicast datagram, it sets the IP\_FORWARDING flag so that ip\_output does not pass the datagram back to ip\_mforward, which would create an infinite loop.

ip\_mloopback was described with Figure 12.42. ip\_mforward is described in Section 14.8.

mrouted Daemon 401

Section 14.4

## 14.4 mrouted Daemon

Multicast routing is enabled and managed by a user-level process: the mrouted daemon. mrouted implements the router portion of the IGMP protocol and communicates with other multicast routers to implement multicast routing between networks. The routing algorithms are implemented in mrouted, but the multicast routing tables are maintained in the kernel, which forwards the datagrams.

In this text we describe only the kernel data structures and functions that support mrouted—we do not describe mrouted itself. We describe the Truncated Reverse Path Broadcast (TRPB) algorithm [Deering and Cheriton 1990], used to select routes for multicast datagrams, and the Distance Vector Multicast Routing Protocol (DVMRP), used to convey information between multicast routers, in enough detail to make sense of the kernel multicast code.

RFC 1075 [Waitzman, Partridge, and Deering 1988] describes an old version of DVMRP. mrouted implements a newer version of DVMRP, which is not yet documented in an RFC. The best documentation for the current algorithm and protocol is the source code release for mrouted. Appendix B describes where the source code can be obtained.

The mrouted daemon communicates with the kernel by setting options on an IGMP socket (Chapter 32). The options are summarized in Figure 14.8.

optname	optval type	Function	Description
DVMRP_INIT		ip_mrouter_init	mrouted is starting
DVMRP_DONE		ip_mrouter_done	mrouted is shutting down
DVMRP_ADD_VIF	struct vifctl	add_vif	add virtual interface
DVMRP_DEL_VIF	vifi_t	del_vif	delete virtual interface
DVMRP ADD LGRP	struct lgrplctl	add_lgrp	add multicast group entry for an interface
DVMRP DEL LGRP	struct lgrplctl	del_lgrp	delete multicast group entry for an interface
DVMRP_ADD_MRT	struct mrtctl	add mrt	add multicast route
 DVMRP_DEL_MRT	struct in_addr	del_mrt	delete multicast route

Figure 14.8 Multicast routing socket options.

The socket options shown in Figure 14.8 are passed to rip\_ctloutput (Section 32.8) by the setsockopt system call. Figure 14.9 shows the portion of rip\_ctloutput that handles the DVMRP\_xxx options.

When setsockopt is called, op equals PRCO\_SETOPT and all the options are passed to the ip\_mrouter\_cmd function. For the getsockopt system call, op equals PRCO\_GETOPT and EINVAL is returned for all the options.

Figure 14.10 shows the ip\_mrouter\_cmd function.

These "options" are more like commands, since they cause the kernel to update various data structures. We use the term *command* throughout the rest of this chapter to emphasize this fact.

173-187

402 IP Multicast Routing

Chapter 14

19.20

は 読みを見ていた たちの ちちん からの なる

1400 (A)

ip\_mroute.c

173	case DVMRP INIT:	raw_ip.c
174	case DVMRP_DONE:	
175	case DVMRP_ADD_VIF:	
176	case DVMRP_DEL_VIF:	
177	case DVMRP_ADD_LGRP:	
178	case DVMRP_DEL_LGRP:	
179	case DVMRP_ADD_MRT:	
180	case DVMRP_DEL_MRT:	
181	if $(op == PRCO\_SETOPT)$ {	
182	error = ip_mrouter_cmd(optname, so, *m);	
183	if (*m)	
184	<pre>(void) m_free(*m);</pre>	
185	} else	
186	error = EINVAL;	
187	return (error);	

Figure 14.9 rip\_ctloutput function: DVMRP\_xxx socket options.

84 int 85 ip\_mrouter\_cmd(cmd, so, m) 86 int cmd; 87 struct socket \*so; 88 struct mbuf \*m; 89 { 90 error = 0;int 91 if (cmd != DVMRP\_INIT && so != ip\_mrouter) 92 error = EACCES; 93 else 94 switch (cmd) { 95 case DVMRP\_INIT: 96 error = ip\_mrouter\_init(so); 97 break; 98 case DVMRP\_DONE: 99 error = ip\_mrouter\_done(); 100 break; 101 case DVMRP\_ADD\_VIF: if (m == NULL || m->m\_len < sizeof(struct vifctl))</pre> 102 103 error = EINVAL; 104 else 105 error = add\_vif(mtod(m, struct vifctl \*)); 106 break; 107 case DVMRP\_DEL\_VIF: 108 if (m == NULL || m->m\_len < sizeof(short))</pre> 109 error = EINVAL; 110 else 111 error = del\_vif(mtod(m, vifi\_t \*)); 112 break;

**INTEL Ex.1013.428** 

Section 14.4

mrouted Daemon

403

```
113
              case DVMRP_ADD_LGRP:
 114
                  if (m == NULL |! m->m_len < sizeof(struct lgrplctl))</pre>
 115
                               error = EINVAL;
 116
                  else
 117
                      error = add_lgrp(mtod(m, struct lgrplctl *));
 118
                  break;
 119
             case DVMRP_DEL_LGRP:
 120
                  if (m == NULL || m->m_len < sizeof(struct lgrplctl))
 121
                               error = EINVAL;
122
                  else
123
                      error = del_lgrp(mtod(m, struct lgrplctl *));
124
                 break;
125
             case DVMRP_ADD_MRT:
126
                 if (m == NULL || m->m_len < sizeof(struct mrtctl))</pre>
127
                              error = EINVAL;
128
                 else
129
                     error = add_mrt(mtod(m, struct mrtctl *));
130
                 break;
131
             case DVMRP_DEL_MRT:
132
                 if (m == NULL || m->m_len < sizeof(struct in_addr))</pre>
133
                              error = EINVAL;
134
                 else
135
                     error = del_mrt(mtod(m, struct in_addr *));
136
                 break;
137
             default:
138
                 error = EOPNOTSUPP;
139
                 break:
140
             }
141
        return (error);
142 }
```

ip\_mroute.c

Figure 14.10 ip\_mrouter\_cmd function.

The first command issued by mrouted must be DVMRP\_INIT. Subsequent com-84-92 mands must come from the same socket as the DVMRP\_INIT command. EACCES is returned when other commands are issued on a different socket. 94-142

Each case in the switch checks to see if the right amount of data was included with the command and then calls the matching function. If the command is not recognized, EOPNOTSUPP is returned. Any error returned from the matching function is posted in error and returned at the end of the function.

Figure 14.11 shows ip\_mrouter\_init, which is called when mrouted issues the DVMRP\_INIT command during initialization.

If the command is issued on something other than a raw IGMP socket, or if 146-157 DVMRP\_INIT has already been set, EOPNOTSUPP or EADDRINUSE are returned respectively. A pointer to the socket on which the initialization command is issued is saved in the global ip\_mrouter. Subsequent commands must be issued on this socket. This prevents the concurrent operation of more than one instance of mrouted.

404 IP Multicast Routing

Chapter 14

いいのない 日本ないのでからう

```
ip_mroute.c
146 static int
147 ip_mrouter_init(so)
148 struct socket *so;
149
   {
        if (so->so_type != SOCK_RAW | {
150
            so->so_proto->pr_protocol != IPPROTO_IGMP)
151
152
            return (EOPNOTSUPP);
        if (ip_mrouter != NULL)
153
154
            return (EADDRINUSE);
155
        ip_mrouter = so;
156
        return (0);
157 }
                                                                            ip mroute.c
```

Figure 14.11 ip\_mrouter\_init function: DVMRP\_INIT command.

The remainder of the DVMRP\_*xxx* commands are described in the following sections.

## 14.5 Virtual Interfaces

When operating as a multicast router, Net/3 accepts incoming multicast datagrams, duplicates them and forwards the copies through one or more interfaces. In this way, the datagram is forwarded to other multicast routers on the internet.

An outgoing interface can be a physical interface or it can be a multicast *tunnel*. Each end of the multicast tunnel is associated with a physical interface on a multicast router. Multicast tunnels allow two multicast routers to exchange multicast datagrams even when they are separated by routers that cannot forward multicast datagrams. Figure 14.12 shows two multicast routers connected by a multicast tunnel.

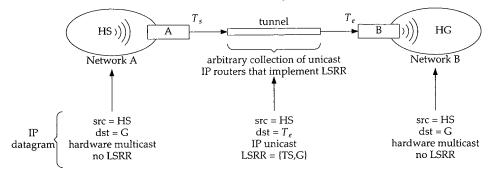


Figure 14.12 A multicast tunnel.

In Figure 14.12, the source host HS on network A is multicasting a datagram to group G. The only member of group G is on network B, which is connected to network A by a multicast tunnel. Router A receives the multicast (because multicast routers receive *all* 

multicasts), consults its multicast routing tables, and forwards the datagram through the multicast tunnel.

The tunnel starts on the *physical* interface on router A identified by the IP unicast address  $T_s$ . The tunnel ends on the *physical* interface on router B identified by the IP unicast address,  $T_e$ . The tunnel itself is an arbitrarily complex collection of networks connected by IP unicast routers that implement the LSRR option. Figure 14.13 shows how an IP LSRR option implements the multicast tunnel.

System	IP header		Source route option		
	ip_src	ip_dst	offset	addresses	Description
$HS \\ T_s \\ T_e \\ T_e$	HS HS HS HS	G T <sub>e</sub> G G	8 12	$T_s \bullet G$ $T_s$ see text $\bullet$	on network A on tunnel after ip_dooptions on router B after ip_mforward on router B

Figure 14.13 LSRR multicast tunnel options.

The first line of Figure 14.13 shows the datagram sent by HS as a multicast on network A. Router A receives the datagram because multicast routers receive all multicasts on their locally attached networks.

To send the datagram through the tunnel, router A inserts an LSRR option in the IP header. The second line shows the datagram as it leaves A on the tunnel. The first address in the LSRR option is the source address of the tunnel and the second address is the destination group. The destination of the datagram is  $T_e$ —the other end of the tunnel. The LSRR offset points to the *destination group*.

The tunneled datagram is forwarded through the internet until it reaches the other end of the tunnel on router B.

The third line of the figure shows the datagram after it is processed by  $ip\_dooptions$  on router B. Recall from Chapter 9 that  $ip\_dooptions$  processes the LSRR option before the destination address of the datagram is examined by ipintr. Since the destination address of the datagram  $(T_e)$  matches one of the interfaces on router B,  $ip\_dooptions$  copies the address identified by the option offset (G in this example) into the destination field of the IP header. In the option, G is replaced with the address returned by  $ip\_rtaddr$ , which normally selects the outgoing interface for the datagram based on the IP destination address (G in this case). This address is irrelevant, since  $ip\_mforward$  discards the entire option. Finally,  $ip\_dooptions$  advances the option offset.

The fourth line in Figure 14.13 shows the datagram after ipintr calls ip\_mforward, where the LSRR option is recognized and removed from the datagram header. The resulting datagram looks like the original multicast datagram and is processed by ip\_mforward, which in our example forwards it onto network B as a multicast datagram where it is received by HG.

Multicast tunnels constructed with LSRR options are obsolete. Since the March 1993 release of mrouted, tunnels have been constructed by prepending another IP header to the IP multicast datagram. The protocol in the new IP header is set to 4 to indicate that the contents of the packet is another IP packet. This value is documented

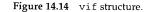
Contraction of the second s

in RFC 1700 as the "IP in IP" protocol. LSRR tunnels are supported in newer versions of mrouted for backward compatibility.

#### Virtual Interface Table

For both physical interfaces and tunnel interfaces, the kernel maintains an entry in a *virtual interface* table, which contains information that is used only for multicasting. Each virtual interface is described by a vif structure (Figure 14.14). The global variable viftable is an array of these structures. An index to the table is stored in a vifi\_t variable, which is an unsigned short integer.

105	struct vif	{		ip_mroute.h
106		v flags;	/*	VIFF flags */
107	u_char	v_threshold;		min ttl required to forward on vif */
108	struct	in_addr v_lcl_addr;		local interface address */
109	struct	in_addr v_rmt_addr;	/*	remote address (tunnels only) */
110	struct	ifnet *v_ifp;	/*	pointer to interface */
111	struct	in_addr *v_lcl_grps;	/*	list of local grps (phyints only) */
112	int	v_lcl_grps_max;	/*	malloc'ed number of v_lcl_grps */
113	int	v_lcl_grps_n;	/*	used number of v_lcl_grps */
114	u_long	v_cached_group;	/*	last grp looked-up (phyints only) */
115	int	v_cached_result;	/*	last look-up result (phyints only) */
116	};			
				ip_mroute.h



105-110 The only flag defined for v\_flags is VIFF\_TUNNEL. When set, the interface is a tunnel to a remote multicast router. When not set, the interface is a physical interface on the local system. v\_threshold is the multicast threshold, which we described in Section 12.9. v\_lcl\_addr is the unicast IP address of the local interface associated with this virtual interface. v\_rmt\_addr is the unicast IP address of the remote end of an IP multicast tunnel. Either v\_lcl\_addr or v\_rmt\_addr is nonzero, but never both. For physical interfaces, v\_ifp is nonnull and points to the ifnet structure of the local interface. For tunnels, v\_ifp is null.

The list of groups with members on the attached interface is kept as an array of IP multicast group addresses pointed to by v\_lcl\_grps, which is always null for tunnels. The size of the array is in v\_lcl\_grps\_max, and the number of entries that are used is in v\_lcl\_grps\_n. The array grows as needed to accommodate the group membership list. v\_cached\_group and v\_cached\_result implement a one-entry cache, which contain the group and result of the previous lookup.

Figure 14.15 illustrates the viftable, which has 32 (MAXVIFS) entries. viftable[2] is the last entry in use, so numvifs is 3. The size of the table is fixed when the kernel is compiled. Several members of the vif structure in the first entry of the table are shown. v\_ifp points to an ifnet structure, v\_lcl\_grps points to an array of in\_addr structures. The array has 32 (v\_lcl\_grps\_max) entries, of which only 4 (v\_lcl\_grps\_n) are in use.

Section 14.5

.4

۱S

а

g.

le \_t

:h

:h

а

n cth

[P

or

al

IP

ls.

is

ip 2h

۶s.

эd

of

an

ch

Virtual Interfaces 407

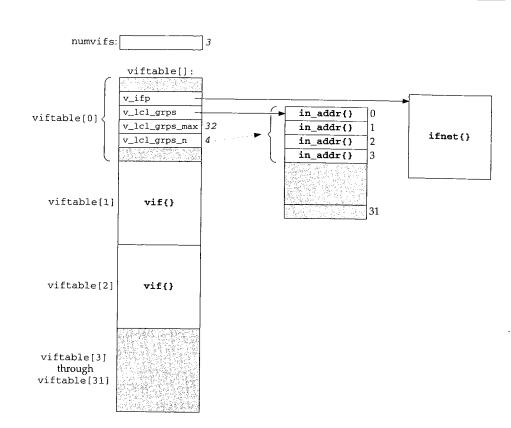


Figure 14.15 viftable array.

mrouted maintains viftable through the DVMRP\_ADD\_VIF and DVMRP\_DEL\_VIF commands. Normally all multicast-capable interfaces on the local system are added to the table when mrouted begins. Multicast tunnels are added when mrouted reads its configuration file, usually /etc/mrouted.conf. Commands in this file can also delete physical interfaces from the virtual interface table or change the multicast information associated with the interfaces.

A vifctl structure (Figure 14.16) is passed by mrouted to the kernel with the DVMRP\_ADD\_VIF command. It instructs the kernel to add an interface to the table of virtual interfaces.

ip\_mroute.h 76 struct vifctl { 77 vifi\_t vifc\_vifi; /\* the index of the vif to be added \*/ 78 u\_char vifc\_flags; /\* VIFF\_ flags (Figure 14.14) \*/ 79 u\_char vifc\_threshold; /\* min ttl required to forward on vif \*/ 80 struct in\_addr vifc\_lcl\_addr; /\* local interface address \*/ 81 struct in\_addr vifc\_rmt\_addr; /\* remote address (tunnels only) \*/ 82 };

#### Figure 14.16 vifctl structure.

76-82 vifc\_vifi identifies the index of the virtual interface within viftable. The remaining four members, vifc\_flags, vifc\_threshold, vifc\_lcl\_addr, and vifc\_rmt\_addr, are copied into the vif structure by the add\_vif function.

## add\_vif Function

Figure 14.17 shows the add\_vif function.

```
ip_mroute.c
 202 static int
 203 add_vif(vifcp)
 204 struct vifctl *vifcp;
205 {
206
         struct vif *vifp = viftable + vifcp->vifc_vifi;
207
         struct ifaddr *ifa;
208
         struct ifnet *ifp;
209
         struct ifreq ifr;
210
         int
                 error, s;
211
         static struct sockaddr_in sin =
212
         {sizeof(sin), AF_INET};
213
         if (vifcp->vifc_vifi >= MAXVIFS)
214
             return (EINVAL);
215
         if (vifp->v_lcl_addr.s_addr != 0)
216
             return (EADDRINUSE);
         /* Find the interface with an address in AF_INET family */
217
218
        sin.sin_addr = vifcp->vifc_lcl_addr;
219
        ifa = ifa_ifwithaddr((struct sockaddr *) &sin);
220
        if (ifa == 0)
221
             return (EADDRNOTAVAIL);
222
        s = splnet();
        if (vifcp->vifc_flags & VIFF_TUNNEL)
223
224
            vifp->v_rmt_addr = vifcp->vifc_rmt_addr;
225
        else {
226
             /* Make sure the interface supports multicast */
227
            ifp = ifa->ifa_ifp;
228
            if ((ifp->if_flags & IFF_MULTICAST) == 0) {
229
                 splx(s);
230
                 return (EOPNOTSUPP);
231
            }
232
            /*
             * Enable promiscuous reception of all IP multicasts
233
234
             * from the interface.
235
             */
            satosin(&ifr.ifr_addr)->sin_family = AF_INET;
236
237
            satosin(&ifr.ifr_addr)->sin_addr.s_addr = INADDR_ANY;
238
            error = (*ifp->if_ioctl) (ifp, SIOCADDMULTI, (caddr_t) & ifr);
239
            if (error) {
240
                splx(s);
241
                return (error);
242
            }
243
        }
```

**NUMBER** 

Section 14.5

- ip\_mroute.c

## 09

244 vifp->v\_flags = vifcp->vifc\_flags; vifp->v\_threshold = vifcp->vifc\_threshold; 245 246 vifp->v\_lcl\_addr = vifcp->vifc\_lcl\_addr; 247 vifp->v\_ifp = ifa->ifa\_ifp; 248 /\* Adjust numvifs up if the vifi is higher than numvifs \*/ 249 if (numvifs <= vifcp->vifc\_vifi) 250 numvifs = vifcp->vifc\_vifi + 1; 251 splx(s); return (0); 252 253 }

Figure 14.17 add\_vif function: DVMRP\_ADD\_VIF command.

#### Validate index

202-216 If the table index specified by mrouted in vifc\_vifi is too large, or the table entry is already in use, EINVAL or EADDRINUSE is returned respectively.

#### Locate physical interface

217-221 ifa\_ifwithaddr takes the unicast IP address in vifc\_lcl\_addr and returns a pointer to the associated ifnet structure. This identifies the physical interface to be used for this virtual interface. If there is no matching interface, EADDRNOTAVAIL is returned.

#### Configure tunnel interface

For a tunnel, the remote end of the tunnel is copied from the vifctl structure to the vif structure in the interface table.

#### Configure physical interface

225-243 For a physical interface, the link-level driver must support multicasting. The SIOCADDMULTI command used with INADDR\_ANY configures the interface to begin receiving all IP multicast datagrams (Figure 12.32) because it is a multicast router. Incoming datagrams are forwarded when ipintr passes them to ip\_mforward.

#### Save multicast information

244-253 The remaining interface information is copied from the vifctl structure to the vif structure. If necessary, numvifs is updated to record the number of virtual interfaces in use.

#### del\_vif Function

The function del\_vif, shown in Figure 14.18, deletes entries from the virtual interface table. It is called when mrouted sets the DVMRP\_DEL\_VIF command.

#### Validate index

If the index passed to del\_vif is greater than the largest index in use or it references an entry that is not in use, EINVAL or EADDRNOTAVAIL is returned respectively.

14

he

nd

Chapter 14

「東京市大学校の ほうやう

Section \_\_\_\_\_

#### ip\_mroute.c 257 static int 258 del\_vif(vifip) 259 vifi\_t \*vifip; 260 { struct vif \*vifp = viftable + \*vifip; 261 262 struct ifnet \*ifp; i, s; 263 int 264 struct ifreq ifr; if (\*vifip >= numvifs) 265 return (EINVAL); 266 267 if (vifp->v\_lcl\_addr.s\_addr == 0) return (EADDRNOTAVAIL); 268 269 s = splnet(); 270 if (!(vifp->v\_flags & VIFF\_TUNNEL)) { 271 if (vifp->v\_lcl\_grps) 272 free(vifp->v\_lcl\_grps, M\_MRTABLE); 273 satosin(&ifr.ifr\_addr)->sin\_family = AF\_INET; satosin(&ifr.ifr\_addr)->sin\_addr.s\_addr = INADDR\_ANY; 274 275 ifp = vifp->v\_ifp; 276 (\*ifp->if\_ioctl) (ifp, SIOCDELMULTI, (caddr\_t) & ifr); 277 3 278 bzero((caddr\_t) vifp, sizeof(\*vifp)); 279 /\* Adjust numvifs down \*/ for (i = numvifs - 1; i >= 0; i--) 280 281 if (viftable[i].v\_lcl\_addr.s\_addr != 0) 282 break; numvifs = i + 1; 283 284 splx(s): 285 return (0); 286 } - ip mroute.c

Figure 14.18 del\_vif function: DVMRP\_DEL\_VIF command.

#### **Delete interface**

<sup>269–278</sup> For a physical interface, the local group table is released, and the reception of all multicast datagrams is disabled by SIOCDELMULTI. The entry in viftable is cleared by bzero.

#### Adjust interface count

The for loop searches for the first active entry in the table starting at the largest previously active entry and working back toward the first entry. For unused entries, the s\_addr member of v\_lcl\_addr (an in\_addr structure) is 0. numvifs is updated accordingly and the function returns.

87-9

add\_:

291-36

302-32

IGMP Revisited 411

Section 14.6

## 14.6 IGMP Revisited

Chapter 13 focused on the host part of the IGMP protocol. mrouted implements the router portion of this protocol. For every physical interface, mrouted must keep track of which multicast groups have members on the attached network. mrouted multicasts an IGMP\_HOST\_MEMBERSHIP\_QUERY datagram every 120 seconds and compiles the resulting IGMP\_HOST\_MEMBERSHIP\_REPORT datagrams into a membership array associated with each network. This array is *not* the same as the membership list we described in Chapter 13.

From the information collected, mrouted constructs the multicast routing tables. The list of groups is also used to suppress multicasts to areas of the multicast internet that do not have members of the destination group.

The membership array is maintained only for physical interfaces. Tunnels are point-to-point interfaces to another multicast router, so no group membership information is needed.

We saw in Figure 14.14 that v\_lcl\_grps points to an array of IP multicast groups. mrouted maintains this list with the DVMRP\_ADD\_LGRP and DVMRP\_DEL\_LGRP commands. An lgrplctl (Figure 14.19) structure is passed with both commands.

87 s	truct lgrplctl {	ip_mroute.h
88	vifi_t lgc_vifi;	
89	struct in_addr lgc_gaddr;	
90 }		
		ip mroute.h

#### Figure 14.19 lgrplctl structure.

<sup>87-90</sup> The {interface, group} pair is identified by lgc\_vifi and lgc\_gaddr. The interface index (lgc\_vifi, an unsigned short) identifies a *virtual* interface, not a physical interface.

When an IGMP\_HOST\_MEMBERSHIP\_REPORT datagram is received, the functions shown in Figure 14.20 are called.

#### add\_lgrp Function

mrouted examines the source address of an incoming IGMP report to determine which subnet and therefore which interface the report arrived on. Based on this information, mrouted sets the DVMRP\_ADD\_LGRP command for the interface to update the membership table in the kernel. This information is also fed into the multicast routing algorithm to update the routing tables. Figure 14.21 shows the add\_lgrp function.

#### Validate add request

291-301

If the request identifies an invalid interface, EINVAL is returned. If the interface is not in use or is a tunnel, EADDRNOTAVAIL is returned.

## If needed, expand group array

302-326

If the new group won't fit in the current group array, a new array is allocated. The first time add\_lgrp is called for an interface, an array is allocated to hold 32 groups.

412 IP Multicast Routing

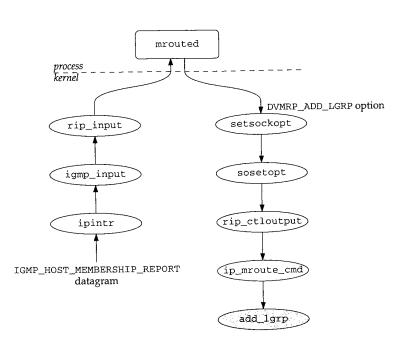


Figure 14.20 IGMP report processing.

Each time the array fills, add\_lgrp allocates a new array of twice the previous size. The new array is allocated by malloc, cleared by bzero, and filled by copying the old array into the new one with bcopy. The maximum number of entries, v\_lcl\_grps\_max, is updated, the old array (if any) is released, and the new array is attached to the vif entry with v\_lcl\_grps.

The "paranoid" comment points out there is no guarantee that the memory allocated by malloc contains all 0s.

#### Add new group

<sup>327–332</sup> The new group is copied into the next available entry and if the cache already contains the new group, the cache is marked as valid.

The lookup cache contains an address, v\_cached\_group, and a cached lookup result, v\_cached\_result. The grplst\_member function always consults the cache before searching the membership array. If the given group matches v\_cached\_group, the cached result is returned; otherwise the membership array is searched.

## del\_lgrp Function

Group information is expired for each interface when no membership report has been received for the group within 270 seconds. mrouted maintains the appropriate timers and issues the DVMRP\_DEL\_LGRP command when the information expires. Figure 14.22 shows del\_lgrp.

#### Chapter 14

ないないないで

「「「「「「「「「「」」」」」

「「「「「「「「「「」」」を見ていていていていた」

Section 14.6

\_\_\_\_\_

- ip\_mroute.c 291 static int 292 add\_lgrp(gcp) 293 struct lgrplctl \*gcp; 294 { 295 struct vif \*vifp; 296 int s; 297 if (gcp->lgc\_vifi >= numvifs) 298 return (EINVAL); 299 vifp = viftable + gcp->lgc\_vifi; 300 if (vifp->v\_lcl\_addr.s\_addr == 0 || (vifp->v\_flags & VIFF\_TUNNEL)) 301 return (EADDRNOTAVAIL); 302 /\* If not enough space in existing list, allocate a larger one \*/ 303 s = splnet(); if (vifp->v\_lcl\_grps\_n + 1 >= vifp->v\_lcl\_grps\_max) { 304 305 int num; 306 struct in\_addr \*ip; 307 num = vifp->v\_lcl\_grps\_max; 308 if (num <= 0)309 num = 32;/\* initial number \*/ 310 else 311 /\* double last number \*/ num += num; 312 ip = (struct in\_addr \*) malloc(num \* sizeof(\*ip), 313 M\_MRTABLE, M\_NOWAIT); 314 if (ip == NULL) { 315 splx(s); 316 return (ENOBUFS); 317 } 318 bzero((caddr\_t) ip, num \* sizeof(\*ip)); /\* XXX paranoid \*/ 319 bcopy((caddr\_t) vifp->v\_lcl\_grps, (caddr\_t) ip, 320 vifp->v\_lcl\_grps\_n \* sizeof(\*ip)); 321 vifp->v\_lcl\_grps\_max = num; 322 if (vifp->v\_lcl\_grps) 323 free(vifp->v\_lcl\_grps, M\_MRTABLE); 324 vifp->v\_lcl\_grps = ip; 325 splx(s); 326 } 327 vifp->v\_lcl\_grps[vifp->v\_lcl\_grps\_n++] = gcp->lgc\_gaddr; 328 if (gcp->lgc\_gaddr.s\_addr == vifp->v\_cached\_group) 329 vifp->v\_cached\_result = 1; 330 splx(s); 331 return (0); 332 } – ip\_mroute.c

 $\label{eq:figure14.21} add\_lgrp\ function:\ process\ {\tt DVMRP\_ADD\_LGRP\ command}.$ 

₂e.

١d

es,

is

by

)n-

.up

:he

up,

een

iers 1.22 414 IP Multicast Routing

Chapter 14

5

```
ip_mroute.c
337 static int
338 del_lgrp(gcp)
339 struct lgrplctl *gcp;
340 {
        struct vif *vifp;
341
342
        int
                i, error, s;
        if (gcp->lgc_vifi >= numvifs)
343
            return (EINVAL);
344
        vifp = viftable + gcp->lgc_vifi;
345
        if (vifp->v_lcl_addr.s_addr == 0 || (vifp->v_flags & VIFF_TUNNEL))
346
            return (EADDRNOTAVAIL);
347
348
        s = splnet();
        if (gcp->lgc_gaddr.s_addr == vifp->v_cached_group)
349
350
            vifp->v_cached_result = 0;
        error = EADDRNOTAVAIL;
351
352
        for (i = 0; i < vifp->v_lcl_grps_n; ++i)
            if (same(&gcp->lgc_gaddr, &vifp->v_lcl_grps[i])) (
353
354
                error = 0;
                vifp->v_lcl_grps_n--;
355
                bcopy((caddr_t) & vifp->v_lcl_grps[i + 1],
35б
                       (caddr_t) & vifp->v_lcl_grps[i],
357
                       (vifp->v_lcl_grps_n - i) * sizeof(struct in_addr));
358
359
                 error = 0;
                break:
360
            }
361
362
        splx(s);
        return (error);
363
364 }
                                                                          ip_mroute.c
```

Figure 14.22 del\_lgrp function: process DVMRP\_DEL\_LGRP command.

#### Validate interface index

337–347 If the request identifies an invalid interface, EINVAL is returned. If the interface is not in use or is a tunnel, EADDRNOTAVAIL is returned.

#### Update lookup cache

348–350 If the group to be deleted is in the cache, the lookup result is set to 0 (false).

#### Delete group

EADDRNOTAVAIL is posted in error in case the group is not found in the membership list. The for loop searches the membership array associated with the interface. If same (a macro that uses bcmp to compare the two addresses) is true, error is cleared and the group count is decremented. bcopy shifts the subsequent array entries down to delete the group and del\_lgrp breaks out of the loop.

If the loop completes without finding a match, EADDRNOTAVAIL is returned; otherwise 0 is returned.

Section 14.6

- ip\_mroute.c

#### grplst\_member Function

During multicast forwarding, the membership array is consulted to avoid sending datagrams on a network when no member of the destination group is present. grplst\_member, shown in Figure 14.23, searches the list looking for the given group address.

368 st	atic int	— ip_mroute.c			
369 grplst_member(vifp, gaddr)					
	ruct vif *vifp;				
371 st:	ruct in_addr gaddr;				
372 {					
373	int i, s;				
374	u_long addr;				
375	<pre>mrtstat.mrts_grp_lookups++;</pre>				
376	addr = gaddr.s_addr;				
377	if (addr == vifp->v_cached_group)				
378	return (vifp->v_cached_result);				
379	<pre>mrtstat.mrts_grp_misses++;</pre>				
380	<pre>for (i = 0; i &lt; vifp-&gt;v_lcl_grps_n; ++i)</pre>				
381	if (addr == vifp->v_lcl_grps(i).s_addr) {				
382	s = splnet();				
383	vifp->v_cached_group = addr;				
384	<pre>vifp-&gt;v_cached_result = 1;</pre>				
385	<pre>splx(s);</pre>				
386	return (1);				
387	}				
388	s = splnet();				
389	vifp->v_cached_group = addr;				
390	vifp->v_cached_result = 0;				
391	splx(s);				
392	return (0);				
393 }					

ace is

mber-

ice. If

eared

down

other-

route.c

Figure 14.23 grplst\_member function.

#### Check the cache

<sup>368–379</sup> If the requested group is located in the cache, the cached result is returned and the membership array is not searched.

## Search the membership array

<sup>380-393</sup> A linear search determines if the group is in the array. If it is found, the cache is updated to record the match and one is returned. If it is not found, the cache is updated to record the miss and 0 is returned.

er 14 -----

ute.c

- 14 (No.

S

## 14.7 Multicast Routing

As we mentioned at the start of this chapter, we will not be presenting the TRPB algorithm implemented by mrouted, but we do need to provide a general overview of the mechanism to describe the multicast routing table and the multicast routing functions in the kernel. Figure 14.24 shows the sample multicast network that we use to illustrate the algorithms.

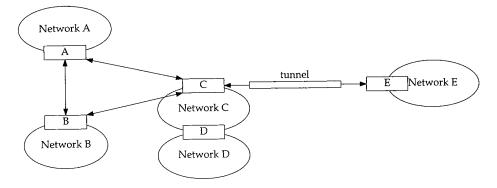


Figure 14.24 Sample multicast network.

In Figure 14.24, routers are shown as boxes and the ellipses are the multicast networks attached to the routers. For example, router D can multicast on network D and C. Router C can multicast to network C, to routers A and B through point-to-point interfaces, and to E through a multicast tunnel.

The simplest approach to multicast routing is to select a subset of the internet topology that forms a *spanning tree*. If each router forwards multicasts along the spanning tree, every router eventually receives the datagram. Figure 14.25 shows one spanning tree for our sample network, where host S on network A represents the source of a multicast datagram.

## For a discussion of spanning trees, see [Tanenbaum 1989] or [Perlman 1992].

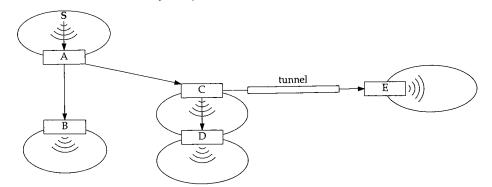


Figure 14.25 Spanning tree for network A.

Multicast Routing 417

Section 14.7

dialgoof the ons in Istrate

pter 14

st net-D and -point

topolinning inning re of a We constructed the tree based on the shortest *reverse path* from every network back to the source in network A. In Figure 14.25, the link between routers B and C is omitted to form the spanning tree. The arrows between the source and router A, and between router C and D, emphasize that the multicast network is part of the spanning tree.

If the same spanning tree were used to forward a datagram from network C, the datagram would be forwarded along a longer path than needed to get to a recipient on network B. The algorithm described in RFC 1075 computes a separate spanning tree for each potential source network to avoid this problem. The routing tables contain a network number and subnet mask for each route, so that a single route applies to any host within the source subnet.

Because each spanning tree is constructed to provide the shortest reverse path to the source of the datagram, and every network receives every multicast datagram, this process is called *reverse path broadcasting* or RPB.

The RPB protocol has no knowledge of multicast group membership, so many datagrams are unnecessarily forwarded to networks that have no members in the destination group. If, in addition to computing the spanning trees, the routing algorithm records which networks are *leaves* and is aware of the group membership on each network, then routers attached to leaf networks can avoid forwarding datagrams onto the network when there there is no member of the destination group present. This is called *truncated reverse path broadcasting* (TRPB), and is implemented by version 2.0 of mrouted with the help of IGMP to keep track of membership in the leaf networks.

Figure 14.26 shows TRPB applied to a multicast sent from a source on network C and with a member of the destination group on network B.

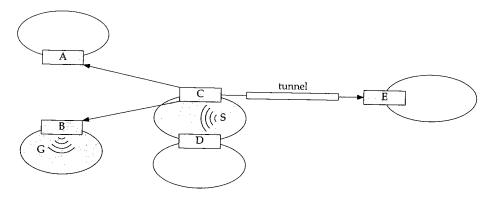


Figure 14.26 TRPB routing for network C.

We'll use Figure 14.26 to illustrate the terms used in the Net/3 multicast routing table. In this example, the shaded networks and routers receive a copy of the multicast datagram sent from the source on network C. The link between A and B is not part of the spanning tree and C does not have a link to D, since the multicast sent by the source is received directly by C and D.

In this figure, networks A, B, D, and E are leaf networks. Router C receives the multicast and forwards it through the interfaces attached to routers A, B, and E—even

00.

though sending it to A and E is wasted effort. This is a major weakness of the TRPB algorithm.

The interface associated with network C on router C is called the *parent* because it is the interface on which router C expects to receive multicasts originating from network C. The interfaces from router C to routers A, B, and E, are *child* interfaces. For router A, the point-to-point interface is the parent for the source packets from C and the interface for network A is a child. Interfaces are identified as a parent or as a child relative to the source of the datagram. Multicast datagrams are forwarded only to the associated child interfaces, and never to the parent interface.

Continuing with the example, networks A, D, and E are not shaded because they are leaf networks without members of the destination group, so the spanning tree is truncated at the routers and the datagram is not forwarded onto these networks. Router B forwards the datagram onto network B, since there is a member of the destination group on the network. To implement the truncation algorithm, each multicast router that receives the datagram consults the group table associated with every virtual interface in the router's viftable.

The final refinement to the multicast routing algorithm is called *reverse path multicasting* (RPM). The goal of RPM is to *prune* each spanning tree and avoid sending datagrams along branches of the tree that do not contain a member of the destination group. In Figure 14.26, RPM would prevent router C from sending a datagram to A and E, since there is no member of the destination group in those branches of the tree. Version 3.3 of mrouted implements RPM.

Figure 14.27 shows our example network, but this time only the routers and networks reached when the datagram is routed by RPM are shaded.

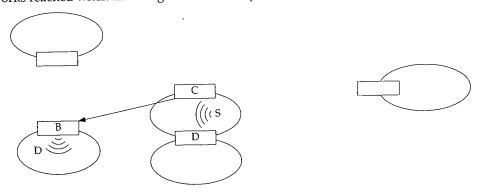


Figure 14.27 RPM routing for network C.

To compute the routing tables corresponding to the spanning trees we described, the multicast routers communicate with adjacent multicast routers to discover the multicast internet topology and the location of multicast group members. In Net/3, DVMRP is used for this communication. DVMRP messages are transmitted as IGMP datagrams and are sent to the multicast group 224.0.0.4, which is reserved for DVMRP communication (Figure 12.1).

In Figure 12.39, we saw that incoming IGMP packets are always accepted by a

Multicast Routing 419

Section 14.7

multicast router. They are passed to igmp\_input, to rip\_input, and then read by mrouted on a raw IGMP socket. mrouted sends DVMRP messages to other multicast routers on the same raw IGMP socket.

For more information about RPB, TRPB, RPM, and the DVMRP messages that are needed to implement these algorithms, see [Deering and Cheriton 1990] and the source code release of mrouted.

There are other multicast routing protocols in use on the Internet. Proteon routers implement the MOSPF protocol described in RFC 1584 [Moy 1994]. PIM (Protocol Independent Multicasting) is implemented by Cisco routers, starting with Release 10.2 of their operating software. PIM is described in [Deering et al. 1994].

## **Multicast Routing Table**

We can now describe the implementation of the multicast routing tables in Net/3. The kernel's multicast routing table is maintained as a hash table with 64 entries (MRTHASHSIZ). The table is kept in the global array mrttable, and each entry points to a linked list of mrt structures, shown in Figure 14.28.

```
- ip_mroute.h
120 struct mrt {
121
        struct in_addr mrt_origin; /* subnet origin of multicasts */
122
        struct in_addr mrt_originmask; /* subnet mask for origin */
123
       vifi_t mrt_parent;
                             /* incoming vif */
124
       vifbitmap_t mrt_children;
                                   /* outgoing children vifs */
125
       vifbitmap_t mrt leaves;
                                   /* subset of outgoing children vifs */
126
       struct mrt *mrt_next;
                                   /* forward link */
127 };
```

## Figure 14.28 mrt structure.

- ip\_mroute.h

120-127

127 mrtc\_origin and mrtc\_originmask identify an entry in the table. mrtc\_parent is the index of the virtual interface on which all multicast datagrams from the origin are expected. The outgoing interfaces are identified within mrtc\_children, which is a bitmap. Outgoing interfaces that are also leaves in the multicast routing tree are identified in mrtc\_leaves, which is also a bitmap. The last member, mrt\_next, implements a linked list in case multiple routes hash to the same array entry.

Figure 14.29 shows the organization of the multicast routing table. Each mrt structure is placed in the hash chain that corresponds to return value from the nethash function shown in Figure 14.31.

The multicast routing table maintained by the kernel is a subset of the routing table maintained within mrouted and contains enough information to support multicast forwarding within the kernel. Updates to the kernel table are sent with the DVMRP\_ADD\_MRT command, which includes the mrtctl structure shown in Figure 14.30.

14

PB

: is

۳k

A, ice

:he

ild

ley

is ks.

naast

ual

ath

ing

ion

ind

ee.

1et-

vy a

and the second second

and the second sec

Chapter 14 420 IP Multicast Routing mrt{} mrt{} mrt{} mrttable[]: mrt\_origin 0 mrt\_originmask 1 mrt\_parent 2 mrt\_children 3 mrt\_leaves 4 mrt next mrt next mrt\_next mrt{} mrt{} mrttable[5] through mrttable[62] mrt, \_next mrt\_next 63 Figure 14.29 Multicast routing table. - ip\_mroute.h 95 struct mrtctl { struct in\_addr mrtc\_origin; /\* subnet origin of multicasts \*/ 96 /\* subnet mask for origin \*/ struct in\_addr mrtc\_originmask; 97 vifi\_t mrtc\_parent; /\* incoming vif \*/ 98 vifbitmap\_t mrtc\_children; /\* outgoing children vifs \*/ 99 /\* subset of outgoing children vifs \*/ vifbitmap\_t mrtc\_leaves; 100 101 }; - ip\_mroute.h Figure 14.30 mrtctl structure. The five members of the mrtctl structure carry the information we have already 95–101 described (Figure 14.28) between mrouted and the kernel. The multicast routing table is keyed by the source IP address of the multicast datagram. nethash (Figure 14.31) implements the hashing algorithm used for the table. It のないので、「ない」のないで、 accepts the source IP address and returns a value between 0 and 63 (MRTHASHSIZ - 1). ip\_mroute.c 398 static u\_long 399 nethash(in) 400 struct in\_addr in; 401 { 402 u\_long n; 403 n = in\_netof(in); 10.00 while ((n & 0xff) == 0)404 n >>= 8; 405 return (MRTHASHMOD(n)); 406

Figure 14.31 nethash function.

407 }

ip\_mroute.c

Section 14.7

<sup>398-407</sup> in\_netof returns in with the host portion set to all 0s leaving only the class A, B, or C network of the sending host in n. The result is shifted to the right until the loworder 8 bits are nonzero. MRTHASHMOD is

#define MRTHASHMOD(h) ((h) & (MRTHASHSIZ - 1))

The low-order 8 bits are logically ANDed with 63, leaving only the low-order 6 bits, which is an integer in the range 0 to 63.

Doing two function calls (nethash and in\_netof) to calculate a hash value is an expensive algorithm to compute a hash for a 32-bit address.

#### del\_mrt Function

.h

?.h

łу

:a-It . :e.c

te.c

The mrouted daemon adds and deletes entries in the kernel's multicast routing table through the DVMRP\_ADD\_MRT and DVMRP\_DEL\_MRT commands. Figure 14.32 shows the del\_mrt function.

451 st	atic int ip_mroute.c
452 de	el_mrt(origin)
	ruct in_addr *origin;
454 {	
455	<pre>struct mrt *rt, *prev_rt;</pre>
456	<pre>u_long hash = nethash(*origin);</pre>
457	int s;
458	<pre>for (prev_rt = rt = mrttable[hash]; rt; prev_rt = rt, rt = rt-&gt;mrt_next)</pre>
459	if (origin->s_addr == rt->mrt_origin.s_addr)
460	break;
461	if (!rt)
462	return (ESRCH);
463	s = splnet();
464	if (rt == cached_mrt)
465	cached_mrt = NULL;
466	if (prev_rt == rt)
467	<pre>mrttable[hash] = rt-&gt;mrt_next;</pre>
468	else
469	<pre>prev_rt-&gt;mrt_next = rt-&gt;mrt_next;</pre>
470	<pre>free(rt, M_MRTABLE);</pre>
471	<pre>splx(s);</pre>
472	return (0);
473 }	

Figure 14.32 del\_mrt function: process DVMRP\_DEL\_MRT command.

## Find route entry

<sup>451–462</sup> The for loop starts at the entry identified by hash (initialized in its declaration from nethash). If the entry is not located, ESRCH is returned.

*ip\_mroute.c* 

**IP Multicast Routing** 422

Chapter 14

語等語

のないないない

の変換

2 

「「「「「「「「」」」」

「「「「「「「「」」」」」

No. No. of Concession, Name

小市市御史の湯を

のないので、中国語の語を見たいである。

#### **Delete route entry**

If the entry was stored in the cache, the cache is invalidated. The entry is unlinked 463-473 from the hash chain and released. The if statement is needed to handle the special case when the matched entry is at the front of the list.

#### add\_mrt Function

The add\_mrt function is shown in Figure 14.33.

```
- ip_mroute.c
411 static int
412 add_mrt(mrtcp)
413 struct mrtctl *mrtcp;
414 {
        struct mrt *rt;
415
        u_long hash;
416
        int
                s;
417
        if (rt = mrtfind(mrtcp->mrtc_origin)) {
418
            /* Just update the route */
419
            s = splnet();
420
            rt->mrt_parent = mrtcp->mrtc_parent;
421
            VIFM_COPY(mrtcp->mrtc_children, rt->mrt_children);
422
            VIFM_COPY(mrtcp->mrtc_leaves, rt->mrt_leaves);
423
            splx(s);
424
             return (0);
425
426
         }
        s = splnet();
427
        rt = (struct mrt *) malloc(sizeof(*rt), M_MRTABLE, M_NOWAIT);
428
         if (rt == NULL) {
429
             splx(s);
430
             return (ENOBUFS);
431
         }
432
 433
         /*
          * insert new entry at head of hash chain
 434
          */
 435
         rt->mrt_origin = mrtcp->mrtc_origin;
 436
         rt->mrt_originmask = mrtcp->mrtc_originmask;
 437
         rt->mrt_parent = mrtcp->mrtc_parent;
 438
         VIFM_COPY(mrtcp->mrtc_children, rt->mrt_children);
 439
         VIFM_COPY(mrtcp->mrtc_leaves, rt->mrt_leaves);
 440
         /* link into table */
 441
         hash = nethash(mrtcp->mrtc_origin);
 442
         rt->mrt_next = mrttable[hash];
 443
         mrttable[hash] = rt;
 444
         splx(s);
 445
         return (0);
 446
 447 }
                                                                           ip_mroute.c
```

Figure 14.33 add\_mrt function: process DVMRP\_ADD\_MRT command.

## Section 14.7

#### Update existing route

411-427 If the requested route is already in the routing table, the new information is copied into the route and add\_mrt returns.

#### Allocate new route

428-447 An mrt structure is constructed in a newly allocated mbuf with the information from mrtctl structure passed with the add request. The hash index is computed from mrtc\_origin, and the new route is inserted as the first entry on the hash chain.

#### mrtfind Function

The multicast routing table is searched with the mrtfind function. The source of the datagram is passed to mrtfind, which returns a pointer to the matching mrt structure, or a null pointer if there is no match.

```
ip_mroute.c
477 static struct mrt *
478 mrtfind(origin)
479 struct in_addr origin;
480 {
481
        struct mrt *rt;
482
        u_int
               hash;
483
        int
                s:
484
        mrtstat.mrts_mrt_lookups++;
485
        if (cached_mrt != NULL &&
486
            (origin.s_addr & cached_originmask) == cached_origin)
487
            return (cached_mrt);
488
        mrtstat.mrts_mrt_misses++;
489
        hash = nethash(origin);
490
        for (rt = mrttable[hash]; rt; rt = rt->mrt_next)
491
            if ((origin.s_addr & rt->mrt_originmask.s_addr) ==
492
                rt->mrt_origin.s_addr) {
493
                s = splnet();
494
                cached_mrt = rt;
495
                cached_origin = rt->mrt_origin.s_addr;
                cached_originmask = rt->mrt_originmask.s_addr;
496
497
                splx(s);
498
                return (rt);
499
            }
500
        return (NULL);
501 }
```

#### Figure 14.34 mrtfind function.

#### Check route lookup cache

<sup>477–488</sup> The given source IP address (origin) is logically ANDed with the origin mask in the cache. If the result matches cached\_origin, the cached entry is returned.

'e.c

oute.c

۶d

se

₩. F

- ip\_mroute.c

Se

#### Check the hash table

<sup>489-501</sup> nethash returns the hash index for the route entry. The for loop searches the hash chain for a matching route. When a match is found, the cache is updated and a pointer to the route is returned. If a match is not found, a null pointer is returned.

# 14.8 Multicast Forwarding: ip\_mforward Function

Multicast forwarding is implemented entirely in the kernel. We saw in Figure 12.39 that ipintr passes incoming multicast datagrams to ip\_mforward when ip\_mrouter is nonnull, that is, when mrouted is running.

We also saw in Figure 12.40 that ip\_output can pass multicast datagrams that originate on the local host to ip\_mforward to be routed to interfaces other than the one interface selected by ip\_output.

Unlike unicast forwarding, each time a multicast datagram is forwarded to an interface, a copy is made. For example, if the local host is acting as a multicast router and is connected to three different networks, multicast datagrams originating on the system are duplicated and queued for *output* on all three interfaces. Additionally, the datagram may be duplicated and queued for *input* if the multicast loopback flag was set by the application or if any of the outgoing interfaces receive their own transmissions.

Figure 14.35 shows a multicast datagram arriving on a physical interface.

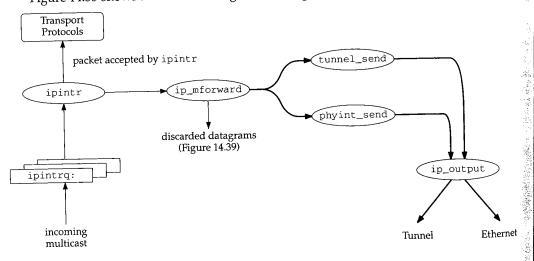


Figure 14.35 Multicast datagram arriving on physical interface.

In Figure 14.35, the interface on which the datagram arrived is a member of the destination group, so the datagram is passed to the transport protocols for input processing. The datagram is also passed to ip\_mforward, where it is duplicated and Section 14.8

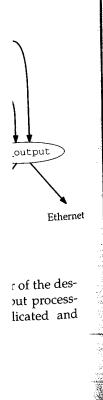
thes the d and a l.

apter 14

2.39 that

ams that a the one

an interer and is e system latagram et by the



forwarded to a physical interface and to a tunnel (the thick arrows), both of which must be different from the receiving interface.

Figure 14.36 shows a multicast datagram arriving on a tunnel.

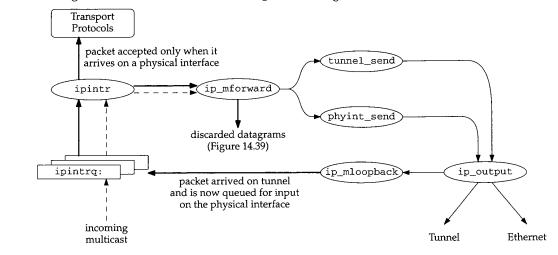


Figure 14.36 Multicast datagram arriving on a multicast tunnel.

In Figure 14.36, the datagram arriving on a physical interface associated with the local end of the tunnel is represented by the dashed arrows. It is passed to ip\_mforward, which as we'll see in Figure 14.37 returns a nonzero value because the packet arrived on a tunnel. This causes ipintr to not pass the packet to the transport protocols.

ip\_mforward strips the tunnel options from the packet, consults the multicast routing table, and, in this example, forwards the packet on another tunnel and on the same *physical* interface on which it arrived, as shown by the thin arrows. This is OK because the multicast routing tables are based on the *virtual* interfaces, not the physical interfaces.

In Figure 14.36 we assume that the physical interface is a member of the destination group, so ip\_output passes the datagram to ip\_mloopback, which queues it for processing by ipintr (the thick arrows). The packet is passed to ip\_mforward again, where it is discarded (Exercise 14.4). ip\_mforward returns 0 this time (because the packet arrived on a physical interface), so ipintr considers and accepts the datagram for input processing.

We show the multicast forwarding code in three parts:

- tunnel input processing (Figure 14.37),
- forwarding eligibility (Figure 14.39), and
- forward to outgoing interfaces (Figure 14.40).

426 IP Multicast Routing

、たちまいたちの主要の構成ない

516 int 517 ip\_mforward(m, ifp) 518 struct mbuf \*m; 519 struct ifnet \*ifp; 520 { struct ip \*ip = mtod(m, struct ip \*); 521 struct mrt \*rt; 522 struct vif \*vifp; 523 524 int vifi; u\_char \*ipoptions; 525 u\_long tunnel\_src; 526 if (ip->ip\_h1 < (IP\_HDR\_LEN + TUNNEL\_LEN) >> 2 | | 527 (ipoptions = (u\_char \*) (ip + 1))[1] != IPOPT\_LSRR) { 528 /\* Packet arrived via a physical interface. \*/ 529 530 tunnel\_src = 0; 531 } else { 532 /\* \* Packet arrived through a tunnel. 533 \* A tunneled packet has a single NOP option and a 534 \* two-element loose-source-and-record-route (LSRR) 535 \* option immediately following the fixed-size part of 536 \* the IP header. At this point in processing, the IP 537 \* header should contain the following IP addresses: 538 539 - in the source address field 540 \* original source \* destination group - in the destination address field 541 \* remote tunnel end-point - in the first element.of LSRR 542 \* one of this host's addrs - in the second element of LSRR 543 544 \* NOTE: RFC-1075 would have the original source and 545 \* remote tunnel end-point addresses swapped. However, 546 \* that could cause delivery of ICMP error messages to 547 \* innocent applications on intermediate routing 548 \* hosts! Therefore, we hereby change the spec. 549 550 \*/ /\* Verify that the tunnel options are well-formed. \*/ 551 if (ipoptions[0] != IPOPT\_NOP || 552 ipoptions[2] != 11 || /\* LSRR option length \*/ 553 ipoptions[3] != 12 || /\* LSRR address pointer \*/ 554 (tunnel\_src = \*(u\_long \*) (&ipoptions[4])) == 0) { 555 mrtstat.mrts\_bad\_tunnel++; 556 return (1); 557 558 } /\* Delete the tunnel options from the packet. \*/ 559 ovbcopy((caddr\_t) (ipoptions + TUNNEL\_LEN), (caddr\_t) ipoptions, 560 (unsigned) (m->m\_len - (IP\_HDR\_LEN + TUNNEL\_LEN))); 561 m->m\_len -= TUNNEL\_LEN; 562 ip->ip\_len -= TUNNEL\_LEN; 563 564 ip->ip\_hl -= TUNNEL\_LEN >> 2; 565 } ip\_mroute.c

Figure 14.37 ip\_mforward function: tunnel arrival.

Chapter 14

ip mroute.c

Section 14.8

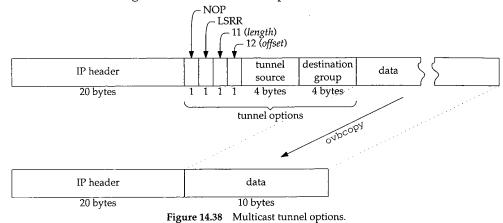
<sup>516–526</sup> The two arguments to ip\_mforward are a pointer to the mbuf chain containing the datagram; and a pointer to the ifnet structure of the receiving interface.

#### Arrival on physical interface

<sup>527-530</sup> To distinguish between a multicast datagram arriving on a physical interface and a tunneled datagram arriving on the same physical interface, the IP header is examined for the characteristic LSRR option. If the header is too small to contain the option, or if the options don't start with a NOP followed by an LSRR option, it is assumed that the datagram arrived on a physical interface and tunnel\_src is set to 0.

#### Arrival on a tunnel

<sup>531–558</sup> If the datagram looks as though it arrived on a tunnel, the options are verified to make sure they are well formed. If the options are not well formed for a multicast tunnel, ip\_mforward returns 1 to indicate that the datagram should be discarded. Figure 14.38 shows the organization of the tunnel options.



In Figure 14.38 we assume there are no other options in the datagram, although that is not required. Any other IP options will appear after the LSRR option, which is always inserted before any other options by the multicast router at the start of the tunnel.

#### **Delete tunnel options**

If the options are OK, they are removed from the datagram by shifting the remaining options and data forward and adjusting m\_len in the mbuf header and ip\_len and ip\_hl in the IP header (Figure 14.38).

ip\_mforward often uses tunnel\_source as its return value, which is only nonzero when the datagram arrives on a tunnel. When ip\_mforward returns a nonzero value, the caller discards the datagram. For ipintr this means that a datagram that arrives on a tunnel is passed to ip\_mforward and discarded by ipintr. The forwarding code strips out the tunnel information, duplicates the datagram, and sends the datagrams with ip\_output, which calls ip\_mloopback if the interface is a member of the destination group.

iroute.c

r 14

**IP Multicast Routing** 

Chapter 14

Sec

580

The next part of ip\_mforward, shown in Figure 14.39, discards the datagram if it is ineligible for forwarding.

0	ip_mroute.c
566	/*
567	* Don't forward a packet with time-to-live of zero or one,
568	* or a packet destined to a local-only group.
569	*/
570	if (ip->ip_ttl <= 1
571	<pre>ntohl(ip-&gt;ip_dst.s_addr) &lt;= INADDR_MAX_LOCAL_GROUP)</pre>
572	return ((int) tunnel_src);
573	/*
574	* Don't forward if we don't have a route for the packet's origin.
575	*/
576	if (!(rt = mrtfind(ip->ip_src))) {
577	<pre>mrtstat.mrts_no_route++;</pre>
578	<pre>return ((int) tunnel_src);</pre>
579	}
580	/*
581	' * Don't forward if it didn't arrive from the parent vif for its origin.
582	*/
583	vifi = rt->mrt_parent;
584	if (tunnel_src == 0) {
585	if ((viftable[vifi].v_flags & VIFF_TUNNEL)
586	<pre>viftable[vifi].v_ifp != ifp)</pre>
587	<pre>return ((int) tunnel_src);</pre>
588	} else {
589	if (!(viftable[vifi].v_flags & VIFF_TUNNEL)
590	viftable[vifi].v_rmt_addr.s_addr != tunnel_src)
591	return ((int) tunnel_src);
592	) ip_mroute.c

Figure 14.39 ip\_mforward function: forwarding eligibility checks.

## **Expired TTL or local multicast**

If ip\_ttl is 0 or 1, the datagram has reached the end of its lifetime and is not for-566-572 warded. If the destination group is less than or equal to INADDR\_MAX\_LOCAL\_GROUP (the 224.0.0.x groups, Figure 12.1), the datagram is not allowed beyond the local network and is not forwarded. In either case, tunnel\_src is returned to the caller.

> Version 3.3 of mrouted supports administrative scoping of certain destination groups. An interface can be configured to discard datagrams addressed to these groups, similar to the automatic scoping of the 224.0.0.x groups.

#### No route available

428

If mrtfind cannot locate a route based on the source address of the datagram, the 573--579 function returns. Without a route, the multicast router cannot determine to which interfaces the datagram should be forwarded. This might occur, for example, when the multicast datagrams arrive before the multicast routing table has been updated by mrouted.

593-

Section 14.8

#### Arrived on unexpected interface

<sup>580–592</sup> If the datagram arrived on a physical interface but was expected to arrive on a tunnel or on a different physical interface, ip\_mforward returns. If the datagram arrived on a tunnel but was expected to arrive on a physical interface or on a different tunnel, ip\_mforward returns. A datagram may arrive on an unexpected interface when the routing tables are in transition because of changes in the group membership or in the physical topology of the network.

The final part of ip\_mforward (Figure 14.40) sends the datagram on each of the outgoing interfaces specified in the multicast route entry.

ip\_mroute.c 593 /\*  $\star$  For each vif, decide if a copy of the packet should be forwarded. 594 595 Forward if: 596 - the ttl exceeds the vif's threshold AND 597 - the vif is a child in the origin's route AND 598 - ( the vif is not a leaf in the origin's route OR 599 the destination group has members on the vif ) 600 601 \* (This might be speeded up with some sort of cache -- someday.) \*/ 602 603 for (vifp = viftable, vifi = 0; vifi < numvifs; vifp++, vifi++) {</pre> if (ip->ip\_ttl > vifp->v\_threshold && 604 605 VIFM\_ISSET(vifi, rt->mrt\_children) && 606 (!VIFM\_ISSET(vifi, rt->mrt\_leaves) || grplst\_member(vifp, ip->ip\_dst))) { 607 608 if (vifp->v\_flags & VIFF\_TUNNEL) 609 tunnel\_send(m, vifp); 610 else 611 phyint\_send(m, vifp); 612 } 613 } 614 return ((int) tunnel\_src); 615 }

Figure 14.40 ip\_mforward function: forwarding.

593-615

2

n

ρ

e

е

y

For each interface in viftable, a datagram is sent on the interface if

the datagram's TTL is greater than the multicast threshold for the interface,

- the interface is a child interface for the route, and
- the interface is not connected to a leaf network.

If the interface is a leaf, the datagram is output only if there is a member of the destination group on the network (i.e., grplst\_member returns a nonzero value).

tunnel\_send forwards the datagram on tunnel interfaces; phyint\_send is used for physical interfaces.

ip mroute.c

## phyint\_send Function

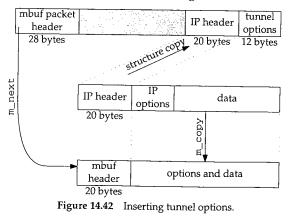
To send a multicast datagram on a physical interface, phyint\_send (Figure 14.41) specifies the output interface explicitly in the ip\_moptions structure it passes to ip\_output.

```
616 static void
617 phyint_send(m, vifp)
618 struct mbuf *m;
619 struct vif *vifp;
620 (
621
        struct ip *ip = mtod(m, struct ip *);
622
        struct mbuf *mb_copy;
623
        struct ip_moptions *imo;
624
        int
                error:
625
        struct ip_moptions simo;
626
        mb_copy = m_copy(m, 0, M_COPYALL);
627
        if (mb_copy == NULL)
628
            return;
629
        imo = &simo;
630
        imo->imo_multicast_ifp = vifp->v_ifp;
        imo->imo_multicast_ttl = ip->ip_ttl - 1;
631
632
        imo->imo_multicast_loop = 1;
633
        error = ip_output(mb_copy, NULL, NULL, IP_FORWARDING, imo);
634 }
```

616-634

34 m\_copy duplicates the outgoing datagram. The ip\_moptions structure is set to force the datagram to be transmitted on the selected interface. The TTL value is decremented, and multicast loopback is enabled.

The datagram is passed to ip\_output. The IP\_FORWARDING flag avoids an infinite loop, where ip\_output calls ip\_mforward again.



Chapter 14

*ip\_mroute.c* 

ip\_mroute.c

Ť,

Figure 14.41 phyint\_send function.

Section 14.8

#### tunnel\_send Function

To send a datagram on a tunnel\_tunnel\_send (Figure 14.43) must construct the appropriate tunnel options and insert them in the header of the outgoing datagram. Figure 14.42 shows how tunnel\_send prepares a packet for the tunnel.

```
635 static void
636 tunnel_send(m, vifp)
637 struct mbuf *m;
638 struct vif *vifp;
639 {
640
        struct ip *ip = mtod(m, struct ip *);
641
        struct mbuf *mb_copy, *mb_opts;
        struct ip *ip_copy;
642
643
        int
                error;
        u_char *cp;
644
645
        /*
646
         * Make sure that adding the tunnel options won't exceed the
647
         * maximum allowed number of option bytes.
         */
648
649
        if (ip->ip_hl > (60 - TUNNEL_LEN) >> 2) {
650
            mrtstat.mrts_cant_tunnel++;
651
            return;
652
        3
653
        /*
         \star Get a private copy of the IP header so that changes to some
654
655
         \ast of the IP fields don't damage the original header, which is
656
         * examined later in ip_input.c.
         */
657
658
        mb_copy = m_copy(m, IP_HDR_LEN, M_COPYALL);
659
        if (mb_copy == NULL)
660
            return;
661
        MGETHDR(mb_opts, M_DONTWAIT, MT_HEADER);
        if (mb_opts == NULL) {
662
663
            m_freem(mb_copy);
664
            return:
665
        }
        /*
666
         * Make mb_opts be the new head of the packet chain.
667
668
         * Any options of the packet were left in the old packet chain head
669
         */
670
        mb_opts->m_next = mb_copy;
671
        mb_opts->m_len = IP_HDR_LEN + TUNNEL_LEN;
672
        mb_opts->m_data += MSIZE - mb_opts->m_len;
                                                                          ip_mroute.c
```

Figure 14.43 tunnel\_send function: verify and allocate new header.

#### Will the tunnel options fit?

635-652 If there is no room in the IP header for the tunnel options, tunnel\_send returns immediately and the datagram is not forwarded on the tunnel. It may be forwarded on other interfaces.

- ip\_mroute.c

e.c

ι4

1)

to

2.C

to ·e-

fi-

Chapter 14

Sec

66

14

# Duplicate the datagram and allocate mbuf for new header and tunnel options

<sup>653–672</sup> In the call to m\_copy, the starting offset for the copy is 20 (IP\_HDR\_LEN). The resulting mbuf chain contains the options and data for the datagram but not the IP header. mb\_opts points to a new datagram header allocated by MGETHDR. The datagram header is prepended to mb\_copy. Then m\_len and m\_data are adjusted to accommodate an IP header and the tunnel options.

The second half of tunnel\_send, shown in Figure 14.44, modifies the headers of the outgoing packet and sends the packet.

```
ip mroute.c
        ip_copy = mtod(mb_opts, struct ip *);
673
674
        /*
         * Copy the base ip header to the new head mbuf.
675
676
         */
        *ip_copy = *ip;
677
        ip_copy->ip_ttl--;
678
        ip_copy->ip_dst = vifp->v_rmt_addr;
                                                  /* remote tunnel end-point */
679
        /*
680
         * Adjust the ip header length to account for the tunnel options.
681
         */
682
        ip_copy->ip_hl += TUNNEL_LEN >> 2;
683
        ip_copy->ip_len += TUNNEL_LEN;
684
685
         * Add the NOP and LSRR after the base ip header
686
687
         */
        cp = (u_char *) (ip_copy + 1);
688
        *cp++ = IPOPT_NOP;
689
        *cp++ = IPOPT_LSRR;
690
                                     /* LSRR option length */
691
        *cp++ = 11;
                                     /* LSSR pointer to second element */
692
        *cp++ = 8;
        *(u_long *) cp = vifp->v_lcl_addr.s_addr; /* local tunnel end-point */
693
694
        cp += 4;
                                                  /* destination group */
        *(u_long *) cp = ip->ip_dst.s_addr;
695
        error = ip_output(mb_opts, NULL, NULL, IP_FORWARDING, NULL);
696
697 }
                                                                         — ip_mroute.c
```

Figure 14.44 tunnel\_send function: construct headers and send.

#### Modify IP header

<sup>673–679</sup> The original IP header is copied from the original mbuf chain into the newly allocated mbuf header. The TTL in the header is decremented, and the destination is changed to be the other end of the tunnel.

#### **Construct tunnel options**

ip\_hl and ip\_len are adjusted to accommodate the tunnel options. The tunnel options are placed just after the IP header: a NOP, followed by the LSRR code, the length of the LSRR option (11 bytes), and a pointer to the *second* address in the option (8 bytes). The source route consists of the local tunnel end point followed by the destination group (Figure 14.13).

Section 14.9

er 14

The

e IP

.ata-

1 to

's of

ute.c

## Send the tunneled datagram

<code>ip\_output</code> sends the datagram, which now looks like a unicast datagram with an 665-697 LSRR option since the destination address is the unicast address of the other end of the tunnel. When it reaches the other end of the tunnel, the tunnel options are stripped off and the datagram is forwarded at that point, possibly through additional tunnels.

#### 14.9 Cleanup: ip\_mrouter\_done Function

When mrouted shuts down, it issues the DVMRP\_DONE command, which is handled by the ip\_mrouter\_done function shown in Figure 14.45.

		161 in	nt	— ip mroute.
			p_mrouter_done()	p_moute.
		163 {	j_mroacer_done()	
		164	vifi_t vifi;	
		165	int i;	
		166	struct ifnet *ifp;	
		167		
		168	int s; struct ifreq ifr;	
		169	<b>1</b> , ,	
		170	s = splnet(); /*	
	1	171	,	
		172	* For each phyint in use, free its local group list and	
	1	172	* disable promiscuous reception of all IP multicasts.	
		173	*/	
		174	<pre>for (vifi = 0; vifi &lt; numvifs; vifi++) {     if (vifi &gt; ) = f if (vifi)</pre>	
		175	if (viftable[vifi].v_lcl_addr.s_addr != 0 &&	
			!(viftable[vifi].v_flags & VIFF_TUNNEL)) (	
		177	if (viftable[vifi].v_lcl_grps)	
		178	<pre>free(viftable[vifi].v_lcl_grps, M_MRTABLE);</pre>	
*/	ļ	179	<pre>satosin(&amp;ifr.ifr_addr)-&gt;sin_family = AF_INET;</pre>	
		180	<pre>satosin(&amp;ifr.ifr_addr)-&gt;sin_addr.s_addr = INADDR_ANY;</pre>	
		181	ltp = vittable[vifi].v ifp;	
		182	(*ifp->if_ioctl) (ifp, SIOCDELMULTI, (caddr_t) & ifr)	;
	1	183	J	
2.C		184	}	
		185	<pre>bzero((caddr_t) viftable, sizeof(viftable));</pre>	
		186	numvifs = 0;	
		187	/*	
		188	* Free any multicast route entries.	
		189	* /	
0-		190	for (i = 0; i < MRTHASHSIZ; i++)	
is		191	if (mrttable[i])	
~		192	<pre>free(mrttable[i], M_MRTABLE);</pre>	
		193	<pre>bzero((caddr_t) mrttable, sizeof(mrttable));</pre>	
	erver.	194	cached_mrt = NULL;	
_1	19 TA	195	ip_mrouter = NULL;	
el	1 a	196	<pre>splx(s);</pre>	
ne		197	return (0);	
(8	12.	198 }		

Figure 14.45 ip\_mrouter\_done function: DVMRP\_DONE command.

Chapter 14

i Lie Mari

「「「ない」」

- <sup>161–186</sup> This function runs at splnet to avoid any interaction with the multicast forwarding code. For every physical multicast interface, the list of local groups is released and the SIOCDELMULTI command is issued to stop receiving multicast datagrams (Exercise 14.3). The entire viftable array is cleared by bzero and numvifs is set to 0.
- 187-198 Every active entry in the multicast routing table is released, the entire table is cleared with bzero, the cache is cleared, and ip\_mrouter is reset.

Each entry in the multicast routing table may be the first in a linked list of entries. This code introduces a memory leak by releasing only the first entry in the list.

## 14.10 Summary

In this chapter we described the general concept of internetwork multicasting and the specific functions within the Net/3 kernel that support it. We did not discuss the implementation of mrouted, but the source is readily available for the interested reader.

We described the virtual interface table and the differences between a physical interface and a tunnel, as well as the LSRR options used to implement tunnels in Net/3.

We illustrated the RPB, TRPB, and RPM algorithms and described the kernel tables used to forward multicast datagrams according to TRPB. The concept of parent and leaf networks was also discussed.

## **Exercises**

14.1 In Figure 14.25, how many multicast routes are needed?

- 14.2 Why is the update to the group membership cache in Figure 14.23 protected by splnet and splx?
- 14.3 What happens when SIOCDELMULTI is issued for an interface that has explicitly joined a multicast group with the IP\_ADD\_MEMBERSHIP option?
- 14.4 When a datagram arrives on a tunnel and is accepted by ip\_mforward, it may be looped back by ip\_output when it is forwarded to a physical interface. Why does ip\_mforward discard the looped-back packet when it arrives on the physical interface?
- 14.5 Redesign the group address cache to increase its effectiveness.

## 15.1 Intr

This abst face cuss fron fron crea netv gratacce thro wri

INTEL Ex.1013.460

prot

wor

not ( tion: gran 15

1-

d r-

is

le

e. \_د

ป 3. :s ป

t

а

d

:S

# Socket Layer

## 15.1 Introduction

This chapter is the first of three that cover the socket-layer code in Net/3. The socket abstraction was introduced with the 4.2BSD release in 1983 to provide a uniform interface to network and interprocess communication protocols. The Net/3 release discussed here is based on the 4.3BSD Reno version of sockets, which is slightly different from the earlier 4.2 releases used by many Unix vendors.

As described in Section 1.7, the socket layer maps protocol-independent requests from a process to the protocol-specific implementation selected when the socket was created.

To allow standard Unix I/O system calls such as read and write to operate with network connections, the filesystem and networking facilities in BSD releases are integrated at the system call level. Network connections represented by sockets are accessed through a descriptor (a small integer) in the same way an open file is accessed through a descriptor. This allows the standard filesystem calls such as read and write, as well as network-specific system calls such as sendmsg and recvmsg, to work with a descriptor associated with a socket.

Our focus is on the implementation of sockets and the associated system calls and not on how a typical program might use the socket layer to implement network applications. For a detailed discussion of the process-level socket interface and how to program network applications see [Stevens 1990] and [Rago 1993].

Figure 15.1 shows the layering between the socket interface in a process and the protocol implementation in the kernel.

S

1

INTEL Ex.1013.462

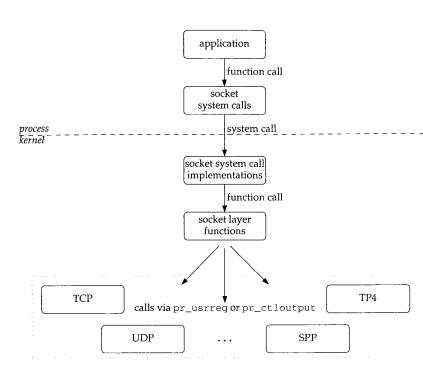


Figure 15.1 The socket layer converts generic requests to specific protocol operations.

#### splnet Processing

The socket layer contains many paired calls to splnet and splx. As discussed in Section 1.12, these calls protect code that accesses data structures shared between the socket layer and the protocol-processing layer. Without calls to splnet, a software interrupt that initiates protocol processing and changes the shared data structures will confuse the socket-layer code when it resumes.

We assume that readers understand these calls and we rarely point them out in our discussion.

## 15.2 Code Introduction

The three files listed in Figure 15.2 are described in this chapter.

## **Global Variables**

The two global variable covered in this chapter are described in Figure 15.3.

socket Structure 437

Section 15.3

File	Description
sys/socketvar.h	socket structure definitions
kern/uipc_syscalls.c kern/uipc_socket.c	system call implementation socket-layer functions

Figure 15.2 Files discussed in this chapter.

Variable	Datatype	Description
socketops	struct fileops	socket implementation of I/O system calls
sysent	struct sysent[]	array of system call entries

Figure 15.3	Global	variable	introduced	in	this	chapter.
-------------	--------	----------	------------	----	------	----------

#### 15.3 socket Structure

A socket represents one end of a communication link and holds or points to all the information associated with the link. This information includes the protocol to use, state information for the protocol (which includes source and destination addresses), queues of arriving connections, data buffers, and option flags. Figure 15.5 shows the definition of a socket and its associated buffers.

so\_type is specified by the process creating a socket and identifies the communica-41-42 tion semantics to be supported by the socket and the associated protocol. so\_type shares the same values as pr\_type shown in Figure 7.8. For UDP, so\_type would be SOCK\_DGRAM and for TCP it would be SOCK\_STREAM.

43

so\_options is a collection of flags that modify the behavior of a socket. Figure 15.4 describes the flags.

so_options	Kernel only	Description
SO_ACCEPTCONN SO_BROADCAST SO_DEBUG SO_DONTROUTE SO_KEEPALIVE SO_OOBINLINE SO_REUSEADDR SO_REUSEPORT SO_USELOOPBACK	•	socket accepts incoming connections socket can send broadcast messages socket records debugging information output operations bypass routing tables socket probes idle connections socket keeps out-of-band data inline socket can reuse a local address socket can reuse a local address and port routing domain sockets only; sending process receives its
		own routing requests

#### Figure 15.4 so\_options values.

A process can modify all the socket options with the getsockopt and setsockopt system calls except SO\_ACCEPTCONN, which is set by the kernel when the listen system call is issued on the socket.

438 Socket Layer

the state of the second state of the

1

Chapter 15

<pre>struct socket {</pre>				
short so_typ	pe;	/*	generic type, Figure 7.8 */	
	tions;		from socket call, Figure 15.4 */	
short so lin			time to linger while closing */	
short so_sta	ate;	/*	internal state flags, Figure 15.6 *	1
caddr_t so_pcl			protocol control block */	
struct protos			protocol handle */	
/*	,			
* Variables for a	connection qu	eueir	na.	
	-		head in all subsidiary sockets.	
	-		celated to an accept.	
			ctially completed connections,	
		-	tions ready to be accepted.	
	-		it has so head set, then	
* it has to be pu				
-			b based on current queue lengths	
			onnections for this socket.	
*/	mber of queue	eu ci	Sineccions for chis socket.	
struct socket	teo hoad.	/*	back pointer to accept socket */	
struct socket			queue of partial connections */	
struct socket			queue of incoming connections */	
short so_q01	-		partials on so q0 */	
short so_qu			number of connections on so_g */	
short so_ql:			max number queued connections */	
short so_q1			connection timeout */	
—				
u_short so_er pid_t so_pg:			error affecting connection */	
u_long so_ool			pgid for signals */ chars to oob mark */	
/*	JIIIALK;	/ "	chars to oob mark */	
* Variables for s	acket buffer			
*/	SOCKEC DUITEL.	ng.		
, struct sockbui	F f			
		/*	petual change in buffers t/	
u_long sł u_long sł			actual chars in buffer */	
u_long si u long si			<pre>max actual char count */ chars of mbufs used */</pre>	
u_long sl long sl			max chars of mbufs to use */	
	o_lowat;		low water mark */	
	uf *sb_mb;		the mbuf chain */	
	linfo sb_sel;		process selecting read/write */	
	o_flags;		Figure 16.5 */	
	o_timeo;	/*	timeout for read/write */	
} so_rcv, so_s				
caddr_t so_tpo			Wisc. protocol control block XXX */	
			<pre>ocket * so, caddr_t arg, int waitf);</pre>	
caddr_t so_upo	callarg;	/*	Arg for above */	

Figure 15.5 struct socket definition.

Section 15.3

so\_linger is the time in clock ticks that a socket waits for data to drain while closing a connection (Section 15.15).

45

44

so\_state represents the internal state and additional characteristics of the socket. Figure 15.6 lists the possible values for so\_state.

so_state	Kernel only	Description
SS_ASYNC SS_NBIO		socket should send asynchronous notification of I/O events socket operations should not block the process
SS_CANTRCVMORE SS_CANTSENDMORE SS_ISCONFIRMING SS_ISCONNECTED SS_ISCONNECTING SS_NOFDREF SS_PRIV SS_RCVATMARK	• • • • •	socket cannot receive more data from peer socket cannot send more data to peer socket is negotiating a connection request socket is connected to a foreign socket socket is connecting to a foreign socket socket is disconnecting from peer socket is not associated with a descriptor socket was created by a process with superuser privileges process has consumed all data received before the most recent out-of-band data was received

Figure 15.6 so\_state values.

In Figure 15.6, the middle column shows that SS\_ASYNC and SS\_NBIO can be changed explicitly by a process by the fcntl and ioctl system calls. The other flags are implicitly changed by the process during the execution of system calls. For example, if the process calls connect, the SS\_ISCONNECTED flag is set by the kernel when the connection is established.

## SS\_NBIO and SS\_ASYNC Flags

By default, a process blocks waiting for resources when it makes an I/O request. For example, a read system call on a socket blocks if there is no data available from the network. When the data arrives, the process is unblocked and read returns. Similarly, when a process calls write, the kernel blocks the process until space is available in the kernel for the data. If SS\_NBIO is set, the kernel does not block a process during I/O on the socket but instead returns the error code EWOULDBLOCK.

If SS\_ASYNC is set, the kernel sends the SIGIO signal to the process or process group specified by so\_pgid when the status of the socket changes for one of the following reasons:

- a connection request has completed,
- a disconnect request has been initiated,
- a disconnect request has completed,
- half of a connection has been shut down,
- data has arrived on a socket,
- data has been sent from a socket (i.e., the output buffer has free space), or
- an asynchronous error has occurred on a UDP or TCP socket.

so\_pcb points to a protocol control block that contains protocol-specific state information and parameters for the socket. Each protocol defines its own control block structure, so so\_pcb is defined to be a generic pointer. Figure 15.7 lists the control block structures that we discuss.

Protocol	Control block	Reference	
UDP	struct inpcb	Section 22.3	
ТСР	struct inpcb struct tcpcb	Section 22.3 Section 24.5	
ICMP, IGMP, raw IP	struct inpcb	Section 22.3	
Route	struct rawcb	Section 20.3	

so_pcb never points to a tcpcb structure di	irectly; see Figure 22.1.
---	---------------------------

Figure 15.7 Protocol control blocks.

47

so\_proto points to the protosw structure of the protocol selected by the process during the socket system call (Section 7.4).

48-64

65

67

Sockets with SO\_ACCEPTCONN set maintain two connection queues. Connections that are not yet established (e.g., the TCP three-way handshake is not yet complete) are placed on the queue so\_q0. Connections that are established and are ready to be accepted (e.g., the TCP three-way handshake is complete) are placed on the queue so\_q. The lengths of the queues are kept in so\_q0len and so\_qlen. Each queued connection is represented by its own socket. so\_head in each queued socket points to the original socket with SO\_ACCEPTCONN set.

The maximum number of queued connections for a particular socket is controlled by so\_qlimit, which is specified by a process when it calls listen. The kernel silently enforces an upper limit of 5 (SOMAXCONN, Figure 15.24) and a lower limit of 0. A somewhat obscure formula shown with Figure 15.29 uses so\_qlimit to control the number of queued connections.

Figure 15.8 illustrates a queue configuration in which three connections are ready to be accepted and one connection is being established.

so\_timeo is a *wait channel* (Section 15.10) used during accept, connect, and close processing.

so\_error holds an error code until it can be reported to a process during the next system call that references the socket.

If SS\_ASYNC is set for a socket, the SIGIO signal is sent to the process (if so\_pgid is greater than 0) or to the progress group (if so\_pgid is less than 0). so\_pgid can be changed or examined with the SIOCSPGRP and SIOCGPGRP ioctl commands. For more information about process groups see [Stevens 1992].

<sup>68</sup> so\_oobmark identifies the point in the input data stream at which out-of-band data was most recently received. Section 16.11 discusses socket support for out-of-band data and Section 29.7 discusses the semantics of out-of-band data in TCP.

<sup>69–82</sup> Each socket contains two data buffers, so\_rcv and so\_snd, used to buffer incoming and outgoing data. These are structures contained within the socket structure, not

Section 15.4

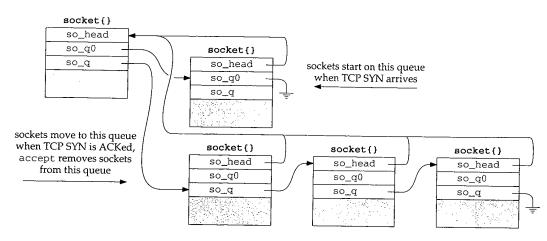


Figure 15.8 Socket connection queues.

pointers to structures. We describe the organization and use of the socket buffers in Chapter 16.

<sup>83-86</sup> so\_tpcb is not used by Net/3. so\_upcall and so\_upcallarg are used only by the NFS software in Net/3.

NFS is unusual. In many ways it is a process-level application that has been moved into the kernel. The so\_upcall mechanism triggers NFS input processing when data is added to a socket receive buffer. The tsleep and wakeup mechanism is inappropriate in this case, since the NFS protocol executes within the kernel, not as a process.

The files socketvar.h and uipc\_socket2.c define several macros and functions that simplify the socket-layer code. Figure 15.9 summarizes them.

## 15.4 System Calls

A process interacts with the kernel through a collection of well-defined functions called *system calls*. Before showing the system calls that support networking, we discuss the system call mechanism itself.

The transfer of execution from a process to the protected environment of the kernel is machine- and implementation-dependent. In the discussion that follows, we use the 386 implementation of Net/3 to illustrate implementation specific operations.

In BSD kernels, each system call is numbered and the hardware is configured to transfer control to a single kernel function when the process executes a system call. The particular system call is identified as an integer argument to the function. In the 386 implementation, syscall is that function. Using the system call number, syscall indexes a table to locate the sysent structure for the requested system call. Each entry in the table is a sysent structure: .

Chapter 15

Name	Description
sosendallatonce	Does the protocol associated with <i>so</i> require each send system call to result in a single protocol request?
	<pre>int sosendallatonce(struct socket *so);</pre>
soisconnecting	Set the socket state to SS_ISCONNECTING.
	int <b>soisconnecting</b> (struct socket *so);
soisconnected	See Figure 15.30.
soreadable	Will a read on so return information without blocking?
	<pre>int soreadable(struct socket *so);</pre>
sowriteable	Will a write on so return without blocking?
	<pre>int sowriteable(struct socket *so);</pre>
socantsendmore	Set the SS_CANTSENDMORE flag. Wake up any processes sleeping on the send buffer.
	<pre>int socantsendmore(struct socket *so);</pre>
socantrcvmore	Set the SS_CANTRCVMORE flag. Wake up processes sleeping on the receive buffer.
	<pre>int socantrcvmore(struct socket *so);</pre>
sodisconnect	Issue the PRU_DISCONNECT request.
	int <b>sodisconnect</b> (struct socket * <i>so</i> );
soisdisconnecting	Clear the SS_ISCONNECTING flag. Set SS_ISDISCONNECTING, SS_CANTRCVMORE, and SS_CANTSENDMORE flags. Wake up any processes selecting on the socket.
	<pre>int soisdisconnecting(struct socket *so);</pre>
soisdisconnected	Clear the SS_ISCONNECTING, SS_ISCONNECTED, and SS_ISDISCONNECTING flags. Set the SS_CANTRCVMORE and SS_CANTSENDMORE flags. Wake up any processes selecting on the socket or waiting for close to complete.
	<pre>int soisdisconnected(struct socket *so);</pre>
soqinsque	Insert <i>so</i> on a queue associated with <i>head</i> . If <i>q</i> is 0, the socket is added to the end of $so_q0$ , which holds incomplete connections. Otherwise, the socket is added to the end of $so_q$ , which holds connections that are ready to be accepted. Net/1 incorrectly placed sockets at the front of the queue.
	int <b>soginsque</b> (struct socket * <i>head</i> , struct socket * <i>so</i> , int <i>q</i> );
sogremque	Remove so from the queue identified by q. The socket queues are located by following so->so_head.
	<pre>int sogremque(struct socket *so, int q);</pre>

Figure 15.9 Socket macros and functions.

Section 15.4

}

.5

```
System Calls 443
```

Here are several entries from the sysent array, which is defined in kern/init\_sysent.c.

```
struct sysent sysent[] = {
    /* ... */
    { 3, recvmsg },    /* 27 = recvmsg */
    { 3, sendmsg },    /* 28 = sendmsg */
    { 6, recvfrom },    /* 29 = recvfrom */
    { 3, accept },    /* 30 = accept */
    { 3, getpeername },    /* 31 = getpeername */
    { 3, getsockname },    /* 32 = getsockname */
    /* ... */
}
```

For example, the recornsg system call is the 27th entry in the system call table, has three arguments, and is implemented by the recornsg function in the kernel.

syscall copies the arguments from the calling process into the kernel and allocates an array to hold the results of the system call, which syscall returns to the process when the system call completes. syscall dispatches control to the kernel function associated with the system call. In the 386 implementation, this call looks like:

```
struct sysent *callp;
error = (*callp->sy_call)(p, args, rval);
```

where callp is a pointer to the relevant sysent structure, p is a pointer to the process table entry for the process that made the system call, args represents the arguments to the system call as an array of 32-bit words, and rval is an array of two 32-bit words to hold the return value of the system call. When we use the term *system call*, we mean the function within the kernel called by syscall, not the function within the process called by the application.

syscall expects the system call function (i.e., what sy\_call points to) to return 0 if no errors occurred and a nonzero error code otherwise. If no error occurs, the kernel passes the values in rval back to the process as the return value of the system call (the one made by the application). If an error occurs, syscall ignores the values in rval and returns the error code to the process in a machine-dependent way so that the error is made available to the process in the external variable errno. The function called by the application returns -1 or a null pointer to indicate that errno should be examined.

The 386 implementation sets the carry bit to indicate that the value returned by syscall is an error code. The system call stub in the process stores the code in errno and returns -1 or a null pointer to the application. If the carry bit is not set, the value returned by syscall is returned by the stub.

To summarize, a function implementing a system call "returns" two values: one for the syscall function, and a second (found in rval) that syscall returns to the calling process when no error occurs.

INTEL Ex.1013.469

Chapter 15

## Example

The prototype for the socket system call is:

int socket(int domain, int type, int protocol);

The prototype for the kernel function that implements the system call is

```
struct socket_args {
    int domain;
    int type;
    int protocol;
};
socket(struct proc *p, struct socket_args *uap, int *retval);
```

When an application calls socket, the process passes three separate integers to the kernel with the system call mechanism. syscall copies the arguments into an array of 32-bit values and passes a pointer to the array as the second argument to the kernel version of socket. The kernel version of socket treats the second argument as a pointer to an socket\_args structure. Figure 15.10 illustrates this arrangement.

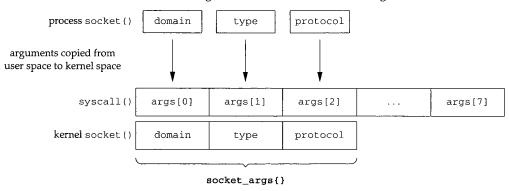


Figure 15.10 socket argument processing.

As illustrated by socket, each kernel function that implements a system call declares args not as a pointer to an array of 32-bit words, but as as a pointer to a structure specific to the system call.

The implicit cast is legal only in traditional K&R C or in ANSI C when a prototype is not in effect. If a prototype is in effect, the compiler generates a warning.

syscall prepares the return value of 0 before executing the kernel system call function. If no error occurs, the system call function can return without clearing \*retval and syscall returns 0 to the process.

## System Call Summary

Figure 15.11 summarizes the system calls relevant to networking.

Section 15.5

Category	Name	Function			
setup	socket	create a new unnamed socket within a specified communication domain			
-	bind	assign a local address to a socket			
server	listen	prepare a socket to accept incoming connections			
	accept	wait for and accept connections			
client	connect	establish a connection to a foreign socket			
	read	receive data into a single buffer			
	readv	receive data into multiple buffers			
input	recv	receive data specifying options			
mput	recvfrom	receive data and address of sender			
	recvmsg	receive data into multiple buffers, control information, and receive the			
		address of sender; specify receive options			
	write	send data from a single buffer			
	writev	send data from multiple buffers			
output	send	send data specifying options			
output	sendto	send data to specified address			
	sendmsg	send data from multiple buffers and control information to a specified			
		address; specify send options			
I/O	select	wait for I/O conditions			
termination	shutdown	terminate connection in one or both directions			
	close	terminate connection and release socket			
administration	fcntl	modify I/O semantics			
	ioctl	miscellaneous socket operations			
	setsockopt	set socket or protocol options			
	getsockopt	get socket or protocol options			
	getsockname	get local address assigned to socket			
	getpeername	get foreign address assigned to socket			

Figure 15.11 Networking system calls in Net/3.

We present the setup, server, client, and termination calls in this chapter. The input and output system calls are discussed in Chapter 16 and the administrative calls in Chapter 17.

Figure 15.12 shows the sequence in which an application might use the calls. The I/O system calls in the large box can be called in any order. This is not a complete state diagram as some valid transitions are not included; just the most common ones are shown.

# 15.5 Processes, Descriptors, and Sockets

Before describing the socket system calls, we need to discuss the data structures that tie together processes, descriptors, and sockets. Figure 15.13 shows the structures and members relevant to our discussion. A more complete explanation of the file structures can be found in [Leffler et al. 1989].

445

446 Socket Layer

ų,

,

in the second seco

10,02

AND TOTAL TOTAL

1111

0.00-000

「東京市 法国际部署 御台」

1.2.1



「「「「「「「」」」」

1

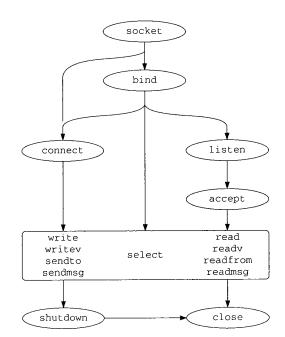
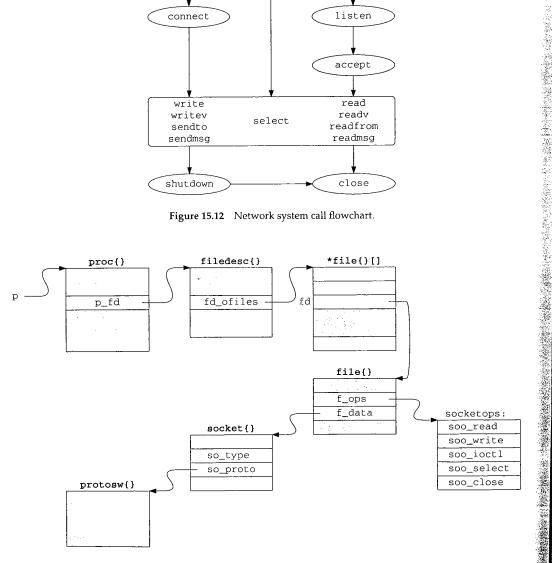


Figure 15.12 Network system call flowchart.





INTEL Ex.1013.472

Section 15.6

socket System Call 447

The first argument to a function implementing a system call is always p, a pointer to the proc structure of the calling process. The proc structure represents the kernel's notion of a process. Within the proc structure, p\_fd points to a filedesc structure, which manages the descriptor table pointed to by fd\_ofiles. The descriptor table is dynamically sized and consists of an array of pointers to file structures. Each file structure describes a single open file and can be shared between multiple processes.

Only a single file structure is shown in Figure 15.13. It is accessed by  $p \rightarrow p_fd \rightarrow fd_ofiles[fd]$ . Within the file structure, two members are of interest to us: f\_ops and f\_data. The implementation of I/O system calls such as read and write varies according to what type of I/O object is associated with a descriptor. f\_ops points to a fileops structure containing a list of function pointers that implement the read, write, ioctl, select, and close system calls for the associated I/O object. Figure 15.13 shows f\_ops pointing to a global fileops structure, socketops, which contains pointers to the functions for sockets.

f\_data points to private data used by the associated I/O object. For sockets, f\_data points to the socket structure associated with the descriptor. Finally, we see that so\_proto in the socket structure points to the protosw structure for the protocol selected when the socket is created. Recall that each protosw structure is shared by all sockets associated with the protocol.

We now proceed to discuss the system calls.

## 15.6 socket System Call

The socket system call creates a new socket and associates it with a protocol as specified by the domain, type, and protocol arguments specified by the process. The function (shown in Figure 15.14) allocates a new descriptor, which identifies the socket in future system calls, and returns the descriptor to the process.

<sup>42–55</sup> Before each system call a structure is defined to describe the arguments passed from the process to the kernel. In this case, the arguments are passed within a socket\_args structure. All the socket-layer system calls have three arguments: p, a pointer to the proc structure for the calling process; uap, a pointer to a structure containing the arguments passed by the process to the system call; and retval, a value-result argument that points to the return value for the system call. Normally, we ignore the p and retval arguments and refer to the contents of the structure pointed to by uap as the arguments to the system call.

<sup>56-60</sup> falloc allocates a new file structure and slot in the fd\_ofiles array (Figure 15.13). fp points to the new structure and fd is the index of the structure in the fd\_ofiles array. socket enables the file structure for read and write access and marks it as a socket. socketops, a global fileops structure shared by all sockets, is attached to the file structure by f\_ops. The socketops variable is initialized at compile time as shown in Figure 15.15.

60-69 socreate allocates and initializes a socket structure. If socreate fails, the error code is posted in error, the file structure is released, and the descriptor slot cleared. If socreate succeeds, f\_data is set to point to the socket structure and establishes

ops: ad ite oct1 elect

pter 15

Socket Layer

Chapter 15

Se

いたい たいたけん 品格 山口花 時間 かいまた 約100 高級 可には登録

		uipc_syscalls.c
42 stru	act socket_args {	
	int domain;	
44	int type;	
45	int protocol;	
46 };		
47 soc}	<pre>ket(p, uap, retval)</pre>	
48 stru	lct proc *p;	
49 stri	uct socket_args *uap;	
50 int		
51 {		
52	struct filedesc *fdp = p->p_fd	;
53	struct socket *so;	
54	struct file *fp;	
55	int fd, error;	
56	if (error = falloc(p, &fp, &fc	() }
57	return (error);	
58	fp->f_flag = FREAD   FWRITE;	
59	fp->f_type = DTYPE_SOCKET;	
60		turne (upp-sprotocol)) {
61	if (error = socreate(uap->dom	ain, &so, uap->type, uap->protocol)) {
62	fdp->fd_ofiles[fd] = 0;	
63	<pre>ffree(fp);</pre>	
64	} else {	
65	fp->f_data = (caddr_t) so	;
66	<pre>*retval = fd;</pre>	
67	}	
68	return (error);	
69 }		uipc_syscalls.c

Figure 15.14 socket system call.

Member	Value
fo_read	soo_read
fo_write	soo_write
fo_ioctl	soo_ioctl
fo_select	soo_select
fo_close	soo_close

Figure 15.15 socketops: the global fileops structure for sockets.

the association between the descriptor and the socket. fd is returned to the process through \*retval. socket returns 0 or the error code returned by socreate.

## socreate Function

Most socket system calls are divided into at least two functions, in the same way that socket and socreate are. The first function retrieves from the process all the data

## 448

required, calls the second soxxx function to do the work, and then returns any results to the process. This split is so that the second function can be called directly by kernel-based network protocols, such as NFS. socreate is shown in Figure 15.16.

43 so	create(dom, aso, type, proto)	– uipc_socket.c
44 in		
45 st	ruct socket **aso;	
46 in		
47 in	t proto;	
48 {		
49	struct proc *p = curproc; /* XXX */	
50	struct protosw *prp;	
51	struct socket *so;	
52	int error;	
53	if (proto)	
54	<pre>prp = pffindproto(dom, proto, type);</pre>	
55	else	
56	<pre>prp = pffindtype(dom, type);</pre>	
57	if (prp == 0 (  prp->pr_usrreq == 0)	
58	return (EPROTONOSUPPORT);	
59	if (prp->pr_type != type)	
60	return (EPROTOTYPE);	
61	MALLOC(so, struct socket *, sizeof(*so), M_SOCKET, M_WAIT);	
62	<pre>bzero((caddr_t) so, sizeof(*so));</pre>	
63	so->so_type = type;	
64	if (p->p_ucred->cr_uid == 0)	
65	so->so_state = SS_PRIV;	
66	so->so_proto = prp;	
67	error =	
68	(*prp->pr_usrreq) (so, PRU_ATTACH,	
69	(struct mbuf *) 0, (struct mbuf *) proto, (struct mb	uf *) 0).
70	if (error) {	ul ") 0);
71	<pre>so-&gt;so_state != SS_NOFDREF;</pre>	
72	sofree(so);	
73	return (error);	
74	}	
75	*aso = so;	
76	return (0);	
77 }		

uipc\_socket.c

#### Figure 15.16 socreate function.

The four arguments to socreate are: dom, the requested protocol domain (e.g., PF\_INET); aso, in which a pointer to a new socket structure is returned; type, the requested socket type (e.g., SOCK\_STREAM); and proto, the requested protocol.

#### Find protocol switch table

<sup>53-60</sup> If proto is nonzero, pffindproto looks for the specific protocol requested by the process. If proto is 0, pffindtype looks for a protocol within the specified domain with the semantics specified by type. Both functions return a pointer to a protosw structure of the matching protocol or a null pointer (Section 7.6).

er 15

# Allocate and initialize socket structure

socreate allocates a new socket structure, fills it with 0s, records the type, and, if the calling process has superuser privileges, turns on SS\_PRIV in the socket structure. 61-66

## PRU\_ATTACH request

67-69

The first example of the protocol-independent socket layer making a protocolspecific request appears in socreate. Recall from Section 7.4 and Figure 15.13 that so->so\_proto->pr\_usrreq is a pointer to the user-request function of the protocol associated with socket so. Every protocol provides this function in order to handle communication requests from the socket layer. The prototype for the function is:

int pr\_usrreg(struct socket \*so, int req, struct mbuf \*m0, \*m1, \*m2);

The first argument, so, is a pointer to the relevant socket and req is a constant identifying the particular request. The next three arguments (m0, m1, and m2) are different for each request. They are always passed as pointers to mbuf structures, even if they have another type. Casts are used when necessary to avoid warnings from the compiler. Figure 15.17 shows the requests available through the pr\_usrreq function. The

semantics of each request depend on the particular protocol servicing the request.

	Arguments			Description	
Request	m0	m1	m2	·	
PRU_ABORT PRU_ACCEPT PRU_ATTACH PRU_BIND PRU_CONNECT PRU_CONNECT2 PRU_DETACH PRU_DISCONNECT PRU_LISTEN PRU_REVD PRU_REVOOB PRU_SEND PRU_SENDOB PRU_SENDOB PRU_SCKADDR	buffer data data	address protocol address address socket2 buffer flags flags address address buffer	control control	abort any existing connection wait for and accept a connection a new socket has been created bind the address to the socket establish association or connection to address connect two sockets together socket is being closed break association between socket and foreign address begin listening for connections return foreign address associated with socket process has accepted some data receive OOB data send regular data send OOB data end communication with foreign address return local address associated with socket	

Figure 15.17 pr\_usrreq requests.

PRU\_CONNECT2 is supported only within the Unix domain, where it connects two local sockets to each other. Unix pipes are implemented in this way.

## **Cleanup and return**

70-77

Returning to socreate, the function attaches the protocol switch table to the new socket and issues the PRU\_ATTACH request to notify the protocol of the new end point. This request causes most protocols, including TCP and UDP, to allocate and initialize any structures required to support the new end point.

#### Superuser Privileges

Figure 15.18 summarizes the networking operations that require superuser access.

	Superuser			
Function	Process	Socket	Description	Reference
in_control		•	interface address, netmask, and destination address assignment	Figure 6.14
in_control		•	broadcast address assignment	Figure 6.22
in_pcbbind	•		binding to an Internet port less than 1024	Figure 22.22
ifioctl	•		interface configuration changes	Figure 4.29
ifioctl	•		multicast address configuration (see text)	Figure 12.11
rip_usrreq	•		creating an ICMP, IGMP, or raw IP socket	Figure 32.10
slopen	•		associating a SLIP device with a tty device	Figure 5.9

Figure 15.18 Superuser privileges in Net/3.

The multicast ioctl commands (SIOCADDMULTI and SIOCDELMULTI) are accessible to nonsuperuser processes when they are invoked indirectly by the IP\_ADD\_MEMBERSHIP and IP\_DROP\_MEMBERSHIP socket options (Sections 12.11 and 12.12).

In Figure 15.18, the "Process" column identifies requests that must be made by a superuser process, and the "Socket" column identifies requests that must be issued on a socket *created* by a superuser process (i.e., the process does not need superuser privileges if it has access to the socket, Exercise 15.1). In Net/3, the suser function determines if the calling process has superuser privileges, and the SS\_PRIV flag determines if the socket was created by a superuser process.

Since rip\_usrreq tests SS\_PRIV immediately after creating the socket with socreate, we show this function as accessible only from a superuser process.

# 15.7 getsock and sockargs Functions

These functions appear repeatedly in the implementation of the socket system calls. getsock maps a descriptor to a file table entry and sockargs copies arguments from the process to a newly allocated mbuf in the kernel. Both functions check for invalid arguments and return a nonzero error code accordingly.

Figure 15.19 shows the getsock function.

754-767

768-783

The function selects the file table entry specified by the descriptor fdes with fdp, a pointer to the filedesc structure. getsock returns a pointer to the open file structure in fpp or an error if the descriptor is out of the valid range, does not point to an open file, or does not have a socket associated with it.

Figure 15.20 shows the sockargs function.

The mechanism described in Section 15.4 copies pointer arguments for a system call from the process to the kernel but does not copy the data referenced by the pointers, since the semantics of each argument are known only by the specific system call and not

759

760

761

762

763

764 765

766

773

774

794 }

Chapter 15

uipc\_syscalls.c 754 getsock(fdp, fdes, fpp) 755 struct filedesc \*fdp; - 7 fdes; 756 int 757 struct file \*\*fpp; 758 { struct file \*fp; if ((unsigned) fdes >= fdp->fd\_nfiles || (fp = fdp->fd\_ofiles[fdes]) == NULL) return (EBADF); if (fp->f\_type != DTYPE\_SOCKET) return (ENOTSOCK); \*fpp = fp; return (0); uipc\_syscalls.c 767 } Figure 15.19 getsock function. uipc\_syscalls.c 768 sockargs(mp, buf, buflen, type) 769 struct mbuf \*\*mp; 770 caddr\_t buf; 日本の変換の buflen, type; 771 int 772 { struct sockaddr \*sa; struct mbuf \*m; int error; 775 if ((u\_int) buflen > MLEN) { 776 return (EINVAL); 777 778 } m = m\_get(M\_WAIT, type); 779 if (m == NULL) 780 return (ENOBUFS); 781 m->m\_len = buflen; error = copyin(buf, mtod(m, caddr\_t), (u\_int) buflen); 782 783 if (error) 784 (void) m\_free(m); 785 else { 「「「「「「「「」」」」 786 \*mp = m;787 if (type == MT\_SONAME) { 788 sa = mtod(m, struct sockaddr \*); 789 sa->sa\_len = buflen; 790 } 791 } 792 return (error); 793 - uipc\_syscalls.c

Figure 15.20 sockargs function.

by the generic system call mechanism. Several system calls use sockargs to follow the pointer arguments and copy the referenced data from the process into a newly allocated mbuf within the kernel. For example, sockargs copies the local socket address pointed to by bind's second argument from the process to an mbuf.

15.

Seci

784

786

If the data does not fit in a single mbuf or an mbuf cannot be allocated, sockargs returns EINVAL or ENOBUFS. Note that a standard mbuf is used and not a packet header mbuf. copyin copies the data from the process into the mbuf. The most common error from copyin is EACCES, returned when the process provides an invalid address.

When an error occurs, the mbuf is discarded and the error code is returned. If there is no error, a pointer to the mbuf is returned in mp, and sockargs returns 0.

786-794

s.c

he

ed ess If type is MT\_SONAME, the process is passing in a sockaddr structure. sockargs sets the internal length, sa\_len, to the length of the argument just copied. This ensures that the size contained within the structure is correct even if the process did not initialize the structure correctly.

Net/3 does include code to support applications compiled on a pre-4.3BSD Reno system, which did not have an sa\_len member in the sockaddr structure, but that code is not shown in Figure 15.20.

## 15.8 bind System Call

The bind system call associates a local network transport address with a socket. A process acting as a client usually does not care what its local address is. In this case, it isn't necessary to call bind before the process attempts to communicate; the kernel selects and implicitly binds a local address to the socket as needed.

A server process almost always needs to bind to a specific well-known address. If so, the process must call bind before accepting connections (TCP) or receiving datagrams (UDP), because the clients establish connections or send datagrams to the wellknown address.

A socket's foreign address is specified by connect or by one of the write calls that allow specification of foreign addresses (sendto or sendmsg).

Figure 15.21 shows bind.

70-82 The arguments to bind (passed within a bind\_args structure) are: s, the socket descriptor; name, a pointer to a buffer containing the transport address (e.g., a sockaddr\_in structure); and namelen, the size of the buffer.

getsock returns the file structure for the descriptor, and sockargs copies the local address from the process into an mbuf, sobind associates the address specified by the process with the socket. Before bind returns sobind's result, the mbuf holding the address is released.

Technically, a descriptor such as s identifies a file structure with an associated socket structure and is not itself a socket structure. We refer to such a descriptor as a socket to simplify our discussion.

We will see this pattern many times: arguments specified by the process are copied into an mbuf and processed as necessary, and then the mbuf is released before the system call returns. Although mbufs were designed explicitly to facilitate processing of network data packets, they are also effective as a general-purpose dynamic memory allocation mechanism. 454 Socket Layer

Chapter 15

```
uipc_syscalls.c
70 struct bind_args {
71
      int
              s;
72
       caddr_t name;
       int namelen;
73
74 };
75 bind(p, uap, retval)
76 struct proc *p;
77 struct bind_args *uap;
          *retval;
78 int
79 {
       struct file *fp;
80
       struct mbuf *nam;
81
       int
                error;
82
       if (error = getsock(p->p_fd, uap->s, &fp))
83
            return (error);
84
       if (error = sockargs(&nam, uap->name, uap->namelen, MT_SONAME))
85
           return (error);
86
        error = sobind((struct socket *) fp->f_data, nam);
 87
        m_freem(nam);
 88
        return (error);
 89

    uipc_syscalls.c

 90 }
```

## Figure 15.21 bind function.

Another pattern illustrated by bind is that retval is unused in many system calls. In Section 15.4 we mentioned that retval is always initialized to 0 before syscall dispatches control to a system call. If 0 is the appropriate return value, the system calls do not need to change retval.

## sobind Function

sobind, shown in Figure 15.22, is a wrapper that issues the PRU\_BIND request to the protocol associated with the socket.

```
78 sobind(so, nam)
79 struct socket *so;
80 struct mbuf *nam;
81 {
82
       int
               s = splnet();
83
       int
               error;
84
       error =
            (*so->so_proto->pr_usrreq) (so, PRU_BIND,
85
                                  (struct mbuf *) 0, nam, (struct mbuf *) 0);
86
87
       splx(s);
       return (error);
88
                                                                        - uipc_socket.c
89 }
```

Figure 15.22 sobind function.

uipc\_socket.c

15

's.c

'ls.c

lls.

ılls

the

:et.c

ket.c

78-89 sobind issues the PRU\_BIND request. The local address, nam, is associated with the socket if the request succeeds; otherwise the error code is returned.

## 15.9 listen System Call

The listen system call, shown in Figure 15.23, notifies a protocol that the process is prepared to accept incoming connections on the socket. It also specifies a limit on the number of connections that can be queued on the socket, after which the socket layer refuses to queue additional connection requests. When this occurs, TCP ignores incoming connection requests. Queued connections are made available to the process when it calls accept (Section 15.11).

```
uipc_syscalls.c
 91 struct listen_args {
 92
        int
                 s;
 93
                 backlog;
        int
 94 };
 95 listen(p, uap, retval)
 96 struct proc *p;
 97 struct listen_args *uap;
 98 int
           *retval;
99 {
100
        struct file *fp;
101
        int
                error;
102
        if (error = getsock(p->p_fd, uap->s, &fp))
103
            return (error);
104
        return (solisten((struct socket *) fp->f_data, uap->backlog));
105 }
                                                                        - uipc_syscalls.c
```

Figure 15.23 listen system call.

*91–98* The two arguments passed to listen specify the socket descriptor and the connection queue limit.

99–105 getsock returns the file structure for the descriptor, s, and solisten passes the listen request to the protocol layer.

## solisten Function

This function, shown in Figure 15.24, issues the PRU\_LISTEN request and prepares the socket to receive connections.

90-109 After solisten issues the PRU\_LISTEN request and pr\_usrreq returns, the socket is marked as ready to accept connections. SS\_ACCEPTCONN is not set if a connection is queued when pr\_usrreq returns.

The maximum queue size for incoming connections is computed and saved in so\_qlimit. Here Net/3 silently enforces a lower limit of 0 and an upper limit of 5 (SOMAXCONN) backlogged connections.

**INTEL Ex.1013.481** 

Chapter 15

```
uipc_socket.c
 90 solisten(so, backlog)
91 struct socket *so;
 92 int
           backlog;
93 {
                s = splnet(), error;
94
        int
 95
        error =
           (*so->so_proto->pr_usrreq) (so, PRU_LISTEN,
 96
                     (struct mbuf *) 0, (struct mbuf *) 0, (struct mbuf *) 0);
 97
 98
        if (error) {
           splx(s);
 99
            return (error);
100
101
        }
        if (so->so_q == 0)
102
            so->so_options |= SO_ACCEPTCONN;
103
        if (backlog < 0)
104
           backlog = 0;
105
        so->so_qlimit = min(backlog, SOMAXCONN);
106
107
        splx(s);
        return (0);
108
109 }
                                                                         uipc_socket.c
```

Figure 15.24 solisten function.

#### 15.10 tsleep and wakeup Functions

When a process executing within the kernel cannot proceed because a kernel resource is unavailable, it waits for the resource by calling tsleep, which has the following proto-type:

int tsleep(caddr\_t chan, int pri, char \*mesg, int timeo);

The first argument to tsleep, *chan*, is called the *wait channel*. It identifies the particular resource or event such as an incoming network connection, for which the process is waiting. Many processes can be sleeping on a single wait channel. When the resource becomes available or when the event occurs, the kernel calls wakeup with the wait channel as the single argument. The prototype for wakeup is:

void wakeup(caddr\_t chan);

All processes waiting for the channel are awakened and set to the run state. The kernel arranges for tsleep to return when each of the processes resumes execution.

The *pri* argument specifies the priority of the process when it is awakened, as well as several optional control flags for tsleep. By setting the PCATCH flag in *pri*, tsleep also returns when a signal arrives. *mesg* is a string identifying the call to tsleep and is included in debugging messages and in ps output. *timeo* sets an upper bound on the sleep period and is measured in clock ticks.

Figure 15.25 summarizes the return values from tsleep.

A process never sees the ERESTART error because it is handled by the syscall function and never returned to a process.

E)

15

Č.

tsleep()	Description
0	The process was awakened by a matching call to wakeup.
EWOULDBLOCK	The process was awakened after sleeping for <i>timeo</i> clock ticks and before the matching call to wakeup.
ERESTART	A signal was handled by the process during the sleep and the pending system call should be restarted.
EINTR	A signal was handled by the process during the sleep and the pending system call should fail.

Figure 15.25 tsleep return values.

Because all processes sleeping on a wait channel are awakened by wakeup, we always see a call to tsleep within a tight loop. Every process must determine if the resource is available before proceeding because another awakened process may have claimed the resource first. If the resource is not available, the process callstsleep once again.

It is unusual for multiple processes to be sleeping on a single socket, so a call to wakeup usually causes only one process to be awakened by the kernel.

For a more detailed discussion of the sleep and wakeup mechanism see [Leffler et al. 1989].

## Example

One use of multiple processes sleeping on the same wait channel is to have multiple server processes reading from a UDP socket. Each server calls recvfrom and, as long as no data is available, the calls block in tsleep. When a datagram arrives on the socket, the socket layer calls wakeup and each server is placed on the run queue. The first server to run receives the datagram while the others call tsleep again. In this way, incoming datagrams are distributed to multiple servers without the cost of starting a new process for each datagram. This technique can also be used to process incoming connection requests in TCP by having multiple processes call accept on the same socket. This technique is described in [Comer and Stevens 1993].

## 15.11 accept System Call

After calling listen, a process waits for incoming connections by calling accept, which returns a descriptor that references a new socket connected to a client. The original socket, s, remains unconnected and ready to receive additional connections. accept returns the address of the foreign system if name points to a valid buffer.

The connection-processing details are handled by the protocol associated with the socket. For TCP, the socket layer is notified when a connection has been established (i.e., when TCP's three-way handshake has completed). For other protocols, such as OSI's TP4, tsleep returns when a connection request has arrived. The connection is completed when explicitly confirmed by the process by reading or writing on the socket.

·ket.c

r 15

ket.c

ce is oto-

parcess the the

. well .eep 1d is

1 the

n and

The

Figure 15.26 shows the implementation of accept.

```
uipc_syscalls.c
 106 struct accept_args {
 107
         int
                  s;
 108
         caddr_t name;
 109
         int
                *anamelen;
 110 };
 111 accept(p, uap, retval)
 112 struct proc *p;
 113 struct accept_args *uap;
 114 int
            *retval;
 115 {
 116
         struct file *fp;
 117
         struct mbuf *nam;
 118
         int
                 namelen, error, s;
 119
         struct socket *so;
 120
         if (uap->name && (error = copyin((caddr_t) uap->anamelen,
 121
                                          (caddr_t) & namelen, sizeof(namelen))))
 122
             return (error);
 123
         if (error = getsock(p->p_fd, uap->s, &fp))
 124
             return (error);
 125
         s = splnet();
 126
         so = (struct socket *) fp->f_data;
 127
         if ((so->so_options & SO_ACCEPTCONN) == 0) {
 128
             splx(s);
 129
             return (EINVAL);
130
         }
131
         if ((so->so_state & SS_NBIO) && so->so_qlen == 0) (
132
             splx(s);
133
             return (EWOULDBLOCK);
134
         }
135
         while (so->so_qlen == 0 && so->so_error == 0) {
136
             if (so->so_state & SS_CANTRCVMORE) {
137
                 so->so_error = ECONNABORTED;
138
                 break:
139
             }
140
             if (error = tsleep((caddr_t) & so->so_timeo, PSOCK | PCATCH,
141
                                netcon, 0)) {
142
                 splx(s);
143
                 return (error);
144
             }
145
        }
146
        if (so->so_error) {
147
            error = so->so_error;
148
            so->so\_error = 0;
149
            splx(s);
150
            return (error);
151
        }
152
        if (error = falloc(p, &fp, retval)) {
153
            splx(s);
154
            return (error);
155
        }
```

**INTEL Ex.1013.484** 

「「「「「「「「」」」」

読んが

「読品」で

AN BERTHALL

uipc\_syscalls.c

156	{ struct socket *aso = so->so_g;
157	if $(sogremque(aso, 1) == 0)$
158	<pre>panic("accept");</pre>
159	SO = aSO;
160	}
161	fp->f_type = DTYPE_SOCKET;
162	fp->f_flag = FREAD   FWRITE;
163	fp->f_ops = &socketops
164	fp->f_data = (caddr_t) so;
165	<pre>nam = m_get(M_WAIT, MT_SONAME);</pre>
166	<pre>(void) soaccept(so, nam);</pre>
167	if (uap->name) {
168	if (namelen > nam->m_len)
169	<pre>namelen = nam-&gt;m_len;</pre>
170	/* SHOULD COPY OUT A CHAIN HERE */
171	if ((error = copyout(mtod(nam, caddr_t), (caddr_t) uap->name,
172	$(u_int)$ namelen $) == 0$
173	error = copyout((caddr_t) & namelen,
174	<pre>(caddr_t) uap-&gt;anamelen, sizeof(*uap-&gt;anamelen));</pre>
175	}
176	<pre>m_freem(nam);</pre>
177	<pre>splx(s);</pre>
178	return (error);
179 }	

Figure 15.26 accept system call.

The three arguments to accept (in the accept\_args structure) are: s, the socket descriptor; name, a pointer to a buffer to be filled in by accept with the transport address of the foreign host; and anamelen, a pointer to the size of the buffer.

#### Validate arguments

116-134 accept copies the size of the buffer (\*anamelen) into namelen, and getsock returns the file structure for the socket. If the socket is not ready to accept connections (i.e., listen has not been called) or nonblocking I/O has been requested and no connections are queued, EINVAL or EWOULDBLOCK are returned respectively.

#### Wait for a connection

<sup>135-145</sup> The while loop continues until a connection is available, an error occurs, or the socket can no longer receive data. accept is not automatically restarted after a signal is caught (tsleep returns EINTR). The protocol layer wakes up the process when it inserts a new connection on the queue with sonewconn.

Within the loop, the process waits in tsleep, which returns 0 when a connection is available. If tsleep is interrupted by a signal or the socket is set for nonblocking semantics, accept returns EINTR or EWOULDBLOCK (Figure 15.25).

#### Asynchronous errors

146-151 If an error occurred on the socket during the sleep, the error code is moved from the socket to the return value for accept, the socket error is cleared, and accept returns.

It is common for asynchronous events to change the state of a socket. The protocol processing layer notifies the socket layer of the change by setting so\_error and waking any process waiting on the socket. Because of this, the socket layer must always examine so\_error after waking to see if an error occurred while the process was sleeping.

## Associate socket with descriptor

<sup>152-164</sup> falloc allocates a descriptor for the new connection; the socket is removed from the accept queue by sogremque and attached to the file structure. Exercise 15.4 discusses the call to panic.

#### Protocol processing

167-179 accept allocates a new mbuf to hold the foreign address and calls soaccept to do protocol processing. The allocation and queueing of new sockets created during connection processing is described in Section 15.12. If the process provided a buffer to receive the foreign address, copyout copies the address from nam and the length from namelen to the process. If necessary, copyout silently truncates the name to fit in the process's buffer. Finally, the mbuf is released, protocol processing enabled, and accept returns.

Because only one mbuf is allocated for the foreign address, transport addresses must fit in one mbuf. Unix domain addresses, which are pathnames in the filesystem (up to 1023 bytes in length), may encounter this limit, but there is no problem with the 16-byte sockaddr\_in structure for the Internet domain. The comment on line 170 indicates that this limitation could be removed by allocating and copying an mbuf chain.

#### soaccept Function

soaccept, shown in Figure 15.27, calls the protocol layer to retrieve the client's address for the new connection.

```
uipc_socket.c
184 soaccept(so, nam)
185 struct socket *so;
186 struct mbuf *nam;
187 {
188
        int
                s = splnet();
189
        int
                error;
190
        if ((so->so_state & SS_NOFDREF) == 0)
191
            panic("soaccept: !NOFDREF");
        so->so_state &= ~SS NOFDREF;
192
193
        error = (*so->so_proto->pr_usrreq) (so, PRU_ACCEPT,
194
                                     (struct mbuf *) 0, nam, (struct mbuf *) 0);
195
        splx(s):
196
        return (error);
197 }
                                                                         uipc socket.c
```

Figure 15.27 soaccept function.

184-197 soaccept ensures that the socket is associated with a descriptor and issues the PRU\_ACCEPT request to the protocol. After pr\_usrreq returns, nam contains the name of the foreign socket.

# 15.12 sonewconn and soisconnected Functions

In Figure 15.26 we saw that accept waits for the protocol layer to process incoming connection requests and to make them available through so\_q. Figure 15.28 uses TCP to illustrate this process.

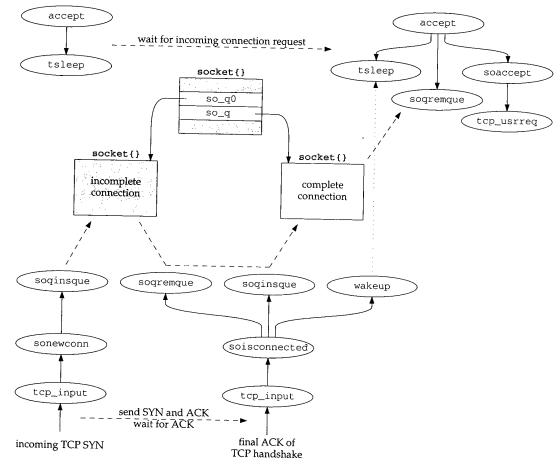


Figure 15.28 Incoming TCP connection processing.

In the upper left corner of Figure 15.28, accept calls tsleep to wait for incoming connections. In the lower left, tcp\_input processes an incoming TCP SYN by calling sonewconn to create a socket for the new connection (Figure 28.7). sonewconn queues the socket on so\_q0, since the three-way handshake is not yet complete.

Chapter 15

uipc\_socket2.c

When the final ACK of the TCP handshake arrives, tcp\_input calls soisconnected (Figure 29.2), which updates the new socket, moves it from so\_q0 to so\_q, and wakes up any processes that had called accept to wait for incoming connections.

The upper right corner of the figure shows the functions we described with Figure 15.26. When tsleep returns, accept takes the connection off so\_q and issues the PRU\_ATTACH request. The socket is associated with a new file descriptor and returned to the calling process.

Figure 15.29 shows the sonewconn function.

123 struc	ct socket *
124 sonew	wconn(head, connstatus)
125 struc	ct socket *head;
126 int	connstatus;
127 {	
128	struct socket *so;
129	int soqueue = connstatus ? 1 : 0;
130	if (head->so_qlen + head->so_q0len > 3 * head->so_qlimit / 2)
	(a + b) = (a + b) + (a + b) + (b) = (a + b) + (b) = (a + b) + (b
131	return ((struct socket , o), M_SOCKET, M_DONTWAIT); MALLOC(so, struct socket *, sizeof(*so), M_SOCKET, M_DONTWAIT);
132	if (so == NULL)
133	return ((struct socket *) 0);
134	<pre>bzero((caddr_t) so, sizeof(*so));</pre>
120	$s_{0}$ type = head->so type;
137	so->so_options = head->so_options & ~SO_ACCEPTCONN;
137	so->so_linger = head->so_linger;
138	so->so_state = head->so_state   SS_NOFDREF;
140	so->so_proto = head->so_proto;
140	so->so_timeo = head->so_timeo;
141	in head ago preid:
143	<pre>so-so_pgid = nead-so_pgid; (void) soreserve(so, head-so_snd.sb_hiwat, head-so_rcv.sb_hiwat);</pre>
144	socinsque (head, so, soqueue);
145	La ann Marrog) (so PBU A'l"'ACH
146	if ((*so->so_proto->pr_usrreq; (so, rho_hermin, (struct mbuf *) 0, (struct mbuf *) 0, (struct mbuf *) 0)) {
140	(void) sogremque(so, soqueue);
148	<pre>(void) free((caddr_t) so, M_SOCKET);</pre>
149	return ((struct socket *) 0);
150	
151	if (connstatus) {
152	sorwakeup(head);
153	<pre>wakeup((caddr_t) &amp; head-&gt;so_timeo);</pre>
154	so->so_state  = connstatus;
155	}
156	return (so);
157 }	uipc_socket2.c



<sup>123–129</sup> The protocol layer passes head, a pointer to the socket that is accepting the incoming connection, and connstatus, a flag to indicate the state of the new connection. For TCP, connstatus is always 0.

13

13

14

14

15

Sec

For TP4, connstatus is always SS\_ISCONFIRMING. The connection is implicitly confirmed when a process begins reading from or writing to the socket.

#### Limit incoming connections

130-131 sonewconn prohibits additional connections when the following inequality is true:

$$so_qlen + so_q0len > \frac{3 \times so_qlimit}{2}$$

This formula provides a fudge factor for connections that never complete and guarantees that listen(fd, 0) allows one connection. See Figure 18.23 in Volume 1 for an additional discussion of this formula.

#### Allocate new socket

132-143

144

A new socket structure is allocated and initialized. If the process calls setsockopt for the listening socket, the connected socket inherits several socket options because so\_options, so\_linger, so\_pgid, and the sb\_hiwat values are copied into the new socket structure.

#### Queue connection

soqueue was set from connstatus on line 129. The new socket is inserted onto so\_q0 if soqueue is 0 (e.g., TCP connections) or onto so\_q if connstatus is nonzero (e.g., TP4 connections).

### Protocol processing

145—150

The PRU\_ATTACH request is issued to perform protocol layer processing on the new connection. If this fails, the socket is dequeued and discarded, and sonewconn returns a null pointer.

#### Wakeup processes

151-157

If connstatus is nonzero, any processes sleeping in accept or selecting for readability on the socket are awakened. connstatus is logically ORed with so\_state. This code is never executed for TCP connections, since connstatus is always 0 for TCP.

Protocols, such as TCP, that put incoming connections on so\_q0 first, call soisconnected when the connection establishment phase completes. For TCP, this happens when the second SYN is ACKed on the connection.

Figure 15.30 shows soisconnected.

#### Queue incomplete connections

78-87

The socket state is changed to show that the connection has completed. When soisconnected is called for incoming connections, (i.e., when the local process is calling accept), head is nonnull.

If sogremque returns 1, the socket is queued on so\_q and sorwakeup wakes up any processes using select to monitor the socket for connection arrival by testing for readability. If a process is blocked in accept waiting for the connection, wakeup causes the matching tsleep to return. 464 Socket Layer

Chapter 15

```
uipc_socket2.c
78 soisconnected(so)
79 struct socket *so;
80 {
       struct socket *head = so->so_head;
81
       so->so_state &= ~(SS_ISCONNECTING | SS_ISDISCONNECTING | SS_ISCONFIRMING);
82
       so->so_state |= SS_ISCONNECTED;
83
       if (head && sogremque(so, 0)) {
84
           soginsque(head, so, 1);
85
           sorwakeup(head);
86
           wakeup((caddr_t) & head->so_timeo);
87
88
       } else {
           wakeup((caddr_t) & so->so_timeo);
89
90
           sorwakeup(so);
           sowwakeup(so);
91
92
       }
93 }
                                                                        uipc_socket2.c
```

Figure 15.30 soisconnected function.

## Wakeup processes waiting for new connection

<sup>88–93</sup> If head is null, sogremque is not called since the process initiated the connection with the connect system call and the socket is not on a queue. If head is nonnull and sogremque returns 0, the socket is already on so\_q. This happens with protocols such as TP4, which place connections on so\_q before they are complete. wakeup awakens any process blocked in connect, and sorwakeup and sowwakeup take care of any processes that are using select to wait for the connection to complete.

## 15.13 connect System call

A server process calls the listen and accept system calls to wait for a remote process to initiate a connection. If the process wants to initiate a connection itself (i.e., a client), it calls connect.

For connection-oriented protocols such as TCP, connect establishes a connection to the specified foreign address. The kernel selects and implicitly binds an address to the local socket if the process has not already done so with bind.

For connectionless protocols such as UDP or ICMP, connect records the foreign address for use in sending future datagrams. Any previous foreign address is replaced with the new address.

Figure 15.31 shows the functions called when connect is used for UDP or TCP.

The left side of the figure shows connect processing for connectionless protocols, such as UDP. In this case the protocol layer calls soisconnected and the connect system call returns immediately.

The right side of the figure shows connect processing for connection-oriented protocols, such as TCP. In this case, the protocol layer begins the connection establishment and calls soisconnecting to indicate that the connection will complete some time in the future. Unless the socket is nonblocking, soconnect calls tsleep to wait for the



5

С

Ġ);

С

i l

3

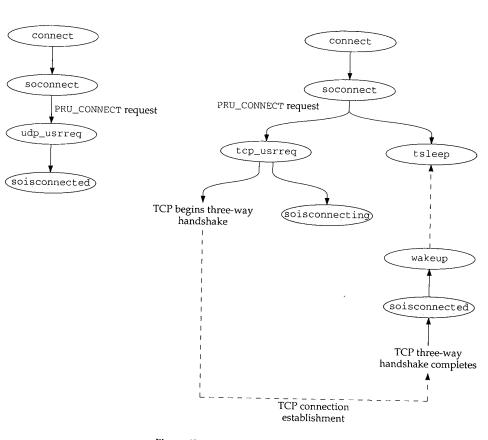


Figure 15.31 connect processing.

connection to complete. For TCP, when the three-way handshake is complete, the protocol layer calls soisconnected to mark the socket as connected and then calls wakeup to awaken the process and complete the connect system call.

Figure 15.32 shows the connect system call.

180-188 The three arguments to connect (in the connect\_args structure) are: s, the socket descriptor; name, a pointer to a buffer containing the foreign address; and namelen, the length of the buffer.

<sup>189–200</sup> get sock returns the socket as usual. A connection request may already be pending on a nonblocking socket, in which case EALREADY is returned. sockargs copies the foreign address from the process into the kernel.

## Start connection processing

201-208

The connection attempt is started by calling soconnect. If soconnect reports an error, connect jumps to bad. If a connection has not yet completed by the time soconnect returns and nonblocking I/O is enabled, EINPROGRESS is returned immediately to avoid waiting for the connection to complete. Since connection establishment

466 Socket Layer

181

180 struct connect\_args {

Chapter 15

uipc\_syscalls.c

- uipc\_syscalls.c

20

Se

21

so

int s; 182 caddr\_t name; 183 int namelen; 184 ); 185 connect(p, uap, retval) 186 struct proc \*p; 187 struct connect\_args \*uap; 188 int \*retval; 189 { 190 struct file \*fp; 191 struct socket \*so; 192 struct mbuf \*nam; 193 int error, s; 194 if (error = getsock(p->p\_fd, uap->s, &fp)) 195 return (error); 196 so = (struct socket \*) fp->f\_data; 197 if ((so->so\_state & SS\_NBIO) && (so->so\_state & SS\_ISCONNECTING)) 198 return (EALREADY); if (error = sockargs(&nam, uap->name, uap->namelen, MT\_SONAME)) 199 200 return (error); 201 error = soconnect(so, nam); 202 if (error) 203 goto bad; 204 if ((so->so\_state & SS\_NBIO) && (so->so\_state & SS\_ISCONNECTING)) { 205 m\_freem(nam); 206 return (EINPROGRESS); 207 } 208 s = splnet(); while ((so->so\_state & SS\_ISCONNECTING) && so->so\_error == 0) 209 210 if (error = tsleep((caddr\_t) & so->so\_timeo, PSOCK | PCATCH, 211 netcon, 0)) 212 break; 213 if (error == 0) { 214 error = so->so\_error; 215  $so->so\_error = 0;$ 216 } 217 splx(s); 218 bad: 219 so->so\_state &= ~SS\_ISCONNECTING; 220 m\_freem(nam); 221 if (error == ERESTART) 222 error = EINTR; 223 return (error); 224 }

Figure 15.32 connect system call.

INTEL Ex.1013.492

normally involves exchanging

normally involves exchanging several packets with the remote system, it may take a while to complete. Further calls to connect return EALREADY until the connection completes. EISCONN is returned when the connection is complete.

### Wait for connection establishment

208-217 The while loop continues until the connection is established or an error occurs. splnet prevents connect from missing a wakeup between testing the state of the socket and the call to tsleep. After the loop, error contains 0, the error code from tsleep, or the error from the socket.

<sup>218–224</sup> The SS\_ISCONNECTING flag is cleared since the connection has completed or the attempt has failed. The mbuf containing the foreign address is released and any error is returned.

#### soconnect Function

This function ensures that the socket is in a valid state for a connection request. If the socket is not connected or a connection is not pending, then the connection request is always valid. If the socket is already connected or a connection is pending, the new connection request is rejected for connection-oriented protocols such as TCP. For connectionless protocols such as UDP, multiple connection requests are OK but each new request replaces the previous foreign address.

Figure 15.33 shows the soconnect function.

```
·uipc_socket.c
198 soconnect(so, nam)
199 struct socket *so;
200 struct mbuf *nam;
201 {
202
        int
                s;
203
        int
                error;
204
        if (so->so_options & SO_ACCEPTCONN)
205
            return (EOPNOTSUPP);
206
        s = splnet();
207
        /*
208
        * If protocol is connection-based, can only connect once.
209
        * Otherwise, if connected, try to disconnect first.
         * This allows user to disconnect by connecting to, e.g.,
210
211
         * a null address.
         */
212
213
        if (so->so_state & (SS_ISCONNECTED | SS_ISCONNECTING) &&
214
            ((so->so_proto->pr_flags & PR_CONNREQUIRED) ||
215
             (error = sodisconnect(so))))
216
            error = EISCONN;
217
        else
218
            error = (*so->so_proto->pr_usrreq) (so, PRU_CONNECT,
219
                                    (struct mbuf *) 0, nam, (struct mbuf *) 0);
220
        splx(s);
221
        return (error);
222 1
```

Figure 15.33 soconnect function.

— uipc\_socket.c

calls.c

alls.c

ter 15

198-222

soconnect returns EOPNOTSUPP if the socket is marked to accept connections, since a process cannot initiate connections if listen has already been called for the socket. EISCONN is returned if the protocol is connection oriented and a connection has already been initiated. For a connectionless protocol, any existing association with a foreign address is broken by sodisconnect.

The PRU\_CONNECT request starts the appropriate protocol processing to establish the connection or the association.

#### **Breaking a Connectionless Association**

For connectionless protocols, the foreign address associated with a socket can be discarded by calling connect with an invalid name such as a pointer to a structure filled with 0s or a structure with an invalid size. sodisconnect removes a foreign address associated with the socket, and PRU\_CONNECT returns an error such as EAFNOSUPPORT or EADDRNOTAVAIL, leaving the socket with no foreign address. This is a useful, although obscure, way of breaking the association between a connectionless socket and a foreign address without replacing it.

## 15.14 shutdown System Call

The shutdown system call, shown in Figure 15.34, closes the write-half, read-half, or both halves of a connection. For the read-half, shutdown discards any data the process hasn't yet read and any data that arrives after the call to shutdown. For the write-half, shutdown lets the protocol specify the semantics. For TCP, any remaining data will be sent followed by a FIN. This is TCP's half-close feature (Section 18.5 of Volume 1).

To destroy the socket and release the descriptor, close must be called. close can also be called directly without first calling shutdown. As with all descriptors, close is called by the kernel for sockets that have not been closed when a process terminates.

```
uipc_syscalls.c
550 struct shutdown_args {
551
        int
                 s;
552
        int
                 how;
553 };
554 shutdown(p, uap, retval)
555 struct proc *p;
556 struct shutdown_args *uap;
557 int
           *retval;
558 {
559
        struct file *fp;
560
        int
                 error;
561
        if (error = getsock(p->p_fd, uap->s, &fp))
562
            return (error);
563
        return (soshutdown((struct socket *) fp->f data, uap->how));
564 }
```

– uipc\_syscalls.c

Figure 15.34 shutdown system call.

15

s,

۱e

is a

;h

s-:d

35 .T

ıI,

d

r

3S

f,

۰e

n

is

.С

С

<sup>550–557</sup> In the shutdown\_args structure, s is the socket descriptor and how specifies which halves of the connection are to be closed. Figure 15.35 shows the expected values for how and how++ (which is used in Figure 15.36).

how	how++	Description
0	FREAD	shut down the read-half of the connection
1	FWRITE	shut down the write-half of the connection
2	FREAD FWRITE	shut down both halves of the connection

Figure 15.35 shutdown system call options.

Notice that there is an implicit numerical relationship between how and the constants  $\ensuremath{\mathsf{FREAD}}$  and  $\ensuremath{\mathsf{FWRITE}}$ .

558-564 shutdown is a wrapper function for soshutdown. The socket associated with the descriptor is returned by getsock, soshutdown is called, and its value is returned.

## soshutdown and sorflush Functions

The shut down of the read-half of a connection is handled in the socket layer by sorflush, and the shut down of the write-half of a connection is processed by the PRU\_SHUTDOWN request in the protocol layer. The soshutdown function is shown in Figure 15.36.

uipc\_socket.c 720 soshutdown(so, how) 721 struct socket \*so; 722 int how: 723 { 724 struct protosw \*pr = so->so\_proto; 725 how++; 726 if (how & FREAD) 727 sorflush(so); 728 if (how & FWRITE) 729 return ((\*pr->pr\_usrreq) (so, PRU\_SHUTDOWN, 730 (struct mbuf \*) 0, (struct mbuf \*) 0, (struct mbuf \*) 0)); 731 return (0); 732 }

Figure 15.36 soshutdown function.

720-732 If the read-half of the socket is being closed, sorflush, shown in Figure 15.37, discards the data in the socket's receive buffer and disables the read-half of the connection. If the write-half of the socket is being closed, the PRU\_SHUTDOWN request is issued to the protocol.

733-747 The process waits for a lock on the receive buffer. Because of SB\_NOINTR, sblock does not return when an interrupt occurs. splimp blocks network interrupts and protocol processing while the socket is modified, since the receive buffer may be accessed by the protocol layer as it processes incoming packets.

uipc socket.c

Chapter 15

「ここの学術の問題

and the second se

「日本の読書」

uipc\_socket.c

```
uipc_socket.c
733 sorflush(so)
734 struct socket *so;
735 {
        struct sockbuf *sb = &so->so_rcv;
736
        struct protosw *pr = so->so_proto;
737
738
        int
                s:
        struct sockbuf asb;
739
        sb->sb_flags |= SB_NOINTR;
740
        (void) sblock(sb, M_WAITOK);
741
742
        s = splimp();
        socantrcvmore(so);
743
744
        sbunlock(sb);
745
        asb = *sb;
        bzero((caddr_t) sb, sizeof(*sb));
746
747
        splx(s);
        if (pr->pr_flags & PR_RIGHTS && pr->pr_domain->dom_dispose)
748
             (*pr->pr_domain->dom_dispose) (asb.sb_mb);
749
        sbrelease(&asb);
750
751 }
```

### Figure 15.37 sorflush function.

socantrcvmore marks the socket to reject incoming packets. A copy of the sockbuf structure is saved in asb to be used after interrupts are restored by splx. The original sockbuf structure is cleared by bzero, so that the receive queue appears to be empty.

#### **Release control mbufs**

748-751

Some kernel resources may be referenced by control information present in the receive queue when shutdown was called. The mbuf chain is still available through sb\_mb in the copy of the sockbuf structure.

If the protocol supports access rights and has registered a dom\_dispose function, it is called here to release these resources.

> In the Unix domain it is possible to pass descriptors between processes with control messages. These messages contain pointers to reference counted data structures. The dom\_dispose function takes care of discarding the references and the data structures if necessary to avoid creating an unreferenced structure and introducing a memory leak in the kernel. For more information on passing file descriptors within the Unix domain, see [Stevens 1990] and [Leffler et al. 1989].

Any input data pending when shutdown is called is discarded when sbrelease releases any mbufs on the receive queue.

Notice that the shut down of the read-half of the connection is processed entirely by the socket layer (Exercise 15.6) and the shut down of the write-half of the connection is handled by the protocol through the PRU\_SHUTDOWN request. TCP responds to the PRU\_SHUTDOWN by sending all queued data and then a FIN to close the write-half of the TCP connection.

471 close System Call

Section 15.15

## 15.15 close System Call

The close system call works with any type of descriptor. When fd is the last descriptor that references the object, the object-specific close function is called:

error = (\*fp->f\_ops->fo\_close)(fp, p);

As shown in Figure 15.13, fp->f\_ops->fo\_close for a socket is the function soo\_close.

#### soo\_close Function

This function, shown in Figure 15.38, is a wrapper for the soclose function.

```
152 soo_close(fp, p)
153 struct file *fp;
154 struct proc *p;
155 {
156
        int
                 error = 0;
157
        if (fp->f_data)
             error = soclose((struct socket *) fp->f_data);
158
159
        fp \rightarrow f data = 0:
160
        return (error);
161 }
```

-sys\_socket.c

- sys\_socket.c

If a socket structure is associated with the file structure, soclose is called, 152-161 f\_data is cleared, and any posted error is returned.

Figure 15.38 soo\_close function.

#### soclose Function

This function aborts any connections that are pending on the socket (i.e., that have not yet been accepted by a process), waits for data to be transmitted to the foreign system, and releases the data structures that are no longer needed.

soclose is shown in Figure 15.39.

#### **Discard pending connections**

129-141

If the socket was accepting connections, soclose traverses the two connection queues and calls soabort for each pending connection. If the protocol control block is null, the protocol has already been detached from the socket and soclose jumps to the cleanup code at discard.

> soabort issues the PRU\_ABORT request to the socket's protocol and returns the result. soabort is not shown in this text. Figures 23.38 and 30.7 discuss how UDP and TCP handle this request.

472 Socket Layer

Chapter 15

```
uipc_socket.c
129 soclose(so)
130 struct socket *so;
131 {
                                     /* conservative */
                s = splnet();
132
        int
                error = 0;
133
        int
        if (so->so_options & SO_ACCEPTCONN) {
134
            while (so->so_q0)
135
                 (void) soabort(so->so_q0);
136
            while (so->so_q)
137
                 (void) soabort(so->so_q);
138
139
        }
        if (so->so_pcb == 0)
140
            goto discard;
141
        if (so->so_state & SS_ISCONNECTED) {
142
             if ((so->so_state & SS_ISDISCONNECTING) == 0) {
143
                 error = sodisconnect(so);
144
                 if (error)
145
                     goto drop;
146
147
             }
             if (so->so_options & SO_LINGER) {
148
                 if ((so->so_state & SS_ISDISCONNECTING) &&
149
                     (so->so_state & SS_NBIO))
150
151
                     goto drop;
                 while (so->so_state & SS_ISCONNECTED)
152
                     if (error = tsleep((caddr_t) & so->so_timeo,
153
                                         PSOCK | PCATCH, netcls, so->so_linger))
 154
                          break;
 155
 156
             }
         }
 157
 158
       drop:
         if (so->so_pcb) {
 159
                     error2 =
 160
             int
             (*so->so_proto->pr_usrreq) (so, PRU_DETACH,
 161
                      (struct mbuf *) 0, (struct mbuf *) 0, (struct mbuf *) 0);
 162
             if (error == 0)
 163
                 error = error2;
 164
 165
         }
       discard:
 166
         if (so->so_state & SS_NOFDREF)
 167
             panic("soclose: NOFDREF");
 168
         so->so_state != SS_NOFDREF;
 169
 170
         sofree(so);
 171
         splx(s);
 172
         return (error);
 173 }
                                                                          uipc_socket.c
```

Figure 15.39 soclose function.

「「「「「

## Break established connection or association

142-157

If the socket is not connected, execution continues at drop; otherwise the socket must be disconnected from its peer. If a disconnect is not in progress, sodisconnect starts the disconnection process. If the SO\_LINGER socket option is set, soclose may need to wait for the disconnect to complete before returning. A nonblocking socket never waits for a disconnect to complete, so soclose jumps immediately to drop in that case. Otherwise, the connection termination is in progress and the SO\_LINGER option indicates that soclose must wait some time for it to complete. The while loop continues until the disconnect completes, the linger time (so\_linger) expires, or a signal is delivered to the process.

> If the linger time is set to 0, tsleep returns only when the disconnect completes (perhaps because of an error) or a signal is delivered.

#### **Release data structures**

If the socket still has an attached protocol, the PRU\_DETACH request breaks the con-158-173 nection between this socket and the protocol. Finally the socket is marked as not having an associated file descriptor, which allows sofree to release the socket.

The sofree function is shown in Figure 15.40.

```
uipc_socket.c
110 `sofree(so)
111 struct socket *so;
112 {
        if (so->so_pcb || (so->so_state & SS_NOFDREF) == 0)
113
            return;
114
115
        if (so->so head) {
            if (!sogremque(so, 0) && !sogremque(so, 1))
116
117
                panic("sofree dq");
            so->so_head = 0;
118
119
        }
120
        sbrelease(&so->so_snd);
121
        sorflush(so);
        FREE (so, M_SOCKET);
122
123 }

    uipc_socket.c
```

Figure 15.40 sofree function.

### Return if socket still in use

If a protocol is still associated with the socket, or if the socket is still associated with 110-114 a descriptor, sofree returns immediately.

#### **Remove from connection queues**

115-119

If the socket is on a connection queue (so\_head is nonnull), sogremque is called to remove the socket. An attempt is made to remove the socket from the incomplete connection queue and if this fails, then from the completed connection queue. One of the removals must succeed or the kernel panics, since so\_head was nonnull. so\_head is cleared.

## Discard send and receive queues

120-123 sbrelease discards any buffers in the send queue and sorflush discards any buffers in the receive queue. Finally, the socket itself is released.

## 15.16 Summary

In this chapter we looked at all the system calls related to network operations. The system call mechanism was described, and we traced the calls until they entered the protocol processing layer through the pr\_usrreg function.

While looking at the socket layer, we avoided any discussion of address formats, protocol semantics, or protocol implementations. In the upcoming chapters we tie together the link-layer processing and socket-layer processing by looking in detail at the implementation of the Internet protocols in the protocol processing layer.

## **Exercises**

- **15.1** How can a process *without* superuser privileges gain access to a socket created by a superuser process?
- **15.2** How can a process determine if the sockaddr buffer it provides to accept was too small to hold the foreign address returned by the call?
- 15.3 A feature proposed for IPv6 sockets is to have accept and recvfrom return a source route as an array of 128-bit IPv6 addresses instead of a single peer address. Since the array will not fit in a single mbuf, modify accept and recvfrom to handle an mbuf chain from the protocol layer instead of a single mbuf. Will the existing code work if the protocol layer returns the array in an mbuf cluster instead of a chain of mbufs?
- 15.4 Why is panic called when sogremque returns a null pointer in Figure 15.26?
- 15.5 Why does sorflush make a copy of the receive buffer?
- **15.6** What happens when additional data is received after sorflush has zeroed the socket's receive buffer? Read Chapter 16 before attempting this exercise.

1

国大学学校研究の行

# 16

# Socket I/O

## 16.1 Introduction

In this chapter we discuss the system calls that read and write data on a network connection. The chapter is divided into three parts.

The first part covers the four system calls for sending data: write, writev, sendto, and sendmsg. The second part covers the four system calls for receiving data: read, readv, recvfrom, and recvmsg. The third part of the chapter covers the select system call, which provides a standard way to monitor the status of descriptors in general and sockets in particular.

The core of the socket layer is the sosend and soreceive functions. They handle all I/O between the socket layer and the protocol layer. As we'll see, the semantics of the various types of protocols overlap in these functions, making the functions long and complex.

## 16.2 Code Introduction

The three headers and four C files listed in Figure 16.1 are covered in this chapter.

## **Global Variables**

The first two global variables shown in Figure 16.2 are used by the select system call. The third global variable controls the amount of memory allocated to a socket.

475

Socket I/O 476

File	Description
sys/socket.h sys/socketvar.h sys/uio.h	structures and macro for sockets API socket structure and macros uio structure definition
<pre>kern/uipc_syscalls.c kern/uipc_socket.c kern/sys_generic.c kern/sys_socket.c</pre>	socket system calls socket layer processing select system call select processing for sockets

Figure 16.1 Files discussed in this chapter.

Variable	Datatype	Description
selwait	int	wait channel for select
nselcoll	int	flag used to avoid race conditions in select
sb_max	u_long	maximum number of bytes to allocate for a socket receive or send buffer

Figure 16.2 Global variables introduced in this chapter.

#### 16.3 Socket Buffers

Section 15.3 showed that each socket has an associated send and receive buffer. The sockbuf structure definition from Figure 15.5 is repeated in Figure 16.3.

		socketvar.h
72	struct sockbuf {	
73	u_long sb_cc;	/* actual chars in buffer */
74	u_long sb_hiwat;	/* max actual char count */
75	u_long sb_mbcnt;	/* chars of mbufs used */
76	u_long sb_mbmax;	/* max chars of mbufs to use */
77	long sb_lowat;	/* low water mark */
78	<pre>struct mbuf *sb_mb;</pre>	/* the mbuf chain */
79	<pre>struct selinfo sb_sel;</pre>	<pre>/* process selecting read/write */</pre>
80	<pre>short sb_flags;</pre>	/* Figure 16.5 */
81	<pre>short sb_timeo;</pre>	/* timeout for read/write */
82	} so_rcv, so_snd;	socketvar.h

Figure 16.3 sockbuf structure.

Each buffer contains control information as well as pointers to data stored in mbuf 72-78 chains. sb\_mb points to the first mbuf in the chain, and sb\_cc is the total number of data bytes contained within the mbufs. sb\_hiwat and sb\_lowat regulate the socket flow control algorithms. sb\_mbcnt is the total amount of memory allocated to the mbufs in the buffer.

Recall that each mbuf may store from 0 to 2048 bytes of data (if an external cluster is used). sb\_mbmax is an upper bound on the amount of memory to be allocated as

INTEL Ex.1013.502

S

79

80

mbufs for each socket buffer. Default limits are specified by each protocol when the PRU\_ATTACH request is issued by the socket system call. The high-water and low-water marks may be modified by the process as long as the kernel-enforced hard limit of 262,144 bytes per socket buffer (sb\_max) is not exceeded. The buffering algorithms are described in Sections 16.7 and 16.12. Figure 16.4 shows the default settings for the Internet protocols.

Protocol	so_snd			so_rcv		
	sb_hiwat	sb_lowat	sb_mbmax	sb_hiwat	sb_lowat	sb_mbmax
UDP TCP	9×1024 8×1024	2048 (ignored) 2048	2×sb_hiwat 2×sb_hiwat	$40 \times (1024 + 16)$ $8 \times 1024$	1 1	2×sb_hiwat 2×sb hiwat
raw IP ICMP IGMP	8×1024	2048 (ignored)	2×sb_hiwat	8×1024	1	2×sb_hiwat

Figure 16.4 Default socket buffer limits for the Internet protocols.

Since the source address of each incoming UDP datagram is queued with the data (Section 23.8), the default UDP value for sb\_hiwat is set to accommodate 40 1K datagrams and their associated sockaddr\_in structures (16 bytes each).

sb\_sel is a selinfo structure used to implement the select system call (Section 16.13).

Figure 16.5 lists the possible values for sb\_flags.

sb_flags	Description
SB_LOCK	a process has locked the socket buffer
SB_WANT	a process is waiting to lock the buffer
SB_WAIT	a process is waiting for data (receive) or space (send) in this buffer
SB_SEL	one or more processes are selecting on this buffer
SB_ASYNC	generate asynchronous I/O signal for this buffer
SB_NOINTR	signals do not cancel a lock request
SB_NOTIFY	(SB_WAIT SB_SEL SB_ASYNC)
	a process is waiting for changes to the buffer and should be notified by wakeup when any changes occur

Figure 16.5 sb\_flags values.

\$1-82 sb\_timeo is measured in clock ticks and limits the time a process blocks during a read or write call. The default value of 0 causes the process to wait indefinitely. sb\_timeo may be changed or retrieved by the SO\_SNDTIMEO and SO\_RCVTIMEO socket options.

## Socket Macros and Functions

There are many macros and functions that manipulate the send and receive buffers associated with each socket. The macros and functions in Figure 16.6 handle buffer locking and synchronization.

Chapter 16

Name	Description
sblock	Acquires a lock for sb. If wf is M_WAITOK, the process sleeps waiting for the lock; otherwise EWOULDBLOCK is returned if the buffer cannot be locked immediately. EINTR or ERESTART is returned if the sleep is interrupted by a signal; 0 is returned otherwise.
	<pre>int <b>sblock</b>(struct sockbuf *sb, int wf);</pre>
sbunlock	Releases the lock on sb. Any other process waiting to lock sb is awakened.
	void <b>sbunlock</b> (struct sockbuf *sb);
sbwait	Calls tsleep to wait for protocol activity on sb. Returns result of tsleep.
	<pre>int sbwait(struct sockbuf *sb);</pre>
sowakeup	Notifies socket of protocol activity. Wakes up matching call to sbwait or to tsleep if any processes are selecting on <i>sb</i> .
	<pre>void sowakeup(struct socket *so, struct sockbuf *sb);</pre>
sorwakeup	Wakes up any process waiting for read events on <i>so</i> and sends the SIGIO signal if a process requested asynchronous notification of I/O.
	void <b>sorwakeup</b> (struct socket * <i>so</i> );
sowwakeup	Wakes up any process waiting for write events on <i>so</i> and sends the SIGIO signal if a process requested asynchronous notification of I/O.
	void <b>sowwakeup</b> (struct socket *50);

Figure 16.6 Macros and functions for socket buffer locking and synchronization.

Figure 16.7 includes the macros and functions used to set the resource limits for socket buffers and to append and delete data from the buffers. In the table, m, m0, n, and *control* are all pointers to mbuf chains. *sb* points to the send or receive buffer for a socket.

Name	Description
sbspace	The number of bytes that may be added to <i>sb</i> before it is considered full: min((sb_hiwat - sb_cc), (sb_mbmax - sb_mbcnt)).
	long <b>sbspace</b> (struct sockbuf *sb);
sballoc	<pre>m has been added to sb. Adjust sb_cc and sb_mbcnt in sb accordingly. void sballoc(struct sockbuf *sb, struct mbuf *m);</pre>
sbfree	<pre>m has been removed from sb. Adjust sb_cc and sb_mbcnt in sb accordingly. int <b>sbfree</b>(struct sockbuf *sb, struct mbuf *m);</pre>

(10 m) (10 m)

4

うちょう 日本語のの 学校の いままの

Name	Description
sbappend	Append the mbufs in $m$ to the end of the last record in $sb$ . Call sbcompress.
	<pre>int sbappend(struct sockbuf *sb, struct mbuf *m);</pre>
sbappendrecord	Append the record in <i>ni0</i> after the last record in <i>sb</i> . Call sbcompress.
	<pre>int sbappendrecord(struct sockbuf *sb, struct mbuf *m0);</pre>
sbappendaddr	Put address from <i>asa</i> in an mbuf. Concatenate address, <i>control</i> , and <i>m0</i> . Append the resulting mbuf chain after the last record in <i>sb</i> .
	int <b>sbappendaddr</b> (struct sockbuf * <i>sb</i> , struct sockaddr * <i>asa</i> , struct mbuf * <i>m0</i> , struct mbuf * <i>control</i> );
sbappendcontrol	Concatenate <i>control</i> and <i>m0</i> . Append the resulting mbuf chain after the last record in <i>sb</i> .
	<pre>int sbappendcontrol(struct sockbuf *sb, struct mbuf *m0,</pre>
sbinsertoob	Insert m0 before first record in sb without out-of-band data. Call sbcompress.
	<pre>int sbinsertoob(struct sockbuf *sb, struct mbuf *m0);</pre>
sbcompress	Append <i>m</i> to <i>n</i> squeezing out any unused space.
	<pre>void sbcompress(struct sockbuf *sb, struct mbuf *m, struct mbuf *n);</pre>
sbdrop	Discard len bytes from the front of sb.
	<pre>void sbdrop(struct sockbuf *sb, intlen);</pre>
sbdroprecord	Discard the first record in <i>sb</i> . Move the next record to the front.
	void <b>sbdroprecord</b> (struct sockbuf * <i>sb</i> );
sbrelease	Call sbflush to release all mbufs in <i>sb</i> . Reset sb_hiwat and sb_mbmax values to 0.
	void <b>sbrelease</b> (struct sockbuf * <i>sb</i> );
sbflush	Release all mbufs in <i>sb</i> .
ĺ	void <b>sbflush</b> (struct sockbuf * <i>sb</i> );
soreserve	Set high-water and low-water marks. For the send buffer, call sbreserve with <i>sndcc</i> . For the receive buffer, call sbreserve with <i>rcvcc</i> . Initialize sb_lowat in both buffers to default values, Figure 16.4. ENOBUFS is returned if any limits are exceeded.
	<pre>int soreserve(struct socket *so, int sndcc, int rcvcc);</pre>
sbreserve	Set high-water mark for <i>sb</i> to <i>cc</i> . Also drop low-water mark to <i>cc</i> . No memory is allocated by this function.
	<pre>int sbreserve(struct sockbuf *sb, int cc);</pre>

Figure 16.7 Macros and functions for socket buffer allocation and manipulation.

Se

# 16.4 write, writev, sendto, and sendmsg System Calls

These four system calls, which we refer to collectively as the *write system calls*, send data on a network connection. The first three system calls are simpler interfaces to the most general request, sendmsg.

All the write system calls, directly or indirectly, call sosend, which does the work of copying data from the process to the kernel and passing data to the protocol associated with the socket. Figure 16.8 summarizes the flow of control.

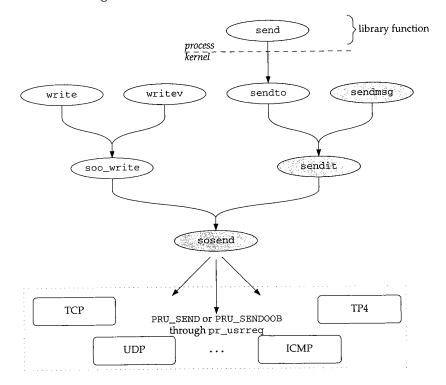


Figure 16.8 All socket output is handled by sosend.

In the following sections, we discuss the functions shaded in Figure 16.8. The other four system calls and soo\_write are left for readers to investigate on their own. Figure 16.9 shows the features of these four system calls and a related library function (send).

In Net/3, send is implemented as a library function that calls sendto. For binary compatibility with previously compiled programs, the kernel maps the old send system call to the function osend, which is not discussed in this text.

From the second column in Figure 16.9 we see that the write and writev system calls are valid with any descriptor, but the remaining system calls are valid only with socket descriptors.

### INTEL Ex.1013.506

l6

ta st

:k 'i-

er

<u>۲</u>-

il-IC-

m

:h

## write, writev, sendto, and sendmsg System Calls

481

Function	Type of descriptor	Number of buffers	Specify destination address?	Flags?	Control information?
write	any	1			
writev	any	[1UIO_MAXIOV]			
send	socket only	1		í .	
sendto	socket only	1	•	•	
sendmsg	socket only	[1UIO_MAXIOV]	•	•	•

## Figure 16.9 Write system calls.

The third column shows that writev and sendmsg accept data from multiple buffers. Writing from multiple buffers is called gathering. The analogous read operation is called scattering. In a gather operation the kernel accepts, in order, data from each buffer specified in an array of iovec structures. The array can have a maximum of UIO\_MAXIOV elements. The structure is shown in Figure 16.10.

41 st: 42 43 44 };	ruct iovec { char *iov_base; size_t iov_len;	/* Base address */ /* Length */
<u> </u>		uio.h

Figure 16.10 iovec structure.

iov\_base points to the start of a buffer of iov\_len bytes. 41 - 44

Without this type of interface, a process would have to copy buffers into a single larger buffer or make multiple write system calls to send data from multiple buffers. Both alternatives are less efficient than passing an array of iovec structures to the kernel in a single call. With datagram protocols, the result of one writev is one datagram, which cannot be emulated with multiple writes.

Figure 16.11 illustrates the structures as they are used by writev, where iovp points to the first element of the array and iovent is the size of the array.

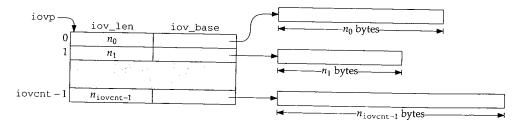


Figure 16.11 iovec arguments to writev.

Datagram protocols require a destination address to be associated with each write call. Since write, writev, and send do not accept an explicit destination, they may be called only after a destination has been associated with a connectionless socket by calling connect. A destination must be provided with sendto or sendmsg, or connect must have been previously called.

「「「「「「「」」」

socket.h

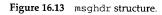
The fifth column in Figure 16.9 shows that the sendxxx system calls accept optional control flags, which are described in Figure 16.12.

flags	Description	Reference
MSG_DONTROUTE	bypass routing tables for this message	Figure 16.23
MSG_DONTWAIT	do not wait for resources during this message	Figure 16.22
MSG_EOR	data marks the end of a logical record	Figure 16.25
MSG_OOB	send as out-of-band data	Figure 16.26

Figure 16.12	sendxxx svs	tem calls:	flags values.
--------------	-------------	------------	---------------

As indicated in the last column of Figure 16.9, only the sendmsg system call supports control information. The control information and several other arguments to sendmsg are specified within a msghdr structure (Figure 16.13) instead of being passed separately.

	uct msghdr {		— socket.h
229	caddr_t msg_name;	/* optional address */	
230	u_int msg_namelen;	/* size of address */	
231	<pre>struct iovec *msg_iov;</pre>	/* scatter/gather array */	
232	u_int msg_iovlen;	/* # elements in msg_iov */	
233	<pre>caddr_t msg_control;</pre>	/* ancillary data, see below */	
234	u_int msg_controllen;	/* ancillary data buffer len */	
235	int msg_flags;	/* Figure 16.33 */	
236 };			



 ${\tt msg\_name}$  should be declared as a pointer to a <code>sockaddr</code> structure, since it contains a network address.

228-236 The msghdr structure contains a destination address (msg\_name and msg\_namelen), a scatter/gather array (msg\_iov and msg\_iovlen), control information (msg\_control and msg\_controllen), and receive flags (msg\_flags). The control information is formatted as a cmsghdr structure shown in Figure 16.14.

251 st	ruct cmsgl	hdr {	socket.h
	u_int int int followed	<pre>cmsg_len; cmsg_level; cmsg_type; by u_char</pre>	<pre>/* data byte count, including hdr */     /* originating protocol */     /* protocol-specific type */ cmsg_data[]; */</pre>
256 };			

Figure 16.14 cmsghdr structure.

<sup>251-256</sup> The control information is not interpreted by the socket layer, but the messages are typed (cmsg\_type) and they have an explicit length (cmsg\_len). Multiple control messages may appear in the control information mbuf.

## Example

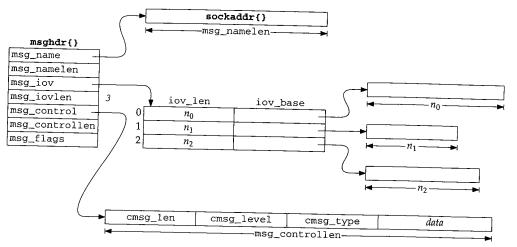


Figure 16.15 shows how a fully specified msghdr structure might look during a call to sendmsg.

Figure 16.15 msghdr structure for sendmsg system call.

# 16.5 sendmsg System Call

Only the sendmsg system call provides access to all the features of the sockets API associated with output. The sendmsg and sendit functions prepare the data structures needed by sosend, which passes the message to the appropriate protocol. For SOCK\_DGRAM protocols, a message is a datagram. For SOCK\_STREAM protocols, a message is a sequence of bytes. For SOCK\_SEQPACKET protocols, a message could be an entire record (implicit record boundaries) or part of a larger record (explicit record boundaries). A message is always an entire record (implicit record boundaries) for SOCK\_RDM protocols.

Even though the general sosend code handles SOCK\_SEQPACKET and SOCK\_RDM protocols, there are no such protocols in the Internet domain.

Figure 16.16 shows the sendmsg code.

<sup>307-321</sup> There are three arguments to sendmsg: the socket descriptor; a pointer to a msghdr structure; and several control flags. The copyin function copies the msghdr structure from user space to the kernel.

## Copy iov array

An iovec array with eight entries (UIO\_SMALLIOV) is allocated automatically on the stack. If this is not large enough, sendmsg calls MALLOC to allocate a larger array. If

484 Socket I/O

Chapter 16

uipc\_syscalls.c 307 struct sendmsg\_args { 308 int s; 309 caddr\_t msg; 310 flags; int 311 }; 312 sendmsg(p, uap, retval) 313 struct proc \*p; 314 struct sendmsg\_args \*uap; 315 int \*retval; 316 { 317 struct msghdr msg; struct iovec aiov[UIO\_SMALLIOV], \*iov; 318 319 int error; if (error = copyin(uap->msg, (caddr\_t) & msg, sizeof(msg))) 320 321 return (error); if ((u\_int) msg.msg\_iovlen >= UIO\_SMALLIOV) { 322 323 if ((u\_int) msg.msg\_iovlen >= UIO\_MAXIOV) 324 return (EMSGSIZE); 325 MALLOC(iov, struct iovec \*, sizeof(struct iovec) \* (u\_int) msg.msg\_iovlen, M\_IOV, 326 327 M WAITOK); 328 } else 329 iov = aiov; 330 if (msg.msg\_iovlen && (error = copyin((caddr\_t) msg.msg\_iov, (caddr\_t) iov, 331 (unsigned) (msg.msg\_iovlen \* sizeof(struct iovec))))) 332 333 goto done; 334 msg.msg\_iov = iov; 335 error = sendit(p, uap->s, &msg, uap->flags, retval); 336 done: 337 if (iov != aiov) 338 FREE(iov, M\_IOV); 339 return (error); 340 } - uipc\_syscalls.c

Figure 16.16 sendmsg system call.

the process specifies an array with more than 1024 (UIO\_MAXIOV) entries, EMSGSIZE is returned. copyin places a copy of the iovec array from user space into either the array on the stack or the larger, dynamically allocated, array.

This technique avoids the relatively expensive call to malloc in the most common case of eight or fewer entries.

#### sendit and cleanup

335-340

When sendit returns, the data has been delivered to the appropriate protocol or an error has occurred. sendmsg releases the iovec array (if it was dynamically allocated) and returns sendit's result.

- uio.h

## 16.6 sendit Function

sendit is the common function called by sendto and sendmsg. sendit initializes a uio structure and copies control and address information from the process into the kernel. Before discussing sosend, we must explain the uiomove function and the uio structure.

## uiomove Function

The prototype for this function is:

int uiomove(caddr\_t cp, int n, struct uio \*uio);

The uiomove function moves n bytes between a single buffer referenced by cp and the multiple buffers specified by an *iovec* array in *uio*. Figure 16.17 shows the definition of the uio structure, which controls and records the actions of the uiomove function.

```
uio.h
45 enum uio_rw {
46
       UIO_READ, UIO_WRITE
47 };
48 enum uio_seg {
                                   /* Segment flag values */
       UIO_USERSPACE,
49
                                    /* from user data space */
50
       UIO_SYSSPACE,
                                    /* from system space */
51
       UIO_USERISPACE
                                    /* from user instruction space */
52 };
53 struct uio {
       struct iovec *uio_iov;  /* an array of iovec structures */
int uio_iovcnt;  /* size of iovec array */
54
55
      off_t uio_offset;
int uio_resid;
56
                                   /* starting position of transfer */
57
                                  /* remaining bytes to transfer */
58
      enum uio_seg uio_segflg; /* location of buffers */
59
      enum uio_rw uio_rw;
                                    /* direction of transfer */
       struct proc *uio_procp;
60
                                   /* the associated process */
61 };
```



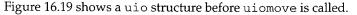
In the uio structure, uio\_iov points to an array of iovec structures, uio\_offset counts the number of bytes transferred by uiomove, and uio\_resid counts the number of bytes remaining to be transferred. Each time uiomove is called, uio\_offset increases by *n* and uio\_resid decreases by *n*. uiomove adjusts the base pointers and buffer lengths in the uio\_iov array to exclude any bytes that uiomove transfers each time it is called. Finally, uio\_iov is advanced through each entry in the array as each buffer is transferred. uio\_segflg indicates the location of the buffers specified by the base pointers in the uio\_iov array and uio\_rw indicates the direction of the transfer. The buffers may be located in the user data space, user instruction space, or kernel data space. Figure 16.18 summarizes the operation of uiomove. The descriptions use the argument names shown in the uiomove prototype.

Chapter 16

uio_segflg	uio_rw	Description	
UIO_USERSPACE	UIO READ	scatter <i>n</i> bytes from a kernel buffer <i>cp</i> to process	
UIO_USERISPACE	010_READ	buffers	
UIO_USERSPACE	UIO WRITE	gather n bytes from process buffers into the kernel	
UIO_USERISPACE	010_WRITE	buffer cp	
UIO_SYSSPACE	UIO_READ	scatter <i>n</i> bytes from the kernel buffer <i>cp</i> to multiple kernel buffers	
010_5155FACE	UIO_WRITE	gather <i>n</i> bytes from multiple kernel buffers into the kernel buffer <i>cp</i>	

Figure 16.18 uiomove operation.

## Example



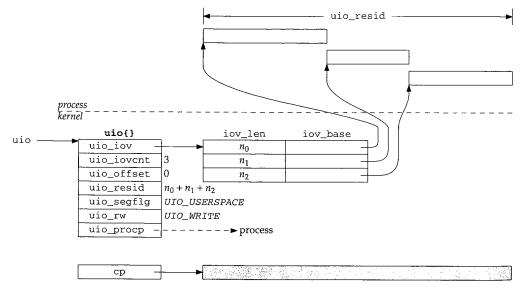


Figure 16.19 uiomove: before.

uio\_iov points to the first entry in the iovec array. Each of the iov\_base pointers point to the start of their respective buffer in the address space of the process. uio\_offset is 0, and uio\_resid is the sum of size of the three buffers. cp points to a buffer within the kernel, typically the data area of an mbuf. Figure 16.20 shows the same data structures after

uiomove(cp, n, uio);

· · · · · · · ·

is executed where n includes all the bytes from the first buffer and only some of the bytes from the second buffer (i.e.,  $n_0 < n < n_0 + n_1$ ).

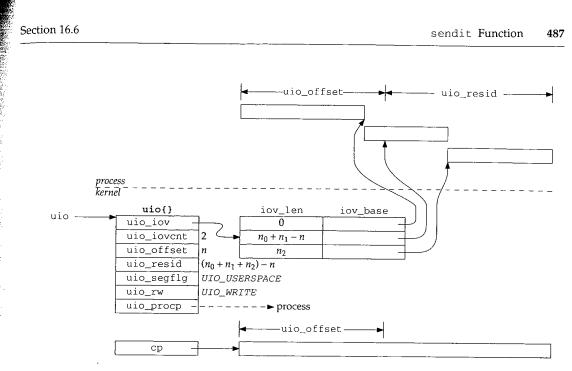


Figure 16.20 uiomove: after.

After uiomove, the first buffer has a length of 0 and its base pointer has been advanced to the end of the buffer. uio\_iov now points to the second entry in the iovec array. The pointer in this entry has been advanced and the length decreased to reflect the transfer of some of the bytes in the buffer. uio\_offset has been increased by *n* and uio\_resid has been decreased by *n*. The data from the buffers in the process has been moved into the kernel's buffer because uio\_rw was UIO\_WRITE.

#### sendit Code

۱t-

35.

) a

he

he

We can now discuss the sendit code shown in Figure 16.21.

## Initialize auio

341-368 sendit calls getsock to get the file structure associated with the descriptor s and initializes the uio structure to gather the output buffers specified by the process into mbufs in the kernel. The length of the transfer is calculated by the for loop as the sum of the buffer lengths and saved in uio\_resid. The first if within the loop ensures that the buffer length is nonnegative. The second if ensures that uio\_resid does not overflow, since uio\_resid is a signed integer and iov\_len is guaranteed to be nonnegative.

## Copy address and control information from the process

<sup>369–385</sup> sockargs makes copies of the destination address and control information into mbufs if they are provided by the process. 488 Socket I/O

Chapter 16

S. Service

小地部後に

の日本のないないとなった

ここの時間、湯いん いちかん 読い

uipc\_syscalls.c 341 sendit(p, s, mp, flags, retsize) 342 struct proc \*p; 343 int s; 344 struct msghdr \*mp; flags, \*retsize; 345 int 346 { struct file \*fp; 347 struct uio auio; 348 struct iovec \*iov; 349 350 int i; struct mbuf \*to, \*control; 351 352 int len, error; if (error = getsock(p->p\_fd, s, &fp)) 353 return (error); 354 auio.uio\_iov = mp->msg\_iov; 355 auio.uio\_iovcnt = mp->msg\_iovlen; 356 auio.uio\_segflg = UIO\_USERSPACE; 357 auio.uio\_rw = UIO\_WRITE; 358 359 auio.uio\_procp = p; /\* XXX \*/ 360 auio.uio\_offset = 0; auio.uio\_resid = 0; 361 iov = mp->msg\_iov; 362 for (i = 0; i < mp->msg\_iovlen; i++, iov++) { 363 if (iov->iov\_len < 0) 364 return (EINVAL); 365 if ((auio.uio\_resid += iov->iov\_len) < 0)</pre> 366 return (EINVAL); 367 368 } 369 if (mp->msg\_name) { if (error = sockargs(&to, mp->msg\_name, mp->msg\_namelen, 370 MT\_SONAME)) 371 return (error); 372 } else 373 374 to = 0;if (mp->msg\_control) { 375 if (mp->msg\_controllen < sizeof(struct cmsghdr) 376 377 ) { error = EINVAL; 378 goto bad; 379 380 if (error = sockargs(&control, mp->msg\_control, 381 mp->msg\_controllen, MT\_CONTROL)) 382 goto bad; 383 384 } else control = 0;385 len = auio.uio\_resid; 386 if (error = sosend((struct socket \*) fp->f\_data, to, &auio, 387 (struct mbuf \*) 0, control, flags)) { 388 if (auio.uio\_resid != len && (error == ERESTART || 389 error == EINTR || error == EWOULDBLOCK)) 390 error = 0;391 392 if (error == EPIPE) psignal(p, SIGPIPE); 393

**INTEL Ex.1013.514** 

```
sosend Function 489
```

```
394  }
395  if (error == 0)
396      *retsize = len - auio.uio_resid;
397  bad:
398  if (to)
399      m_freem(to);
400    return (error);
401 }
```

— uipc\_syscalls.c

### Send data and cleanup

386-401 uio\_resid is saved in len so that the number of bytes transferred can be calculated if sosend does not accept all the data. The socket, destination address, uio structure, control information, and flags are all passed to sosend. When sosend returns, sendit responds as follows:

Figure 16.21 sendit function.

- If sosend transfers some data and is interrupted by a signal or a blocking condition, the error is discarded and the partial transfer is reported.
- If sosend returns EPIPE, the SIGPIPE signal is sent to the process. error is not set to 0, so if a process catches the signal and the signal handler returns, or if the process ignores the signal, the write call returns EPIPE.
- If no error occurred (or it was discarded), the number of bytes transferred is calculated and saved in \*retsize. Since sendit returns 0, syscall (Section 15.4) returns \*retsize to the process instead of returning the error code.
- If any other error occurs, the error code is returned to the process.

Before returning, sendit releases the mbuf containing the destination address. sosend is responsible for releasing the control mbuf.

### 16.7 sosend Function

sosend is one of the most complicated functions in the socket layer. Recall from Figure 16.8 that all five write calls eventually call sosend. It is sosend's responsibility to pass the data and control information to the pr\_usrreq function of the protocol associated with the socket according to the semantics supported by the protocol and the buffer limits specified by the socket. sosend never places data in the send buffer; it is the protocol's responsibility to store and remove the data.

The interpretation of the send buffer's sb\_hiwat and sb\_lowat values by sosend depends on whether the associated protocol implements reliable or unreliable data transfer semantics.

### Reliable Protocol Buffering

For reliable protocols, the send buffer holds both data that has not yet been transmitted and data that has been sent, but has not been acknowledged.  $sb_cc$  is the number of bytes of data that reside in the send buffer, and  $0 \le sb_cc \le sb_hiwat$ .

### sb\_cc may temporarily exceed sb\_hiwat when out-of-band data is sent.

It is sosend's responsibility to ensure that there is enough space in the send buffer before passing any data to the protocol layer through the pr\_usrreq function. The protocol layer adds the data to the send buffer. sosend transfers data to the protocol in one of two ways:

- If PR\_ATOMIC is set, sosend must preserve the message boundaries between the process and the protocol layer. In this case, sosend waits for enough space to become available to hold the entire message. When the space is available, an mbuf chain containing the entire message is constructed and passed to the protocol in a single call through the pr\_usrreq function. RDP and SPP are examples of this type of protocol.
- If PR\_ATOMIC is not set, sosend passes the message to the protocol one mbuf at a time and may pass a partial mbuf to avoid exceeding the high-water mark. This method is used with SOCK\_STREAM protocols such as TCP and SOCK\_SEQPACKET protocols such as TP4. With TP4, record boundaries are indicated explicitly with the MSG\_EOR flag (Figure 16.12), so it is not necessary for the message boundaries to be preserved by sosend.

TCP applications have no control over the size of outgoing TCP segments. For example, a message of 4096 bytes sent on a TCP socket will be split by the socket layer into two mbufs with external clusters, containing 2048 bytes each, assuming there is enough space in the send buffer for 4096 bytes. Later, during protocol processing, TCP will segment the data according to the maximum segment size for the connection, which is normally less than 2048.

When a message is too large to fit in the available buffer space and the protocol allows messages to be split, sosend still does not pass data to the protocol until the free space in the buffer rises above sb\_lowat. For TCP, sb\_lowat defaults to 2048 (Figure 16.4), so this rule prevents the socket layer from bothering TCP with small chunks of data when the send buffer is nearly full.

### Unreliable Protocol Buffering

With unreliable protocols (e.g., UDP), no data is ever stored in the send buffer and no acknowledgment is ever expected. Each message is passed immediately to the protocol where it is queued for transmission on the appropriate network device. In this case, sb\_cc is always 0, and sb\_hiwat specifies the maximum size of each write and indirectly the maximum size of a datagram.

Figure 16.4 shows that sb\_hiwat defaults to 9216 (9  $\times$  1024) for UDP. Unless the process changes sb\_hiwat with the SO\_SNDBUF socket option, an attempt to write a datagram larger than 9216 bytes returns with an error. Even then, other limitations of the protocol implementation may prevent a process from sending large datagrams. Section 11.10 of Volume 1 discusses these defaults and limits in other TCP/IP implementations.

9216 is large enough for a NFS write, which often defaults to 8192 bytes of data plus protocol headers.

### sosend Code

Figure 16.22 shows an overview of the sosend function. We discuss the four shaded sections separately.

271–278

5. 1

f

The arguments to sosend are: so, a pointer to the relevant socket; addr, a pointer to a destination address; uio, a pointer to a uio structure describing the I/O buffers in user space; top, an mbuf chain that holds data to be sent; control, an mbuf that holds control information to be sent; and flags, which contains options for this write call.

Normally, a process provides data to the socket layer through the uio mechanism and top is null. When the kernel itself is using the socket layer (such as with NFS), the data is passed to sosend as an mbuf chain pointed to by top, and uio is null.

<sup>279–304</sup> The initialization code is described separately.

### Lock send buffer

305-308 sosend's main processing loop starts at restart, where it obtains a lock on the send buffer with sblock before proceeding. The lock ensures orderly access to the socket buffer by multiple processes.

If MSG\_DONTWAIT is set in flags, then SBLOCKWAIT returns M\_NOWAIT, which tells sblock to return EWOULDBLOCK if the lock is not available immediately.

MSG\_DONTWAIT is used only by NFS in Net/3.

The main loop continues until sosend transfers all the data to the protocol (i.e., resid == 0).

### Check for space

<sup>309-341</sup> Before any data is passed to the protocol, various error conditions are checked and sosend implements the flow control and resource control algorithms described earlier. If sosend blocks waiting for more space to appear in the output buffer, it jumps back to restart before continuing.

### Use data from top

<sup>342–350</sup> Once space becomes available and sosend has obtained a lock on the send buffer, the data is prepared for delivery to the protocol layer. If uio is null (i.e., the data is in the mbuf chain pointed to by top), sosend checks MSG\_EOR and sets M\_EOR in the chain to mark the end of a logical record. The mbuf chain is ready for the protocol layer.

Chapter 16 S 492 Socket I/O – uipc\_socket.c 271 sosend(so, addr, uio, top, control, flags) 2 272 struct socket \*so; 273 struct mbuf \*addr; 274 struct uio \*uio; 275 struct mbuf \*top; 276 struct mbuf \*control; flags; 277 int 278 { MALANY MALENCER /\* initialization (Figure 16.23) \*/ 1.000 305 restart: if (error = sblock(&so->so\_snd, SBLOCKWAIT(flags))) 306 goto out; 307 /\* main loop, until resid == 0 \*/ 308 do { /\* wait for space in send buffer (Figure 16,24) \*/ do { 342 if (uio == NULL) { 343 /\* 344 \* Data is prepackaged in "top". 345 \*/ 346 resid = 0; 347 if (flags & MSG\_EOR) 348 top->m\_flags |= M\_EOR; 349 } else 350 do { 351 /\* fill a single mbuf or an mbuf chain (Figure 16.25) \*/ } while (space > 0 && atomic); 396 /\* pass mbuf chain to protocol (Figure 16.26) \*/ 2 } while (resid && space > 0); 2 412 } while (resid); 413 414 release: sbunlock(&so->so\_snd); 415 416 out: 417 if (top) m\_freem(top); 418 if (control) 419 m\_freem(control); 420 421 return (error); - uipc\_socket.c 422 } Figure 16.22 sosend function: overview.

States Barries

**INTEL Ex.1013.518** 

### Copy data from process

351-396 When uio is not null, sosend must transfer the data from the process. When PR\_ATOMIC is set (e.g., UDP), this loop continues until all the data has been stored in a single mbuf chain. A break, which is not shown in Figure 16.22, causes the loop to terminate when all the data has been copied from the process, and sosend passes the entire chain to the protocol.

When PR\_ATOMIC is not set (e.g., TCP), this loop is executed only once, filling a single mbuf with data from uio. In this case, the mbufs are passed one at a time to the protocol.

### Pass data to the protocol

397-413 For PR\_ATOMIC protocols, after the mbuf chain is passed to the protocol, resid is always 0 and control falls through the two loops to release. When PR\_ATOMIC is not set, sosend continues filling individuals mbufs while there is more data to send and while there is still space in the buffer. If the buffer fills and there is still data to send, sosend loops back and waits for more space before filling the next mbuf. If all the data is sent, both loops terminate.

#### Cleanup

414-422 After all the data has been passed to the protocol, the socket buffer is unlocked, any remaining mbufs are discarded, and sosend returns.

The detailed description of sosend is shown in four parts:

- initialization (Figure 16.23),
- error and resource checking (Figure 16.24),
- data transfer (Figure 16.25), and
- protocol dispatch (Figure 16.26).

The first part of sosend shown in Figure 16.23 initializes various variables.

#### Compute transfer size and semantics

atomic is set if sosendallatonce is true (any protocol for which PR\_ATOMIC is set) or the data has been passed to sosend as an mbuf chain in top. This flag controls whether data is passed to the protocol as a single mbuf chain or in separate mbufs.

resid is the number of bytes in the iovec buffers or the number of bytes in the top mbuf chain. Exercise 16.1 discusses why resid might be negative.

#### If requested, disable routing

304

cket.c

298-303 dontroute is set when the routing tables should be bypassed for *this* message only. clen is the number of bytes in the optional control mbuf.

The macro snderr posts the error code, reenables protocol processing, and jumps to the cleanup code at out. This macro simplifies the error handling within the function.

Figure 16.24 shows the part of sosend that checks for error conditions and waits for space to appear in the send buffer.

r 16

Socket I/O 494

uipc\_socket.c

uipc\_socket.c /\* XXX \*/ 279 struct proc \*p = curproc; struct mbuf \*\*mp; 280 281 struct mbuf \*m; long space, len, resid; 282 clen = 0, error, s, dontroute, mlen; 283 int atomic = sosendallatonce(so) || top; 284 int if (uio) 285 resid = uio->uio\_resid; 286 287 else resid = top->m\_pkthdr.len; 288 289 /\* \* In theory resid should be unsigned. 290 \* However, space must be signed, as it might be less than 0 291 \* if we over-committed, and we must use a signed comparison 292 \* of space and resid. On the other hand, a negative resid 293 \* causes us to loop sending 0-length segments to the protocol. 294 295 \*/ if (resid < 0) 296 return (EINVAL); 297 dontroute = 298 (flags & MSG\_DONTROUTE) && (so->so\_options & SO\_DONTROUTE) == 0 && 299 (so->so\_proto->pr\_flags & PR\_ATOMIC); 300 p->p\_stats->p\_ru.ru\_msgsnd++; 301 if (control) 302 303 clen = control->m\_len; { error = errno; splx(s); goto release; }

Figure 16.23 sosend function: initialization.

309

304 #define snderr(errno)

Protocol processing is suspended to prevent the buffer from changing while it is being examined. Before each transfer, sosend checks several conditions:

- If output from the socket is prohibited (e.g., the write-half of a TCP connection 310-311 has been closed), EPIPE is returned.
- If the socket is in an error state (e.g., an ICMP port unreachable may have been 312-313 generated by a previous datagram), so\_error is returned. sendit discards the error if some data has been sent before the error occurs (Figure 16.21, line 389).
- If the protocol requires connections and a connection has not been established or 314-318 a connection attempt has not been started, ENOTCONN is returned. sosend permits a write consisting of control information and no data even when a connection has not been established.

The Internet protocols do not use this feature, but it is used by TP4 to send data with a connection request, to confirm a connection request, and to send data with a disconnect request.

- 319-321
- If a destination address is not specified for a connectionless protocol (e.g., the process calls send without establishing a destination with connect), EDESTADDREQ is returned.

Se

新な世界にはないない

語言語語語を建立行

32

32

uipc socket.c

	uipc_socket.c
309	s = splnet();
310	if (so->so_state & SS_CANTSENDMORE)
311	<pre>snderr(EPIPE);</pre>
312	if (so->so_error)
313	<pre>snderr(so-&gt;so_error);</pre>
314	if ((so->so_state & SS_ISCONNECTED) == 0) {
315	if (so->so_proto->pr_flags & PR_CONNREQUIRED) {
316	if ((so->so_state & SS_ISCONFIRMING) == 0 &&
317	!(resid == 0 && clen != 0))
318	<pre>snderr(ENOTCONN);</pre>
319	} else if (addr == 0)
320	<pre>snderr(EDESTADDRREQ);</pre>
321	}
322	<pre>space = sbspace(&amp;so-&gt;so_snd);</pre>
323	if (flags & MSG_OOB)
324	space += 1024;
325	if (atomic && resid > so->so_snd.sb_hiwat
326	clen > so->so_snd.sb_hiwat)
327	<pre>snderr(EMSGSIZE);</pre>
328	if (space < resid + clen && uio &&
329	(atomic    space < so->so_snd.sb_lowat    space < clen)) {
330	if (so->so_state & SS_NBIO)
331	<pre>snderr(EWOULDBLOCK);</pre>
332	sbunlock(&so->so_snd);
333	error = sbwait(&so->so_snd);
334	<pre>splx(s);</pre>
335	if (error)
336	goto out;
337	goto restart;
338	}
339	<pre>splx(s);</pre>
340	mp = & top;
341	space -= clen; uipc_socket.c

Figure 16.24 sosend function: error and resource checking.

### Compute available space

sbspace computes the amount of free space remaining in the send buffer. This is 322-324 an administrative limit based on the buffer's high-water mark, but is also limited by sb\_mbmax to prevent many small messages from consuming too many mbufs (Figure 16.6). sosend gives out-of-band data some priority by relaxing the limits on the buffer size by 1024 bytes.

### Enforce message size limit

325-327

If atomic is set and the message is larger than the high-water mark, EMSGSIZE is returned; the message is too large to be accepted by the protocol-even if the buffer were empty. If the control information is larger than the high-water mark, EMSGSIZE is also returned. This is the test that limits the size of a datagram or record.

#### Wait for more space?

328-329 If there is not enough space in the send buffer, the data is from a process (versus from the kernel in top), and one of the following conditions is true, then sosend must wait for additional space before continuing:

- the message must be passed to protocol in a single request (atomic is set), or
- the message may be split, but the free space has dropped below the low-water mark, or
- the message may be split, but the control information does not fit in the available space.

When the data is passed to sosend in top (i.e., when uio is null), the data is already located in mbufs. Therefore sosend ignores the high- and low-water marks since no additional mbuf allocations are required to pass the data to the protocol.

If the send buffer low-water mark is not used in this test, an interesting interaction occurs between the socket layer and the transport layer that leads to performance degradation. [Crowcroft et al. 1992] provides details on this scenario.

### Wait for space

330-338 If sosend must wait for space and the socket is nonblocking, EWOULDBLOCK is returned. Otherwise, the buffer lock is released and sosend waits with sbwait until the status of the buffer changes. When sbwait returns, sosend reenables protocol processing and jumps back to restart to obtain a lock on the buffer and to check the error and space conditions again before continuing.

By default, sbwait blocks until data can be sent. By changing sb\_timeo in the buffer through the SO\_SNDTIMEO socket option, the process selects an upper bound for the wait time. If the timer expires, sbwait returns EWOULDBLOCK. Recall from Figure 16.21 that this error is discarded by sendit if some data has already been transferred to the protocol. This timer does not limit the length of the entire call, just the inactivity time between filling mbufs.

At this point, sosend has determined that some data may be passed to the protocol. splx enables interrupts since they should not be blocked during the relatively long time it takes to copy data from the process to the kernel. mp holds a pointer used to construct the mbuf chain. The size of the control information (clen) is subtracted from the space available before sosend transfers any data from the process.

Figure 16.25 shows the section of sosend that moves data from the process to one or more mbufs in the kernel.

#### Allocate packet header or standard mbuf

351-360

When atomic is set, this code allocates a packet header during the first iteration of the loop and standard mbufs afterwards. When atomic is not set, this code always allocates a packet header since top is always cleared before entering the loop.

sosend Function 497 Section 16.7 r 16؛ uipc\_socket.c 351 do { :sus if (top == 0) { 352 າust MGETHDR(m, M\_WAIT, MT\_DATA); 353 354 mlen = MHLEN; m->m\_pkthdr.len = 0; 355 m->m\_pkthdr.rcvif = (struct ifnet \*) 0; 356 } else { 357 ater 358 MGET(m, M\_WAIT, MT\_DATA); mlen = MLEN; 359 360 } /ailif (resid >= MINCLSIZE && space >= MCLBYTES) { 361 MCLGET(m, M\_WAIT); 362 363 if  $((m->m_flags \& M_EXT) == 0)$ a is 364 goto nopages; arks mlen = MCLBYTES; 365 if (atomic && top == 0) { 366 len = min(MCLBYTES - max\_hdr, resid); 367 tion m->m\_data += max\_hdr; 368 ince 369 } else . len = min(MCLBYTES, resid); 370 space -= MCLBYTES; 371 372 } else { K is 373 nopages: ıntil len = min(min(mlen, resid), space); 374 pro-375 space -= len; 376 /\* rror 377 \* For datagram protocols, leave room \* for protocol headers in first mbuf. 378 the \*/ 379 l for if (atomic && top == 0 && len < mlen) 380 Fig-381 MH\_ALIGN(m, len); } 382 ansthe error = uiomove(mtod(m, caddr\_t), (int) len, uio); 383 resid = uio->uio\_resid; 384 m->m\_len = len; 385 oto-386 \*mp = m;long top->m\_pkthdr.len += len; 387 con-388 if (error) ۱ the 389 goto release; mp = &m->m\_next; 390 if (resid <= 0) { 391 if (flags & MSG\_EOR) 392 one top->m\_flags != M\_EOR; 393 394 break; 395 } } while (space > 0 && atomic); 396 ——— uipc\_socket.c n of

### vays

Figure 16.25 sosend function: data transfer.

INTEL Ex.1013.523

#### If possible, use a cluster

361--371

If the message is large enough to make a cluster allocation worthwhile and space is greater than or equal to MCLBYTES, a cluster is attached to the mbuf by MCLGET. When space is less than MCLBYTES, the extra 2048 bytes will break the allocation limit for the buffer since the entire cluster is allocated even if resid is less than MCLBYTES.

If MCLGET fails, sosend jumps to nopages and uses a standard mbuf instead of an external cluster.

The test against MINCLSIZE should use >, not >=, since a write of 208 (MINCLSIZE) bytes fits within two mbufs.

When atomic is set (e.g., UDP), the mbuf chain represents a datagram or record and max\_hdr bytes are reserved at the front of the *first* cluster for protocol headers. Subsequent clusters are part of the same chain and do not need room for the headers.

If atomic is not set (e.g., TCP), no space is reserved since sosend does not know how the protocol will segment the outgoing data.

Notice that space is decremented by the size of the cluster (2048 bytes) and not by len, which is the number of data bytes to be placed in the cluster (Exercise 16.2).

#### Prepare the mbuf

372-382

If a cluster was not used, the number of bytes stored in the mbuf is limited by the smaller of: (1) the space in the mbuf, (2) the number of bytes in the message, or (3) the space in the buffer.

When atomic is set, MH\_ALIGN locates the data at the end of the buffer for the first buffer in the chain. MH\_ALIGN is skipped if the data completely fills the mbuf. This may or may not leave enough room for protocol headers, depending on how much data is placed in the mbuf. When atomic is not set, no space is set aside for the headers.

#### Get data from the process

383-395

uiomove copies len bytes of data from the process to the mbuf. After the transfer, the mbuf length is updated, the previous mbuf is linked to the new mbuf (or top points to the first mbuf), and the length of the mbuf chain is updated. If an error occurred during the transfer, sosend jumps to release.

When the last byte is transferred from the process, M\_EOR is set in the packet if the process set MSG\_EOR, and sosend breaks out of this loop.

MSG\_EOR applies only to protocols with explicit record boundaries such as TP4, from the OSI protocol suite. TCP does not support logical records and ignores the MSG\_EOR flag.

### Fill another buffer?

396

If atomic is set, sosend loops back and begins filling another mbuf.

The test for space > 0 appears to be extraneous. space is irrelevant when atomic is not set since the mbufs are passed to the protocol one at a time. When atomic is set, this loop is entered only when there is enough space for the entire message. See also Exercise 16.2.

The last section of sosend, shown in Figure 16.26, passes the data and control mbufs to the protocol associated with the socket.

400

39

Sec

SO

sosend Function Section 16.7 uipc\_socket.c if (dontroute) 397 so->so\_options |= SO\_DONTROUTE; 398 399 error = (\*so->so\_proto->pr\_usrreq) (so, 400 (flags & MSG\_OOB) ? PRU\_SENDOOB : PRU\_SEND, 401 top, addr, control); 402 403 splx(s); if (dontroute) 404 so->so\_options &= ~SO\_DONTROUTE; 405 clen = 0;406 control = 0;407 top = 0;408 409 mp = & top;if (error) 410 goto release; 411 } while (resid && space > 0); 412 413 } while (resid); - uipc\_socket.c

Figure 16.26 sosend function: protocol dispatch.

The socket's SO\_DONTROUTE option is toggled if necessary before and after passing 397-405 the data to the protocol layer to bypass the routing tables on this message. This is the only option that can be enabled for a single message and, as described with Figure 16.23, it is controlled by the MSG\_DONTROUTE flag during a write.

pr\_usrreq is bracketed with splnet and splx to block interrupts while the protocol is processing the message. This is a paranoid assumption since some protocols (such as UDP) may be able to do output processing without blocking interrupts, but this information is not available at the socket layer.

If the process tagged this message as out-of-band data, sosend issues the PRU\_SENDOOB request; otherwise it issues the PRU\_SEND request. Address and control mbufs are also passed to the protocol at this time.

clen, control, top, and mp are reset, since control information is passed to the 406-413 protocol only once and a new mbuf chain is constructed for the next part of the message. resid is nonzero only when atomic is not set (e.g., TCP). In that case, if space remains in the buffer, sosend loops back to fill another mbuf. If there is no more space, sosend loops back to wait for more space (Figure 16.24).

We'll see in Chapter 23 that unreliable protocols, such as UDP, immediately queue the data for transmission on the network. Chapter 26 describes how reliable protocols, such as TCP, add the data to the socket's send buffer where it remains until it is sent to, and acknowledged by, the destination.

### sosend Summary

>ter 16

pace

LGET.

ı limit

of an

ztes fits

'ecord

aders.

know

not by

by the

(3) the

he first

. This

h data

ansfer,

points

d dur-

t if the

the OSI

s not set ; loop is

control

rs.

ers.

ES.

sosend is a complex function. It is 142 lines long, contains three nested loops, one loop implemented with goto, two code paths based on whether PR\_ATOMIC is set or not, and two concurrency locks. As with much software, some of the complexity has accumulated over the years. NFS added the MSG\_DONTWAIT semantics and the possibility

499

of receiving data from an mbuf chain instead of the buffers in a process. The SS\_ISCONFIRMING state and MSG\_EOR flag were introduced to handle the connection and record semantics of the OSI protocols.

A cleaner approach would be to implement a separate sosend function for each type of protocol and dispatch through a pr\_send pointer in the protosw entry. This idea is suggested and implemented for UDP in [Partridge and Pink 1993].

### Performance Considerations

As described in Figure 16.25, sosend, when possible, passes message in mbuf-sized chunks to the protocol layer. While this results in more calls to the protocol than building and passing an entire mbuf chain, [Jacobson 1988a] reports that it improves performance by increasing parallelism.

Transferring one mbuf at a time (up to 2048 bytes) allows the CPU to prepare a packet while the network hardware is transmitting. Contrast this to sending a large mbuf chain: while the chain is being constructed, the network and the receiving system are idle. On the system described in [Jacobson 1988a], this change resulted in a 20% increase in network throughput.

It is important to make sure the send buffer is always larger than the bandwidthdelay product of a connection (Section 20.7 of Volume 1). For example, if TCP discovers that the connection can hold 20 segments before an acknowledgment is received, the send buffer must be large enough to hold the 20 unacknowledged segments. If it is too small, TCP will run out of data to send before the first acknowledgment is returned and the connection will be idle for some period of time.

# 16.8 read, readv, recvfrom, and recvmsg System Calls

These four system calls, which we refer to collectively as *read system calls*, receive data from a network connection. The first three system calls are simpler interfaces to the most general read system call, recvmsg. Figure 16.27 summarizes the features of the four read system calls and one library function (recv).

Function	Type of descriptor	Number of buffers	Return sender's address?	Flags?	Return control information?
read readv recv	any any sockets only	1 [1uio_maxiov] 1		•	
recvfrom recvmsg	sockets only sockets only	1 [1uio_maxiov]	•	•	•

Figure 16.27 Read system calls.

In Net/3, recv is implemented as a library function that calls recvfrom. For binary compatibility with previously compiled programs, the kernel maps the old recv system call to the function orecv. We discuss only the kernel implementation of recvfrom.

The read and readv system calls are valid with any descriptor, but the remaining calls are valid only with socket descriptors.

As with the write calls, multiple buffers are specified by an array of iovec structures. For datagram protocols, recvfrom and recvmsg return the source address associated with each incoming datagram. For connection-oriented protocols, getpeername returns the address associated with the other end of the connection. The flags associated with the receive calls are shown in Section 16.11.

As with the write calls, the receive calls utilize a common function, in this case soreceive, to do all the work. Figure 16.28 illustrates the flow of control for the read system calls.

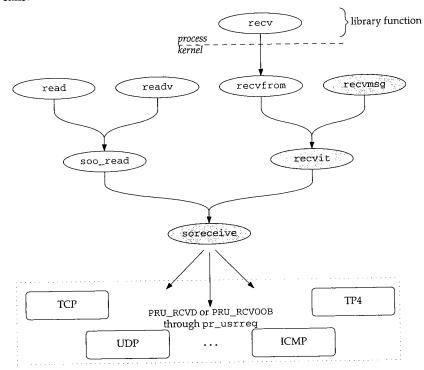


Figure 16.28 All socket input is processed by soreceive.

We discuss only the three shaded functions in Figure 16.28. The remaining functions are left for readers to investigate on their own.

### 16.9 recymsg System Call

The recvmsg function is the most general read system call. Addresses, control information, and receive flags may be discarded without notification if a process uses one of the other read system calls while this information is pending. Figure 16.29 shows the recvmsg function.

Chapter 16

Se

46

16

433 struct recvmsg\_args { uipc\_syscalls.c 434 int s; 435 struct msghdr \*msg; 436 int flags; 437 }; 438 recvmsg(p, uap, retval) 439 struct proc \*p; 440 struct recvmsg\_args \*uap; 441 int \*retval; 442 { 443 struct msghdr msg; 444 struct iovec aiov[UIO\_SMALLIOV], \*uiov, \*iov; 445 int error; 446 if (error = copyin((caddr\_t) uap->msg, (caddr\_t) & msg, sizeof(msg))) 447 return (error); 448 if ((u\_int) msg.msg\_iovlen >= UIO\_SMALLIOV) { 449 if ((u\_int) msg.msg\_iovlen >= UIO\_MAXIOV) 450 return (EMSGSIZE); 451 MALLOC(iov, struct iovec \*, 452 sizeof(struct iovec) \* (u\_int) msg.msg\_iovlen, M\_IOV, 453 M\_WAITOK); 454 } else 455 iov = aiov;456 msg.msg\_flags = uap->flags; 457 uiov = msg.msg\_iov; 458 msg.msg\_iov = iov; 459 if (error = copyin((caddr\_t) uiov, (caddr\_t) iov, 460 (unsigned) (msg.msg\_iovlen \* sizeof(struct iovec)))) 461 aoto done; 462 if ((error = recvit(p, uap->s, &msg, (caddr\_t) 0, retval)) == 0) { 463 msg.msg\_iov = uiov; 464 error = copyout((caddr\_t) & msg, (caddr\_t) uap->msg, sizeof(msg)); 465 } 466 done: 467 if (iov != aiov) 468 FREE(iov, M\_IOV); 469 return (error); 470 } uipc\_syscalls.c

Figure 16.29 recvmsg system call.

The three arguments to recymsg are: the socket descriptor; a pointer to a msghdr 433-445 structure; and several control flags.

### Copy iov array

```
446-461
```

As with sendmsg, recvmsg copies the msghdr structure into the kernel, allocates a larger iovec array if the automatic array aiov is too small, and copies the array entries from the process into the kernel array pointed to by iov (Section 16.4). The flags provided as the third argument are copied into the msghdr structure.

471--5

### recvit and cleanup

462-470 After recvit has received data, the msghdr structure is copied back into the process with the updated buffer lengths and flags. If a larger iovec structure was allocated, it is released before recvmsg returns.

### 16.10 recvit Function

The recvit function shown in Figures 16.30 and 16.31 is called from recv, recvfrom, and recvmsg. It prepares a uio structure for processing by soreceive based on the msghdr structure prepared by the recvxxx calls.

471 re	ecvit(p, s, mp, namelenp, retsize)	uipc_syscalls.c
	ruct proc *p;	
473 in		
474 st	ruct msghdr *mp;	
	ddr_t namelenp;	
476 in	t *retsize;	
477 {		
478	struct file *fp;	
479	struct uio auio;	
480	struct iovec *iov;	
481	int i;	
482	int len, error;	
483	<pre>struct mbuf *from = 0, *control = 0;</pre>	
484	if (error = getsock(p->p_fd, s, &fp))	
485	return (error);	
486	<pre>auio.uio_iov = mp-&gt;msg_iov;</pre>	
487	auio.uio_iovcnt = mp->msg_iovlen;	
488	auio.uio_segflg = UIO_USERSPACE;	
489	auio.uio_rw = UIO_READ;	
490	auio.uio_procp = p;	
491	auio.uio_offset = 0; /* XXX */	
492	auio.uio_resid = 0;	
493	iov = mp->msg_iov;	
494	for (i = 0; i < mp->msg_iovlen; i++, iov++) {	
495	if (iov->iov_len < 0)	
496	return (EINVAL);	
497	if ((auio.uio_resid += iov->iov_len) < 0)	
498	return (EINVAL);	
499	}	
500	len = auio.uio_resid;	

Figure 16.30 recvit function: initialize uio structure.

471-500 getsock returns the file structure for the descriptor s, and then recvit initializes the uio structure to describe a read transfer from the kernel to the process. The number of bytes to transfer is computed by summing the msg\_iovlen members of the iovec array. The total is saved in uio\_resid and in len.

The second half of recvit, shown in Figure 16.31, calls soreceive and copies the results back to the process.

## ılls.c

))

'ls.c

dr

s a

ies

ro-

er 16

4 Socket	I/O Chapter 16		Section
			<u> </u>
501		2.4	501-51
501 502 503	<pre>(struct mbuf **) 0, mp-&gt;msg_control ? &amp;control : (struct mbuf **) 0,</pre>		501 51
503	if (auio.uio_resid != len && (error == ERESTART	a di sana sa	
505	<pre>error == EINTR (  error == EWOULDBLOCK))</pre>		
506	error = 0;		
507	}		
508	if (error)		511-542
509	goto out; *retsize = len - auio.uio_resid;		511-542
510 511	if (mp->msg_name) {		
511	<pre>len = mp-&gt;msg_namelen;</pre>	X.	
513	if (len <= 0    from == 0)	-0 	
514	len = 0;		
515	else {	2 2	
516	if (len > from->m_len)		
517	<pre>len = from-&gt;m_len; /* else if len &lt; from-&gt;m_len ??? */</pre>	5. 1997 -	
518 519	if (error = copyout(mtod(from, caddr_t),	501 101	543 544
520	(caddr_t) mp->msg_name, (unsigned) len))	21. 21.	543-549
521	goto out;	- 14 A.	
522	}		
523	<pre>mp-&gt;msg_namelen = len;</pre>		16.11
524	if (namelenp &&	e.	10.11
525	<pre>(error = copyout((caddr_t) &amp; len, namelenp, sizeof(int)))) {</pre>		
526	goto out; }		
527 528	}		
529	if (mp->msg_control) {		
530	<pre>len = mp-&gt;msg_controllen;</pre>	사람 사람	
531	if (len $\leq 0$    control $= 0$ )	135	
532	len = 0;		
533	else {		
534	if (len >= control->m_len)		
535	<pre>len = control-&gt;m_len;</pre>	- 420) - 414	
536 537	else mp->msg_flags  = MSG_CTRUNC;	- 18 C	
538	<pre>error = copyout((caddr_t) mtod(control, caddr_t),</pre>	1999 1997 1997	
539	<pre>(caddr_t) mp-&gt;msg_control, (unsigned) len);</pre>		
540	}		
541	<pre>mp-&gt;msg_controllen = len;</pre>		
542	}		
543	out:	Zity New York	
544	if (from)		
545 546	<pre>m_freem(from); if (control)</pre>		
540	<pre>m_freem(control);</pre>	in Alexandra Andreas Alexandra	
548	return (error);		Out-of-
549			- at-01-1

INTEL Ex.1013.530

1

#### **Call** soreceive

<sup>501-510</sup> soreceive implements the complex semantics of receiving data from the socket buffers. The number of bytes transferred is saved in \*retsize and returned to the process. When an signal arrives or a blocking condition occurs after some data has been copied to the process (len is not equal to uio\_resid), the error is discarded and the partial transfer is reported.

### Copy address and control information to the process

511-542

If the process provided a buffer for an address or control information or both, the buffers are filled and their lengths adjusted according to what soreceive returned. An address may be truncated if the buffer is too small. This can be detected by the process if it saves the buffer length before the read call and compares it with the value returned by the kernel in the namelenp variable (or in the length field of the sockaddr structure). Truncation of control information is reported by setting MSG\_CTRUNC in msg\_flags. See also Exercise 16.7.

### Cleanup

543-549 At out, the mbufs allocated for the source address and the control information are released.

### 16.11 soreceive Function

This function transfers data from the receive buffer of the socket to the buffers specified by the process. Some protocols provide an address specifying the sender of the data, and this can be returned along with additional control information that may be present. Before examining the code, we need to discuss the semantics of a receive operation, outof-band data, and the organization of a socket's receive buffer.

Figure 16.32 lists the flags that are recognized by the kernel during soreceive.

flags	Description	Reference
MSG_DONTWAIT	do not wait for resources during this call	Figure 16.38
MSG_OOB	receive out-of-band data instead of regular data	Figure 16.39
MSG_PEEK	receive a copy of the data without consuming it	Figure 16.43
MSG_WAITALL	wait for data to fill buffers before returning	Figure 16.50

Figure 16.32 recvxxx system calls: flag values passed to kernel.

recvmsg is the only read system call that returns flags to the process. In the other calls, the information is discarded by the kernel before control returns to the process. Figure 16.33 lists the flags that recvmsg can set in the msghdr structure.

### **Out-of-Band Data**

Out-of-band (OOB) data semantics vary widely among protocols. In general, protocols expedite OOB data along a previously established communication link. The OOB data might not remain in sequence with previously sent regular data. The socket layer

いたいち かいかいかい たいたいに ちょうがいたい ないない

msg_flags	Description	Reference
MSG_CTRUNC	the control information received was larger than the buffer provided	Figure 16.31
MSG_EOR	the data received marks the end of a logical record	Figure 16.48
MSG_OOB	the buffer(s) contains out-of-band data	Figure 16.45
MSG_TRUNC	the message received was larger than the buffer(s) provided	Figure 16.51

Figure 16.33 recvmsg system call: msg\_flag values returned by kernel.

supports two mechanisms to facilitate handling OOB data in a protocol-independent way: tagging and synchronization. In this chapter we describe the abstract OOB mechanisms implemented by the socket layer. UDP does not support OOB data. The relationship between TCP's urgent data mechanism and the socket OOB mechanism is described in the TCP chapters.

A sending process tags data as OOB data by setting the MSG\_OOB flag in any of the sendxxx calls. sosend passes this information to the socket's protocol, which provides any special services, such as expediting the data or using an alternate queueing strategy.

When a protocol receives OOB data, the data is set aside instead of placing it in the socket's receive buffer. A process receives the pending OOB data by setting the MSG\_OOB flag in one of the recvxxx calls. Alternatively, the receiving process can ask the protocol to place OOB data inline with the regular data by setting the SO\_OOBINLINE socket option (Section 17.3). When SO\_OOBINLINE is set, the protocol places incoming OOB data in the receive buffer with the regular data. In this case, MSG\_OOB is not used to receive the OOB data. Read calls return either all regular data or all OOB data. The two types are never mixed in the input buffers of a single input system call. A process that uses recvmsg to receive data can examine the MSG\_OOB flag to determine if the returned data is regular data or OOB data that has been placed inline.

The socket layer supports synchronization of OOB and regular data by allowing the protocol layer to mark the point in the regular data stream at which OOB data was received. The receiver can determine when it has reached this mark by using the SIOCATMARK ioctl command after each read system call. When receiving regular data, the socket layer ensures that only the bytes preceding the mark are returned in a single message so that the receiver does not inadvertently pass the mark. If additional OOB data is received before the receiver reaches the mark, the mark is silently advanced.

#### Example

Figure 16.34 illustrates the two methods of receiving out-of-band data. In both examples, bytes A through I have been received as regular data, byte J as out-of-band data, and bytes K and L as regular data. The receiving process has accepted all data up to but not including byte A.

In the first example, the process can read bytes A through I or, if MSG\_OOB is set, byte J. Even if the length of the read request is more than 9 bytes (A–I), the socket layer

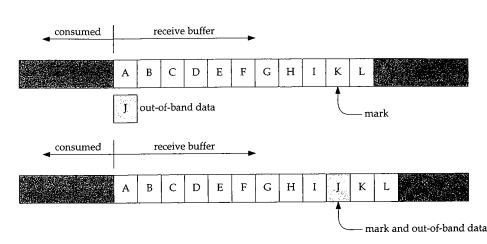


Figure 16.34 Receiving out-of-band data.

returns only 9 bytes to avoid passing the out-of-band synchronization mark. When byte I is consumed, SIOCATMARK is true; it is not necessary to consume byte J for the process to reach the out-of-band mark.

In the second example, the process can read only bytes A through I, at which point SIOCATMARK is true. A second call can read bytes J through L.

In Figure 16.34, byte J is *not* the byte identified by TCP's urgent pointer. The urgent pointer in this example would point to byte K. See Section 29.7 for details.

### **Other Receive Options**

Section 16.11

6

1 8 5

1

1t 1-

۱is

ιe

≥s y. ∖e

۱e

зk

າe ol

ю,

ta

ut 3g

эd

he

as

he

ar

ιa

۱al

tly

m-.ta, vut

set, γer A process can set the MSG\_PEEK flag to retrieve data without consuming it. The data remains on the receive queue until a read system call without MSG\_PEEK is processed.

The MSG\_WAITALL flag indicates that the call should not return until enough data can be returned to fulfill the entire request. Even if soreceive has some data that can be returned to the process, it waits until additional data has been received.

When MSG\_WAITALL is set, soreceive can return without filling the buffer in the following cases:

- the read-half of the connection is closed,
- the socket's receive buffer is smaller than the size of the read,
- an error occurs while the process is waiting for additional data,
- out-of-band data becomes available, or
- the end of a logical record occurs before the read buffer is filled.

NFS is the only software in Net/3 that uses the MSG\_WAITALL and MSG\_DONTWAIT flags. MSG\_DONTWAIT can be set by a process to issue a nonblocking read system call without selecting nonblocking I/O with ioctl or fcntl.

のないので、「ない」ので、

いた 一般語を感染がないため

「「「「

### **Receive Buffer Organization: Message Boundaries**

For protocols that support message boundaries, each message is stored in a single chain of mbufs. Multiple messages in the receive buffer are linked together by m\_nextpkt to form a queue of mbufs (Figure 2.21). The protocol processing layer adds data to the receive queue and the socket layer removes data from the receive queue. The high-water mark for a receive buffer restricts the amount of data that can be stored in the buffer.

When PR\_ATOMIC is not set, the protocol layer stores as much data in the buffer as possible and discards the portion of the incoming data that does not fit. For TCP, this means that any data that arrives and is outside the receive window is discarded. When PR\_ATOMIC is set, the entire message must fit within the buffer. If the message does not fit, the protocol layer discards the entire message. For UDP, this means that incoming datagrams are discarded when the receive buffer is full, probably because the process is not reading datagrams fast enough.

Protocols with PR\_ADDR set use sbappendaddr to construct an mbuf chain and add it to the receive queue. The chain contains an mbuf with the source address of the message, 0 or more control mbufs, followed by 0 or more mbufs containing the data.

For SOCK\_SEQPACKET and SOCK\_RDM protocols, the protocol builds an mbuf chain for each record and calls sbappendrecord to append the record to the end of the receive buffer if PR\_ATOMIC is set. If PR\_ATOMIC is not set (OSI's TP4), a new record is started with sbappendrecord. Additional data is added to the record with sbappend.

It is not correct to assume that PR\_ATOMIC indicates the buffer organization. For example, TP4 does not have PR\_ATOMIC set, but supports record boundaries with the M\_EOR flag.

Figure 16.35 illustrates the organization of a UDP receive buffer consisting of 3 mbuf chains (i.e., three datagrams). The m\_type value for each mbuf is included.

In the figure, the third datagram has some control information associated with it. Three UDP socket options can cause control information to be placed in the receive buffer. See Figure 22.5 and Section 23.7 for details.

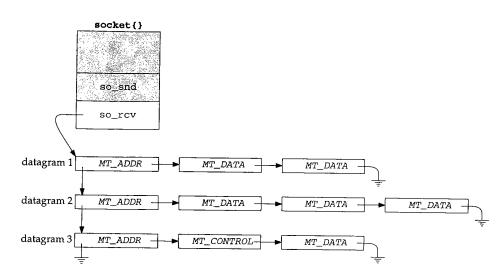
For PR\_ATOMIC protocols, sb\_lowat is ignored while data is being received. When PR\_ATOMIC is not set, sb\_lowat is the smallest number of bytes returned in a read system call. There are some exceptions to this rule, discussed with Figure 16.41.

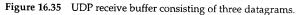
### **Receive Buffer Organization: No Message Boundaries**

When the protocol does not maintain message boundaries (i.e., SOCK\_STREAM protocols such as TCP), incoming data is appended to the end of the last mbuf chain in the buffer with sbappend. Incoming data is trimmed to fit within the receive buffer, and sb\_lowat puts a lower bound on the number of bytes returned by a read system call.

Figure 16.36 illustrates the organization of a TCP receive buffer, which contains only regular data.

「日本で





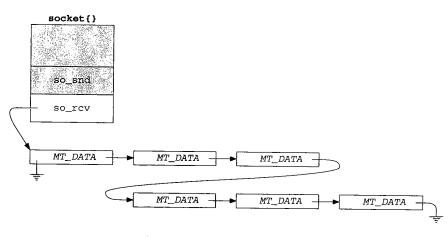


Figure 16.36 so\_rcv buffer for TCP.

### **Control Information and Out-of-band Data**

Unlike TCP, some stream protocols support control information and call sbappendcontrol to append the control information and the associated data as a new mbuf chain in the receive buffer. If the protocol supports inline OOB data, sbinsertoob inserts a new mbuf chain just after any mbuf chain that contains OOB data, but before any mbuf chain with regular data. This ensures that incoming OOB data is queued ahead of any regular data.

510 Socket I/O

Chapter 16

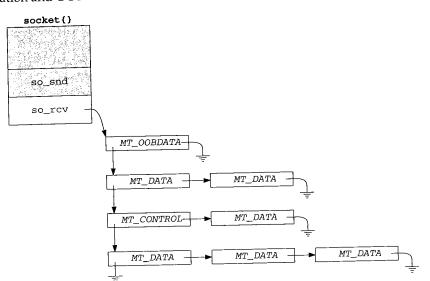


Figure 16.37 illustrates the organization of a receive buffer that contains control information and OOB data.

Figure 16.37 so\_rcv buffer with control and OOB data.

The Unix domain stream protocol supports control information and the OSI TP4 protocol supports MT\_OOBDATA mbufs. TCP does not support control data nor does it support the MT\_OOBDATA form of out-of-band data. If the byte identified by TCP's urgent pointer is stored inline (SO\_OOBINLINE is set), it appears as regular data, not OOB data. TCP's handling of the urgent pointer and the associated byte is described in Section 29.7.

### 16.12 soreceive Code

We now have enough background information to discuss soreceive in detail. While receiving data, soreceive must respect message boundaries, handle addresses and control information, and handle any special semantics identified by the read flags (Figure 16.32). The general rule is that soreceive processes one record per call and tries to return the number of bytes requested. Figure 16.38 shows an overview of the function.

439-446

soreceive has six arguments. so is a pointer to the socket. A pointer to an mbuf to receive address information is returned in \*paddr. If mp0 points to an mbuf pointer, soreceive transfers the receive buffer data to an mbuf chain pointed to by \*mp0. In this case, the uio structure is used only for the count in uio\_resid. If mp0 is null, soreceive copies the data into buffers described by the uio structure. A pointer to the mbuf containing control information is returned in \*controlp, and soreceive returns the flags described in Figure 16.33 in \*flagsp.

soreceive starts by setting pr to point to the socket's protocol switch structure 447-453 and saving uio\_resid (the size of the receive request) in orig\_resid. If control information or addressing information is copied from the kernel to the process, orig\_resid is set to 0. If data is copied, uio\_resid is updated. In either case, orig\_resid will not equal uio\_resid. This fact is used at the end of soreceive (Figure 16.51). \*paddr and \*controlp are cleared. The flags passed to soreceive in \*flagsp 454-461 are saved in flags after the MSG\_EOR flag is cleared (Exercise 16.8). flagsp is a value-result argument, but only the recymsg system call can receive the result flags. If flagsp is null, flags is set to 0. Before accessing the receive buffer, sblock locks the buffer. soreceive waits for 483-487 the lock unless MSG\_DONTWAIT is set in flags. This is another side effect of supporting calls to the socket layer from NFS within the kernel. Protocol processing is suspended, so soreceive is not interrupted while it examines the buffer. m is the first mbuf on the first chain in the receive buffer. If necessary, wait for data soreceive checks several conditions and if necessary waits for more data to arrive 488-541 in the buffer before continuing. If soreceive sleeps in this code, it jumps back to restart when it wakes up to see if enough data has arrived. This continues until the request can be satisfied. soreceive jumps to dontblock when it has enough data to satisfy the request. A 542-545 pointer to the second chain in the receive buffer is saved in nextrecord. Process address and control information Address information and control information are processed before any other data is 546-590 transferred from the receive buffer. Setup data transfer Since only OOB data or regular data is transferred in a single call to soreceive, 591-597 this code remembers the type of data at the front of the queue so soreceive can stop the transfer when the type changes. Mbuf data transfer loop This loop continues as long as there are mbufs in the buffer (m is not null), the 598-692 requested number of bytes has not been transferred (uio\_resid > 0), and no error has occurred. Cleanup 693-719 The remaining code updates various pointers, flags, and offsets; releases the socket buffer lock; enables protocol processing; and returns.

rol

16

3 it P's 10t in

'P4

Juf

:er,

In 111, to ve 512 Socket I/O

Chapter 16

10.00

uipc\_socket.c 439 soreceive(so, paddr, uio, mp0, controlp, flagsp) 440 struct socket \*so; 441 struct mbuf \*\*paddr; 442 struct uio \*uio; 443 struct mbuf \*\*mp0; 444 struct mbuf \*\*controlp; 445 int \*flagsp; 446 { struct mbuf \*m, \*\*mp; 447 int flags, len, error, s, offset; 448 struct protosw \*pr = so->so\_proto; 449 struct mbuf \*nextrecord; 450 int moff, type; 451 orig\_resid = uio->uio\_resid; int 452 mp = mp0;453 if (paddr) 454 \*paddr = 0; 455 if (controlp) 456 \*controlp = 0; 457 if (flagsp) 458 flags = \*flagsp & ~MSG\_EOR; 459 else 460 flags = 0; 461 /\* MSG\_00B processing and \*/ /\* implicit connection confirmation \*/ restart: 483 if (error = sblock(&so->so\_rcv, SBLOCKWAIT(flags))) 484 return (error); 485 s = splnet(); 486 m = so->so\_rcv.sb\_mb; 487 人名法尔 医帕克德氏试验 /\* if necessary, wait for data to arrive \*//\* if necessary, ware ---dontblock: 542 if (uio->uio\_procp) 543 uio->uio\_procp->p\_stats->p\_ru.ru\_msgrcv++; 544 nextrecord = m->m\_nextpkt; 545 /\* process address and control information \*/ if (m) { 591 if ((flags & MSG\_PEEK) == 0) 592 m->m\_nextpkt = nextrecord; 593 594 type = m->m\_type; if (type == MT\_OOBDATA) 595 596 flags |= MSG\_OOB; } 597

### INTEL Ex.1013.538

/\* process data \*/

693

}

/\* while more data and more space to fill \*/

/\* cleanup \*/

715 release: 716 sbunlock(&so->so\_rcv); 717 splx(s); 718 return (error); 719 }

– uipc\_socket.c

Figure 16.38 soreceive function: overview.

In Figure 16.39, soreceive handles requests for OOB data.

	uipc socket.c
462	if (flags & MSG_OOB) (
463	$m = m_get(M_WAIT, MT_DATA);$
464	error = (*pr->pr_usrreq) (so, PRU_RCVOOB,
465	<pre>m, (struct mbuf *) (flags &amp; MSG_PEEK), (struct mbuf *) 0);</pre>
466	if (error)
467	goto bad;
468	do {
469	error = uiomove(mtod(m, caddr_t),
470	<pre>(int) min(uio-&gt;uio_resid, m-&gt;m_len), uio);</pre>
471	$m = m_free(m);$
472	) while (uio->uio_resid && error == 0 && m);
473	bad:
474	if (m)
475	m_freem(m);
476	return (error);
477	}

Figure 16.39 soreceive function: out-of-band data.

#### **Receive OOB data**

「時ににたた」に

- State H

462-477 Since OOB data is not stored in the receive buffer, soreceive allocates a standard mbuf and issues the PRU\_RCVOOB request to the protocol. The while loop copies any data returned by the protocol to the buffers specified by uio. After the copy, soreceive returns 0 or the error code.

UDP always returns EOPNOTSUPP for the PRU\_RCVOOB request. See Section 30.2 for details regarding TCP urgent processing. In Figure 16.40, soreceive handles connection confirmation.

Chapter 16

の言語を意識を認識で

			– uipc_socket.c
478	if	(qm)	, _
479		<pre>*mp = (struct mbuf *) 0;</pre>	
480	if	(so->so_state & SS_ISCONFIRMING && uio->uio_resid)	
481		(*pr->pr_usrreq) (so, PRU_RCVD, (struct mbuf *) 0,	
482		<pre>(struct mbuf *) 0, (struct mbuf *) 0);</pre>	
			– uipc_socket.c

Figure 16.40 soreceive function: connection confirmation.

### **Connection confirmation**

478-482 If the data is to be returned in an mbuf chain, \*mp is initialized to null. If the socket is in the SO\_ISCONFIRMING state, the PRU\_RCVD request notifies the protocol that the process is attempting to receive data.

> The SO\_ISCONFIRMING state is used only by the OSI stream protocol, TP4. In TP4, a connection is not considered complete until a user-level process has confirmed the connection by attempting to send or receive data. The process can reject a connection by calling shutdown or close, perhaps after calling getpeername to determine where the connection came from.

Figure 16.38 showed that the receive buffer is locked before it is examined by the code in Figure 16.41. This part of soreceive determines if the read system call can be satisfied by the data that is already in the receive buffer.

	uipc_socket.c
488	/* , _
489	* If we have less data than requested, block awaiting more
490	* (subject to any timeout) if:
491	* 1. the current count is less than the low water mark, or
492	$\star$ 2. MSG_WAITALL is set, and it is possible to do the entire
493	* receive operation at once if we block (resid <= hiwat).
494	* 3. MSG_DONTWAIT is not set
495	*
496	$\star$ If MSG_WAITALL is set but resid is larger than the receive buffer,
497	* we have to do the receive in sections, and thus risk returning
498	$\star$ a short count if a timeout or signal occurs after we start.
499	*/
500	if (m == 0  } ((flags & MSG_DONTWAIT) == 0 &&
501	so->so_rcv.sb_cc < uio->uio_resid) &&
502	(so->so_rcv.sb_cc < so->so_rcv.sb_lowat
503	((flags & MSG_WAITALL) && uio->uio_resid <= so->so_rcv.sb_hiwat)) &&
504	<pre>m-&gt;m_nextpkt == 0 &amp;&amp; (pr-&gt;pr_flags &amp; PR_ATOMIC) == 0) {     uipc socket.c</pre>

Figure 16.41 soreceive function: enough data?

#### Can the call be satisfied now?

488-504

<sup>4</sup> The general rule for soreceive is that it waits until enough data is in the receive buffer to satisfy the entire read. There are several conditions that cause an error or less data than was requested to be returned.

If any of the following conditions are true, the process is put to sleep to wait for more data to arrive so the call can be satisfied:

soreceive Code 515

Section 16.12

socket c	

apter 16

- : socket.c
- e socket that the

a connecection by hutdown

me from.

l by the ll can be

- c\_socket.c
- ffer,
- g

t)) &&

- >c\_socket.c
- e receive or or less
- wait for

- There is no data in the receive buffer (m equals 0).
- There is not enough data to satisfy the entire read (sb\_cc < uio\_resid and MSG\_DONTWAIT is not set), the minimum amount of data is *not* available (sb\_cc < sb\_lowat), and more data can be appended to this chain when it arrives (m\_nextpkt is 0 and PR\_ATOMIC is *not* set).
- There is not enough data to satisfy the entire read, a minimum amount of data *is* available, data can be added to this chain, but MSG\_WAITALL indicates that soreceive should wait until the entire read can be satisfied.

If the conditions in the last case are met but the read is too large to be satisfied without blocking (uio\_resid > sb\_hiwat), soreceive continues without waiting for more data.

If there is some data in the buffer and MSG\_DONTWAIT is set, soreceive does not wait for more data.

There are several reasons why waiting for more data may not be appropriate. In Figure 16.42, soreceive checks for these conditions and returns, or waits for more data to arrive.

### Wait for more data?

- 505-534 At this point, soreceive has determined that it must wait for additional data to arrive before the read can be satisfied. Before waiting it checks for several additional conditions:
- If the socket is in an error state and *empty* (m is null), soreceive returns the error code. If there is an error and the receive buffer also contains data (m is nonnull), the data is returned and a subsequent read returns the error when there is no more data. If MSG\_PEEK is set, the error is not cleared, since a read system call with MSG\_PEEK set should not change the state of the socket.
- If the read-half of the connection has been closed and data remains in the receive buffer, sosend does not wait and returns the data to the process (at dontblock). If the receive buffer is empty, soreceive jumps to release and the read system call returns 0, which indicates that the read-half of the connection is closed.
  - If the receive buffer contains out-of-band data or the end of a logical record, soreceive does not wait for additional data and jumps to dontblock.
  - If the protocol requires a connection and it does not exist, ENOTCONN is posted and the function jumps to release.
  - If the read is for 0 bytes or nonblocking semantics have been selected, the function jumps to release and returns 0 or EWOULDBLOCK, respectively.

### Yes, wait for more data

519-523

524-528

529-534

535-541 soreceive has now determined that it must wait for more data, and that it is reasonable to do so (i.e., some data will arrive). The receive buffer is unlocked while the process sleeps in sbwait. If sbwait returns because of an error or a signal,

Chapter 16

	——————————————————————————————————————	pc_socket.c
505	if (so->so_error) {	
506	if (m)	
507	goto dontblock;	
508	error = so->so_error;	
509	if ((flags & MSG_PEEK) == 0)	
510	so->so_error = 0;	
511	goto release;	
512	}	
513	if (so->so_state & SS_CANTRCVMORE) {	
514	if (m)	
515	goto dontblock;	
516	else	
517	goto release;	
518	)	
519	for (; m; m = $m - m_next$ )	
520	if (m->m_type == MT_OOBDATA    (m->m_flags & M_EOR)) {	
521	<pre>m = so-&gt;so_rcv.sb_mb;</pre>	
522	goto dontblock;	
523	}	
524	if ((so->so_state & (SS_ISCONNECTED   SS_ISCONNECTING)) == (	33 C
525	(so->so_proto->pr_flags & PR_CONNREQUIRED)) {	
526	error = ENOTCONN;	
527	goto release;	
528	}	
529	if (uio->uio_resid == 0)	
530	goto release;	
531	if ((so->so_state & SS_NBIO)    (flags & MSG_DONTWAIT)) {	
532	error = EWOULDBLOCK;	
533	goto release;	
534	}	
535	<pre>sbunlock(&amp;so-&gt;so_rcv);</pre>	
536	error = sbwait(&so->so_rcv);	
537	<pre>splx(s);</pre>	
538	if (error)	
539	return (error);	
540	goto restart;	
541	}	ipc_socket.

Figure 16.42 soreceive function: wait for more data?

soreceive returns the error; otherwise the function jumps to restart to determine if the read can be satisfied now that more data has arrived.

As in sosend, a process can enable a receive timer for sbwait with the SO\_RCVTIMEO socket option. If the timer expires before any data arrives, sbwait returns EWOULDBLOCK.

The effect of this timer is not what one would expect. Since the timer gets reset every time there is activity on the socket buffer, the timer never expires if at least 1 byte arrives within the timeout interval. This can delay the return of the read system call for more than the value of the timer. sb\_timeo is an inactivity timer and does not put an upper bound on the amount of time that may be required to satisfy the read system call.

### **INTEL Ex.1013.542**

5(

С

.c

?t.C

e if

:he

it

me

the

≥of tof At this point, soreceive is prepared to transfer some data from the receive buffer. Figure 16.43 shows the transfer of any address information.

542	dontblock:	———— uipc_socket.c
543	if (uio->uio_procp)	
544	uio->uio_procp->p_stats->p_ru.ru_msgrcv++;	
545	<pre>nextrecord = m-&gt;m_nextpkt;</pre>	
546	if (pr->pr_flags & PR_ADDR) (	
547	orig_resid = 0;	
548	if (flags & MSG_PEEK) {	
549	if (paddr)	
550	<pre>*paddr = m_copy(m, 0, m-&gt;m_len);</pre>	
551	<pre>m = m-&gt;m_next;</pre>	
552	} else {	
553	<pre>sbfree(&amp;so-&gt;so_rcv, m);</pre>	
554	if (paddr) {	
555	<pre>*paddr = m;</pre>	
556	so->so_rcv.sb_mb = m->m_next;	
557	$m \rightarrow m_next = 0;$	
558	<pre>m = so-&gt;so_rcv.sb_mb;</pre>	
559	} else {	
560	<pre>MFREE(m, so-&gt;so_rcv.sb_mb);</pre>	
561	<pre>m = so-&gt;so_rcv.sb_mb;</pre>	
562	}	
563	}	
564	}	

Figure 16.43 soreceive function: return address information.

### dontblock

<sup>542-545</sup> nextrecord maintains a reference to the next record that appears in the receive buffer. This is used at the end of soreceive to attach the remaining mbufs to the socket buffer after the first chain has been discarded.

### **Return address information**

<sup>546–564</sup> If the protocol provides addresses, such as UDP, the mbuf containing the address is removed from the mbuf chain and returned in \*paddr. If paddr is null, the address is discarded.

Throughout  $\verb"soreceive"$ , if  $\verb"MSG_PEEK"$  is set, the data is not removed from the buffer.

The code in Figure 16.44 processes any control mbufs that are in the buffer.

### Return control information

565-590 Each control mbuf is removed from the buffer (or copied if MSG\_PEEK is set) and attached to \*controlp. If controlp is null, the control information is discarded.

If the process is prepared to receive control information, the protocol has a dom\_externalize function defined, and if the control mbuf contains a SCM\_RIGHTS (access rights) message, the dom\_externalize function is called. This function takes any kernel action associated with receiving the access rights. Only the Unix protocol

518 Socket I/O

Chapter 16

「「「「「「「「「」」」」」

「「「「「「「」」」」

	uipc_socket.c
565	while (m && m->m_type == MT_CONTROL && error == 0) (
566	if (flags & MSG_PEEK) {
567	if (controlp)
568	<pre>*controlp = m_copy(m, 0, m-&gt;m_len);</pre>
569	$m = m - m_next;$
570	} else {
571	<pre>sbfree(&amp;so-&gt;so_rcv, m);</pre>
572	if (controlp) {
573	if (pr->pr_domain->dom_externalize &&
574	<pre>mtod(m, struct cmsghdr *)-&gt;cmsg_type ==</pre>
575	SCM_RIGHTS)
576	error = (*pr->pr_domain->dom_externalize) (m);
577	<pre>*controlp = m;</pre>
578	so->so_rcv.sb_mb = m->m_next;
579	$m-m_next = 0;$
580	<pre>m = so-&gt;so_rcv.sb_mb;</pre>
581	} else {
582	MFREE(m, so->so_rcv.sb_mb);
583	m = so->so_rcv.sb_mb;
584	)
585	}
586	if (controlp) {
587	orig_resid = 0;
588	controlp = &(*controlp)->m_next;
589	}
590	}uipc_socket.c

Figure 16.44 soreceive function: control information.

domain supports access rights, as discussed in Section 7.3. If the process is not prepared to receive control information (controlp is null) the mbuf is discarded.

The loop continues while there are more mbufs with control information and no error has occurred.

For the Unix protocol domain, the dom\_externalize function implements the semantics of passing file descriptors by modifying the file descriptor table of the receiving process.

After the control mbufs are processed, m points to the next mbuf on the chain. If the chain does not contain any mbufs after the address, or after the control information, m is null. This occurs, for example, when a 0-length UDP datagram is queued in the receive buffer. In Figure 16.45 soreceive prepares to transfer the data from the mbuf chain.

### Prepare to transfer data

591-597

After the control mbufs have been processed, the chain should contain regular, outof-band data mbufs or no mbufs at all. If m is null, soreceive is finished with this chain and control drops to the bottom of the while loop. If m is not null, any remaining chains (nextrecord) are reattached to m and the type of the next mbuf is saved in type. If the next mbuf contains OOB data, MSG\_OOB is set in flags, which is later

soreceive Code 519

INTEL Ex.1013.545

		Figure 16 45	uipc_socket.c
597	}		
596		flags (= MSG OOB;	
595		if (type == MT_OOBDATA)	
594		type = m->m_type;	
593		m->m_nextpkt = nextrecord;	
592		if ((flags & MSG_PEEK) == 0)	uipe_socker.e
591	if	(m) {	uipc_socket.c

Figure 16.45 soreceive function: mbuf transfer setup.

returned to the process. Since TCP does not support the MT\_OOBDATA form of out-ofband data, MSG\_OOB will never be returned for reads on TCP sockets.

Figure 16.47 shows the first part of the mbuf transfer loop. Figure 16.46 lists the variables updated within the loop.

Variable	Description
moff offset uio_resid len	the offset of the next byte to transfer when MSG_PEEK is set the offset of the OOB mark when MSG_PEEK is set the number of bytes remaining to be transferred the number of bytes to be transferred from this mbuf; may be less than m_len if uio_resid is small, or if the OOB mark is near

Figure 16.46 soreceive function: loop variables.

<sup>598–600</sup> During each iteration of the while loop, the data in a single mbuf is transferred to the output chain or to the uio buffers. The loop continues while there are more mbufs, the process's buffers are not full, and no error has occurred.

## Check for transition between OOB and regular data

<sup>600–605</sup> If, while processing the mbuf chain, the type of the mbuf changes, the transfer stops. This ensures that regular and out-of-band data are not both returned in the same message. This check does not apply to TCP.

### Update OOB mark

<sup>606-611</sup> The distance to the oobmark is computed and limits the size of the transfer, so the byte before the mark is the last byte transferred. The size of the transfer is also limited by the size of the mbuf. This code does apply to TCP.
 <sup>612-625</sup> If the data is being roturned to the product of the transfer is also limited by the size of the transfer is also limited by the size of the data is being roturned to the product of the transfer is also limited by the size of the t

If the data is being returned to the uio buffers, uiomove is called. If the data is being returned as an mbuf chain, uio\_resid is adjusted to reflect the number of bytes moved.

To avoid suspending protocol processing for a long time, protocol processing is enabled during the call to uiomove. Additional data may appear in the receive buffer because of protocol processing while uiomove is running.

The code in Figure 16.48 adjusts all the pointers and offsets to prepare for the next mbuf.

r 16

ket.c

red

no

s of

:he

ו **is** 

ive

ut-

his

ng in

ter

л÷,

520 Socket I/O

Chapter 16

	uipc_socker
598	moff = 0;
599	offset = 0;
600	while (m && uio->uio_resid > 0 && error == 0) {
601	if $(m \rightarrow m_type == MT_OOBDATA)$ {
602	if (type != MT_OOBDATA)
603	break;
604	} else if (type == MT_OOBDATA)
605	break;
606	so->so_state &= ~SS_RCVATMARK;
607	len = uio->uio_resid;
608	if (so->so_oobmark && len > so->so_oobmark - offset)
609	<pre>len = so-&gt;so_oobmark - offset;</pre>
610	if (len > m->m_len - moff)
611	<pre>len = m-&gt;m_len - moff;</pre>
612	/*
613	* If mp is set, just pass back the mbufs.
614	* Otherwise copy them out via the uio, then free.
615	* Sockbuf must be consistent here (points to current mbuf,
616	* it points to next record) when we drop priority;
617	* we must note any additions to the sockbuf when we
618	* block interrupts again.
619	*/
620	if $(mp == 0)$ {
621	<pre>splx(s);</pre>
622	error = uiomove(mtod(m, caddr_t) + moff, (int) len, uio);
623	s = splnet();
624	} else
625	uio->uio_resid -= len;

Figure 16.47 soreceive function: uiomove.

### Finished with mbuf?

<sup>626–646</sup> If all the bytes in the mbuf have been transferred, the mbuf must be discarded or the pointers advanced. If the mbuf contained the end of a logical record, MSG\_EOR is set. If MSG\_PEEK is set, soreceive skips to the next buffer. If MSG\_PEEK is not set, the buffer is discarded if the data was copied by uiomove, or appended to mp if the data is being returned in an mbuf chain.

### More data to process

<sup>647-657</sup> There may be more data to process in the mbuf if the request didn't consume all the data, if so\_oobmark cut the request short, or if additional data arrived during uiomove. If MSG\_PEEK is set, moff is updated. If the data is to be returned on an mbuf chain, len bytes are copied and attached to the chain. The mbuf pointers and the receive buffer byte count are updated by the amount of data that was transferred.

Figure 16.49 contains the code that handles the OOB offset and the MSG\_EOR processing.

С

е

f

r

g

е

g

n

е

)-

uipc\_socket.c 626 if (len == m->m\_len - moff) { 627 if (m~>m\_flags & M\_EOR) 628 flags |= MSG\_EOR; 629 if (flags & MSG\_PEEK) { 630  $m = m - m_next;$ 631 moff = 0;632 } else { 633 nextrecord = m->m\_nextpkt; 634 sbfree(&so->so\_rcv, m); 635 if (mp) { 636 \*mp = m;637 mp = &m->m\_next; 638 so->so\_rcv.sb\_mb = m = m->m\_next; 639 \*mp = (struct mbuf \*) 0; 640 } else { 641 MFREE(m, so~>so\_rcv.sb\_mb); 642 m = so->so\_rcv.sb\_mb; 643 } 644 if (m) 645 m->m\_nextpkt = nextrecord; 646 } 647 } else { 648 if (flags & MSG\_PEEK) 649 moff += len; 650 else { 651 if (mp) 652 \*mp = m\_copym(m, 0, len, M\_WAIT); 653 m->m\_data += len; 654 m->m\_len -= len; 655 so\_rcv.sb\_cc -= len; 656 } 657 } uipc\_socket.c

Figure 16.48 soreceive function: update buffer.

658	if (so->so_oobmark) {	wipe_socker.e
659	if ((flags & MSG_PEEK) == 0) {	
660	<pre>so-&gt;so_oobmark -= len;</pre>	
661	if $(so->so_obmark == 0)$ {	
662	<pre>so-&gt;so_state  = SS_RCVATMARK;</pre>	
663	break;	
664	}	
665	) else {	
666	offset += len;	
667	if (offset == so->so_oobmark)	
668	break;	
669	}	
670	}	
671	if (flags & MSG_EOR)	
672	break;	

— uipc\_socket.c

Figure 16.49 soreceive function: out-of-band data mark.

### Update OOB mark

658-670 If the out-of-band mark is nonzero, it is decremented by the number of bytes transferred. If the mark has been reached, SS\_RCVATMARK is set and soreceive breaks out of the while loop. If MSG\_PEEK is set, offset is updated instead of so\_oobmark.

### End of logical record

671-672 If the end of a logical record has been reached, soreceive breaks out of the mbuf processing loop so data from the next logical record is not returned with this message.

The loop in Figure 16.50 waits for more data to arrive when MSG\_WAITALL is set and the request is not complete.

673	/*	uipc_socket.c
674	, * If the MSG_WAITALL flag is set (for non-atomic socket),	
675	* we must not quit until "uio->uio_resid == 0" or an erro	
676	* termination. If a signal/timeout occurs, return	1
677	* with a short count but without error.	
678		
679	* Keep sockbuf locked against other readers. */	
680		
	while (flags & MSG_WAITALL && m == 0 && uio->uio_resid > 0	~ & &
681	<pre>!sosendallatonce(so) &amp;&amp; !nextrecord) {</pre>	
682	if (so->so_error    so->so_state & SS_CANTRCVMORE)	
683	break;	
684	error = sbwait(&so->so_rcv);	
685	if (error) {	
686	<pre>sbunlock(&amp;so-&gt;so_rcv);</pre>	
687	<pre>splx(s);</pre>	
688	return (0);	
689	}	
690	if (m = so->so_rcv.sb_mb)	
691	<pre>nextrecord = m-&gt;m_nextpkt;</pre>	
692	)	
693	} /* while more data and more space	to fill */ uipc socket.c

Figure 16.50 soreceive function: MSG\_WAITALL processing.

#### MSG\_WAITALL

673-681 If MSG\_WAITALL is set, there is no more data in the receive buffer (m equals 0), the caller wants more data, sosendallatonce is false, and this is the last record in the receive buffer (nextrecord is null), then soreceive must wait for additional data.

### Error or no more data will arrive

682-683 If an error is pending or the connection is closed, the loop is terminated.

### Wait for data to arrive

684-689 sbwait returns when the receive buffer is changed by the protocol layer. If the wait was interrupted by a signal (error is nonzero), sosend returns immediately.

690-692

693

nates.

707

708 709

710

711

712

713

714

Process next mbuf

: 16

### out

# buf

2

set

# ket.c

1 \*/

cket.c

, the

ι the

f the

a.

694	if (m && pr->pr_flags & PR_ATOMIC) {
695	<pre>flags  = MSG_TRUNC;</pre>
696	if ((flags & MSG_PEEK) == 0)
697	<pre>(void) sbdroprecord(&amp;so-&gt;so_rcv);</pre>
698	}
699	if ((flags & MSG_PEEK) == 0) {
700	if (m == 0)
701	<pre>so-&gt;so_rcv.sb_mb = nextrecord;</pre>
702	if (pr->pr_flags & PR_WANTRCVD && so->so_pcb)
703	(*pr->pr_usrreq) (so, PRU_RCVD, (struct mbuf *) 0,
704	(struct mbuf *) flags, (struct mbuf *
705	<pre>(struct mbuf *) 0);</pre>
706	}

Synchronize m and nextrecord with receive buffer

and no error has occurred, the loop continues.

\*flagsp != flags;

if (orig\_resid == uio->uio\_resid && orig\_resid && (flags & MSG\_EOR) == 0 && (so->so\_state & SS\_CANTRCVMORE) == 0) { sbunlock(&so->so\_rcv); splx(s); goto restart; if (flagsp)

m and nextrecord are updated, since the receive buffer has been modified by the

This is the end of the mbuf processing loop. Control returns to the loop starting on

line 600 (Figure 16.47). As long as there is data in the receive buffer, more space to fill,

protocol layer. If data arrived in the mbuf, m will be nonzero and the while loop termi-

### uipc\_socket.c

uipc\_socket.c

1.0.0.0

Figure 16.51 soreceive function: cleanup.

### Truncated message

}

If the process received a partial message (a datagram or a record) because its receive 694-698 buffer was too small, the process is notified by setting MSG\_TRUNC and the remainder of the message is discarded. MSG\_TRUNC (as with all receive flags) is available only to a process through the recymsg system call, even though soreceive always sets the flags.

### End of record processing

If MSG\_PEEK is not set, the next mbuf chain is attached to the receive buffer and, if 699-706 required, the protocol is notified that the receive operation has been completed by issuing the PRU\_RCVD protocol request. TCP uses this feature to update the receive window for the connection.

「日本の時間のない」

運動がはなない。

### Nothing transferred

- <sup>707-712</sup> If soreceive runs to completion, no data is transferred, the end of a record is not reached, and the read-half of the connection is still active, then the buffer is unlocked and soreceive jumps back to restart to continue waiting for data.
- 713-714 Any flags set during soreceive are returned in \*flagsp, the buffer is unlocked, and soreceive returns.

### Analysis

soreceive is a complex function. Much of the complication is because of the intricate manipulation of pointers and the multiple types of data (out-of-band, address, control, regular) and multiple destinations (process buffers, mbuf chain).

Similar to sosend, soreceive has collected features over the years. A specialized receive function for each protocol would blur the boundary between the socket layer and the protocol layer, but it would simplify the code considerably.

[Partridge and Pink 1993] describe the creation of a custom soreceive function for UDP to checksum datagrams while they are copied from the receive buffer to the process. They note that modifying the generic soreceive function to support this feature would "make the already complicated socket routines even more complex."

### 16.13 select System Call

In the following discussion we assume that the reader is familiar with the basic operation and semantics of select. For a detailed discussion of the application interface to select see [Stevens 1992].

Figure 16.52 shows the conditions detected by using select to monitor a socket.

	Detected by selecting for:			
Description	reading	writing	exceptions	
data available for reading	•			
read-half of connection is closed	•			
listen socket has queued connection	٠	1		
socket error is pending	•			
space available for writing and a connection exists or is not required write-half of connection is closed socket error is pending		•		
OOB synchronization mark is pending			•	

Figure 16.52 select system call: socket events.

We start with the first half of the select system call, shown in Figure 16.53.

Section 16.13

l6

Эť

d

ĺ,

te ›L

d

er

)r )-

٢e

3-

:0

### Validation and setup

Two arrays of three descriptor sets are allocated on the stack: ibits and obits. They are cleared by bzero. The first argument, nd, must be no larger than the maximum number of descriptors associated with the process. If nd is more than the number of descriptors currently allocated to the process, it is reduced to the current allocation. ni is set to the number of bytes needed to store a bit mask with nd bits (1 bit for each descriptor). For example, if the maximum number of descriptors is 256 (FD\_SETSIZE), fd\_set is represented as an array of 32-bit integers (NFDBITS), and nd is 65, then:

 $ni = howmany(65, 32) \times 4 = 3 \times 4 = 12$ 

where how many (x, y) returns the number of y-bit objects required to store x bits.

### Copy file descriptor sets from process

411-418 The getbits macro uses copyin to transfer the file descriptor sets from the process to the three descriptor sets in ibits. If a descriptor set pointer is null, nothing is copied from the process.

### Setup timeout value

419-438 If tv is null, timo is set to 0 and select will wait indefinitely. If tv is not null, the timeout value is copied into the kernel and rounded up to the resolution of the hardware clock by itimerfix. The current time is added to the timeout value by timevaladd. The number of clock ticks until the timeout is computed by hzto and saved in timo. If the resulting timeout is 0, timo is set to 1. This prevents select from blocking and implements the nonblocking semantics of an all-0s timeval structure.

The second half of select, shown in Figure 16.54, scans the file descriptors indicated by the process and returns when one or more become ready, or the timer expires, or a signal occurs.

### Scan file descriptors

439-442 The loop that starts at retry continues until select can return. The current value of the global integer nselcoll is saved and the P\_SELECT flag is set in the calling process's control block. If either of these change while selscan (Figure 16.55) is checking the file descriptors, it indicates that the status of a descriptor has changed because of interrupt processing and select must rescan the descriptors. selscan looks at every descriptor set in the three input descriptor sets and sets the matching descriptor in the output set if the descriptor is ready.

### Error or some descriptors are ready

443–444 Return immediately if an error occurred or if a descriptor is ready.

### **Timeout expired?**

<sup>445–451</sup> If the process supplied a time limit and the current time has advanced beyond the timeout value, return immediately.

and the second second

と言語が

言語後近にという

sys\_generic.c

```
sys_generic.c
390 struct select_args {
391
       u_int nd;
392
        fd_set *in, *ou, *ex;
393
        struct timeval *tv;
394 );
395 select(p, uap, retval)
396 struct proc *p;
397 struct select_args *uap;
398 int
          *retval;
399 {
        fd_set ibits[3], obits[3];
400
        struct timeval atv;
401
       int
                s, ncoll, error = 0, timo;
402
403
       u_int
               ni;
       bzero((caddr_t) ibits, sizeof(ibits));
404
405
       bzero((caddr_t) obits, sizeof(obits));
       if (uap->nd > FD_SETSIZE)
406
407
            return (EINVAL);
408
        if (uap->nd > p->p_fd->fd_nfiles)
           uap->nd = p->p_fd->fd_nfiles;
409
                                            /* forgiving; slightly wrong */
        ni = howmany(uap->nd, NFDBITS) * sizeof(fd_mask);
410
411 #define getbits(name, x) \
412
       if (uap->name && \
413
            (error = copyin((caddr_t)uap->name, (caddr_t)&ibits[x], ni))) \
414
            goto done;
415
        getbits(in, 0);
416
        getbits(ou, 1);
417
        getbits(ex, 2);
418 #undef getbits
419
       if (uap->tv) {
420
            error = copyin((caddr_t) uap->tv, (caddr_t) & atv,
421
                           sizeof(atv));
422
            if (error)
423
                goto done;
424
            if (itimerfix(&atv)) {
425
                error = EINVAL;
426
                goto done;
427
            }
428
            s = splclock();
429
            timevaladd(&atv, (struct timeval *) &time);
430
            timo = hzto(&atv);
431
            /*
432
             * Avoid inadvertently sleeping forever.
             */
433
            if (timo == 0)
434
435
                timo = 1;
436
            splx(s);
437
        } else
438
            timo = 0;
```

Figure 16.53 select function: initialization.

Section 16.13

1eric.c

ter 16

sys\_generic.c 439 retry: 440ncoll = nselcoll; p->p\_flag != P\_SELECT; 441 error = selscan(p, ibits, obits, uap->nd, retval); 442 if (error || \*retval) 443 444 goto done; 445 s = splhigh(); /\* this should be timercmp(&time, &atv, >=) \*/ 446 if (uap->tv && (time.tv\_sec > atv.tv\_sec || 447 time.tv\_sec == atv.tv\_sec && time.tv\_usec >= atv.tv\_usec)) { 448 splx(s); 449 goto done; 450 451 } if ((p->p\_flag & P\_SELECT) == 0  $\parallel$  nselcoll != ncoll) { 452 splx(s); 453 454 goto retry; } 455 p->p\_flag &= ~P\_SELECT; 456 error = tsleep((caddr\_t) & selwait, PSOCK | PCATCH, "select", timo); 457 458 splx(s); if (error == 0)459 460 goto retry; 461 done: p->p\_flag &= ~P\_SELECT; 462 /\* select is not restarted after signals... \*/ 463 if (error == ERESTART) 464 error = EINTR; 465 if (error == EWOULDBLOCK) 466 error = 0;467 468 #define putbits(name, x) \ 469 if (uap->name && \ (error2 = copyout((caddr\_t)&obits[x], (caddr\_t)uap->name, ni))) \ 470 error = error2; 471 if (error == 0) { 472 error2; int 473 474 putbits(in, 0); putbits(ou, 1); 475 putbits(ex, 2); 476 477 #undef putbits 478 } return (error); 479 480 } -sys\_generic.c Figure 16.54 select function: second half.

zeneric.c

INTEL Ex.1013.553

25

#### Status changed during selscan

452-455 selscan can be interrupted by protocol processing. If the socket is modified during the interrupt, P\_SELECT and nselcoll are changed and select must rescan the descriptors.

#### Wait for buffer changes

456-460 All processes calling select use selwait as the wait channel when they call tsleep. With Figure 16.60 we show that this causes some inefficiencies if more than one process is waiting for the same socket buffer. If tsleep returns without an error, select jumps to retry to rescan the descriptors.

#### Ready to return

- 461-480 At done, P\_SELECT is cleared, ERESTART is changed to EINTR, and EWOULDBLOCK is changed to 0. These changes ensure that EINTR is returned when a signal occurs during select and 0 is returned when a timeout occurs.
  - The output descriptor sets are copied back to the process and select returns.

### selscan Function

The heart of select is the selscan function shown in Figure 16.55. For every bit set in one of the three descriptor sets, selscan computes the descriptor associated with the bit and dispatches control to the fo\_select function associated with the descriptor. For sockets, this is the soo\_select function.

### Locate descriptors to be monitored

<sup>481–496</sup> The first for loop iterates through each of the three descriptor sets: read, write, and exception. The second for loop interates within each descriptor set. This loop is executed once for every 32 bits (NFDBITS) in the set.

The inner while loop checks all the descriptors identified by the 32-bit mask extracted from the current descriptor set and stored in bits. The function ffs returns the position within bits of the first 1 bit, starting at the low-order bit. For example, if bits is 1000 (with 28 leading 0s), ffs (bits) is 4.

### Poll descriptor

497-500 From i and the return value of ffs, the descriptor associated with the bit is computed and stored in fd. The bit is cleared in bits (but not in the input descriptor set), the file structure associated with the descriptor is located, and fo\_select is called.

The second argument to fo\_select is one of the elements in the flag array. msk is the index of the outer for loop. So the first time through the loop, the second argument is FREAD, the second time it is FWRITE, and the third time it is 0. EBADF is returned if the descriptor is not valid.

### Descriptor is ready

- 501-504 When a descriptor is found to be ready, the matching bit is set in the output descriptor set and n (the number of matches) is incremented.
- <sup>505-510</sup> The loops continue until all the descriptors are polled. The number of ready descriptors is returned in \*retval.

Section 16.13

sys\_generic.c 481 selscan(p, ibits, obits, nfd, retval) 482 struct proc \*p; 483 fd\_set \*ibits, \*obits; 484 int nfd, \*retval; 485 { 486 struct filedesc \*fdp = p->p\_fd; 487 int msk, i, j, fd; 488 fd\_mask bits; 489 struct file \*fp; 490 int n = 0;491 static int flag[3] = 492 {FREAD, FWRITE, 0}; 493 for (msk = 0; msk < 3; msk++) {</pre> 494 for (i = 0; i < nfd; i += NFDBITS) { 495 bits = ibits[msk].fds\_bits[i / NFDBITS]; 496 while ((j = ffs(bits)) && (fd = i + --j) < nfd){ 497 bits &= ~(1 << j); 498 fp = fdp->fd\_ofiles[fd]; 499 if (fp == NULL) 500 return (EBADF); 501 if ((\*fp->f\_ops->fo\_select) (fp, flag[msk], p)) { 502 FD\_SET(fd, &obits[msk]); 503 n++; 504 } 505 } 506 } 507 } 508 \*retval = n; 509 return (0); 510 } sys\_generic.c

Figure 16.55 selscan function.

### soo\_select Function

For every descriptor that selscan finds in the input descriptor sets, it calls the function referenced by the fo\_select pointer in the fileops structure (Section 15.5) associated with the descriptor. In this text, we are interested only in socket descriptors and the soo\_select function shown in Figure 16.56.

105-112 Each time soo\_select is called, it checks the status of only one descriptor. If the descriptor is ready relative to the conditions specified in which, the function returns 1 immediately. If the descriptor is not ready, selrecord marks either the socket's receive or send buffer to indicate that a process is selecting on the buffer and then soo\_select returns 0.

Figure 16.52 showed the read, write, and exceptional conditions for sockets. Here we see that the macros soreadable and sowriteable are consulted by soo\_select. These macros are defined in sys/socketvar.h.

16

all an

or,

ur-

the

CK

ısk

rns :, if

et),

ısk

zu-

' is

'ip-

ıdy

前推荐

sys\_socket.c

sys\_socket.c

```
105 soo_select(fp, which, p)
106 struct file *fp;
            which;
107 int
108 struct proc *p;
109 {
        struct socket *so = (struct socket *) fp->f_data;
110
                s = splnet();
111
        int
112
        switch (which) {
        case FREAD:
113
            if (soreadable(so)) {
114
                splx(s);
115
                return (1);
116
            }
117
            selrecord(p, &so->so_rcv.sb_sel);
118
            so->so_rcv.sb_flags != SB_SEL;
119
            break;
120
121
        case FWRITE:
            if (sowriteable(so)) {
122
                 splx(s);
123
124
                 return (1);
125
            }
             selrecord(p, &so->so_snd.sb_sel);
126
             so->so_snd.sb_flags |= SB_SEL;
127
            break;
128
        case 0:
129
             if (so->so_oobmark || (so->so_state & SS_RCVATMARK)) {
130
131
                 splx(s);
                 return (1);
132
133
             }
             selrecord(p, &so->so_rcv.sb_sel);
134
             so->so_rcv.sb_flags |= SB_SEL;
135
             break;
136
137
         }
         splx(s);
 138
         return (0);
 139
 140 }
```

Figure 16.56 soo\_select function.

# Is socket readable?

The soreadable macro is:

113-120

#define soreadable(so) \
 ((so)->so\_rcv.sb\_cc >= (so)->so\_rcv.sb\_lowat || \
 ((so)->so\_state & SS\_CANTRCVMORE) || \
 (so)->so\_glen || (so)->so\_error)

Since the receive low-water mark for UDP and TCP defaults to 1 (Figure 16.4), the socket is readable if any data is in the receive buffer, if the read-half of the connection is closed, if any connections are ready to be accepted, or if there is an error pending.

ka sa

### Is socket writeable?

```
121-128
```

```
The sowriteable macro is:
```

```
#define sowriteable(so) \
  (sbspace(&(so)->so_snd) >= (so)->so_snd.sb_lowat && \
  (((so)->so_state&SS_ISCONNECTED) || \
    ((so)->so_proto->pr_flags&PR_CONNREQUIRED)==0) || \
  ((so)->so_state & SS_CANTSENDMORE) || \
  (so)->so_error)
```

The default send low-water mark for UDP and TCP is 2048. For UDP, sowriteable is always true because sbspace is always equal to sb\_hiwat, which is always greater than or equal to sb\_lowat, and a connection is not required.

For TCP, the socket is not writeable when the free space in the send buffer is less than 2048 bytes. The other cases are described in Figure 16.52.

# Are there any exceptional conditions pending?

```
129-140
```

For exceptions, so\_oobmark and the SS\_RCVATMARK flags are examined. An exceptional condition exists until the process has read past the synchronization mark in the data stream.

# selrecord Function

Figure 16.57 shows the definition of the selinfo structure stored with each send and receive buffer (the sb\_sel member from Figure 16.3).

```
      41 struct selinfo {
      select.h

      42 pid_t si_pid;
      /* process to be notified */

      43 short si_flags;
      /* 0 or SI_COLL */

      44 };
      select.h
```

Figure 16.57 selinfo structure.

41-44 When only one process has called select for a given socket buffer, si\_pid is the process ID of the waiting process. When additional processes call select on the same buffer, SI\_COLL is set in si\_flags. This is called a *collision*. This is the only flag currently defined for si\_flags.

The selrecord function shown in Figure 16.58 is called when soo\_select finds a descriptor that is not ready. The function records enough information so that the process is awakened by the protocol processing layer when the buffer changes.

### Already selecting on this descriptor

522-531

The first argument to selrecord points to the proc structure for the selecting process. The second argument points to the selinfo record to update (so\_snd.sb\_sel or so\_rcv.sb\_sel). If this process is already recorded in the selinfo record for this socket buffer, the function returns immediately. For example, the process called select with the read and exception bits set for the same descriptor. 532 Socket I/O

Chapter 16

sys\_generic.c

```
522 void
523 selrecord(selector, sip)
524 struct proc *selector;
525 struct selinfo *sip;
526 {
        struct proc *p;
527
       pid_t mypid;
528
        mypid = selector->p_pid;
529
530
        if (sip->si_pid == mypid)
531
           return;
        if (sip->si_pid && (p = pfind(sip->si_pid)) &&
532
533
            p->p_wchan == (caddr_t) & selwait)
            sip->si_flags {= SI_COLL;
534
535
        else
536
            sip->si_pid = mypid;
537 }
```

– sys\_generic.c

Figure 16.58 selrecord function.

### Select collision with another process?

<sup>532–534</sup> If another process is already selecting on this buffer, SI\_COLL is set.

### No collision

<sup>535–537</sup> If there is no other process already selecting on this buffer, si\_pid is 0 so the ID of the current process is saved in si\_pid.

### selwakeup Function

When protocol processing changes the state of a socket buffer and only one process is selecting on the buffer, Net/3 can immediately put that process on the run queue based on the information it finds in the selinfo structure.

When the state changes and there is more than one process selecting on the buffer (SI\_COLL is set), Net/3 has no way of determining the set of processes interested in the buffer. When we discussed the code in Figure 16.54, we pointed out that *every* process that calls select uses selwait as the wait channel when calling tsleep. This means the corresponding wakeup will schedule *all* the processes that are blocked in select—even those that are not interested in activity on the buffer.

Figure 16.59 shows how selwakeup is called.

The protocol processing layer is responsible for notifying the socket layer by calling one of the functions listed at the bottom of Figure 16.59 when an event occurs that changes the state of a socket. The three functions shown at the bottom of Figure 16.59 cause selwakeup to be called and any process selecting on the socket to be scheduled to run.

selwakeup is shown in Figure 16.60.

541-548 If si\_pid is 0, there is no process selecting on the buffer and the function returns immediately.

Section 16.13

16

c.c

ic.c

of

is

ed

fer

he

2SS

ns

in

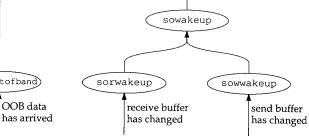
ng

ıat

59

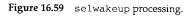
ed

ns



select System Call

533



sys\_generic.c 541 void 542 selwakeup(sip) 543 struct selinfo \*sip; 544 { 545 struct proc \*p; 546 int s; 547 if (sip->si\_pid == 0) 548 return; 549 if (sip->si\_flags & SI\_COLL) { 550 nselcoll++; 551 sip->si\_flags &= ~SI\_COLL; 552 wakeup((caddr\_t) & selwait); 553 } 554 p = pfind(sip->si\_pid); 555 sip->si\_pid = 0; 556 if (p != NULL) { 557 s = splhigh(); 558 if (p->p\_wchan == (caddr\_t) & selwait) { 559 if (p->p\_stat == SSLEEP) 560 setrunnable(p); 561 else 562 unsleep(p); 563 } else if (p->p\_flag & P\_SELECT) 564 p->p\_flag &= ~P\_SELECT; 565 splx(s); 566 } 567 } sys\_generic.c

Figure 16.60 selwakeup function.

INTEL Ex.1013.559

the state of the second se

# Wake all processes during a collision

549-553 If more than one process is selecting on the affected socket, nselcoll is incremented, the collision flag is cleared, and every process blocked in select is awakened. As mentioned with Figure 16.54, nselcoll forces select to rescan the descriptors if the buffers change before the process has blocked in tsleep (Exercise 16.9).

554-567

If the process identified by si\_pid is waiting on selwait, it is scheduled to run. If the process is waiting on some other wait channel, the P\_SELECT flag is cleared. The process can be waiting on some other wait channel if selrecord is called for a valid descriptor and then selscan finds a bad file descriptor in one of the descriptor sets. selscan returns EBADF, but the previously modified selinfo record is not reset. Later, when selwakeup runs, selwakeup may find the process identified by sel\_pid is no longer waiting on the socket buffer so the selinfo information is ignored.

Only one process is awakened during selwakeup unless multiple processes are sharing the same descriptor (i.e., the same socket buffers), which is rare. On the machines to which the authors had access, nselcoll was always 0, which confirms the statement that select collisions are rare.

# 16.14 Summary

In this chapter we looked at the read, write, and select system calls for sockets.

We saw that sosend handles all output between the socket layer and the protocol processing layer and that soreceive handles all input.

The organization of the send buffer and receive buffers was described, as well as the default values and semantics of the high-water and low-water marks for the buffers.

The last part of the chapter discussed the implementation of select. We showed that when only one process is selecting on a descriptor, the protocol processing layer will awaken only the process identified in the selinfo structure. When there is a collision and more than one process is selecting on a descriptor, the protocol layer has no choice but to awaken every process that is selecting on *any* descriptor.

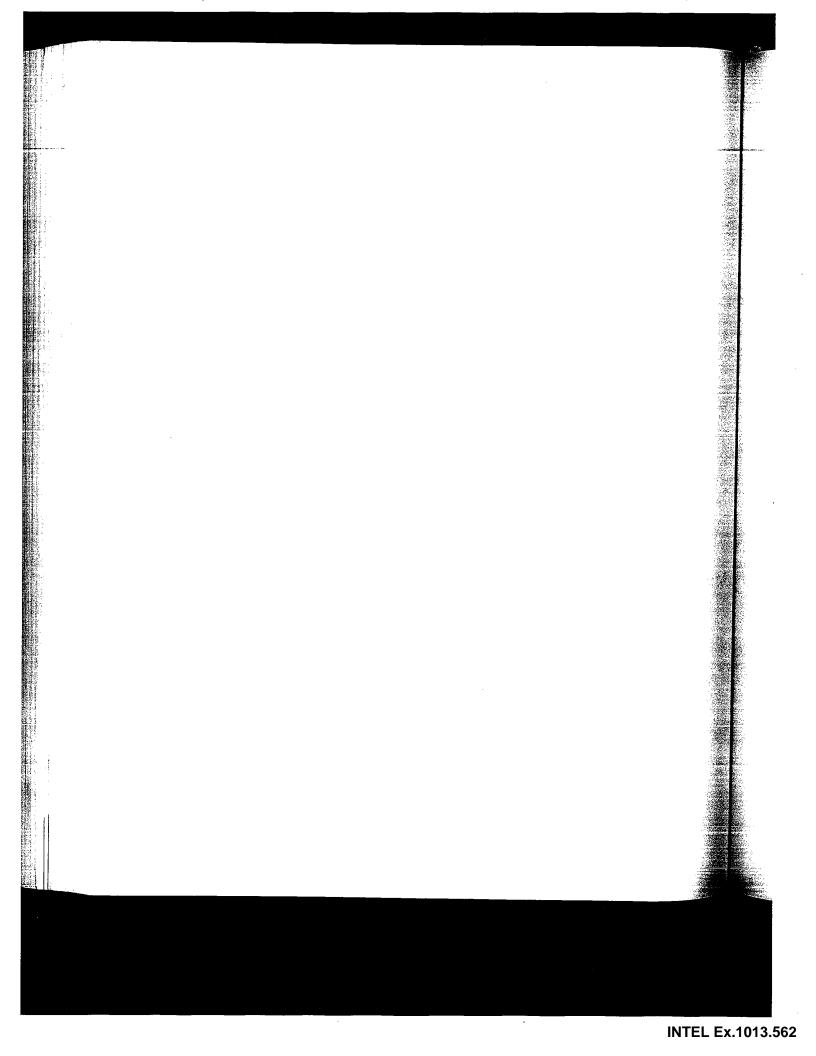
# Exercises

- 16.1 What happens to resid in sosend when an unsigned integer larger than the maximum positive signed integer is passed in the write system call?
- 16.2 When sosend puts less than MCLBYTES of data in a cluster, space is reduced by the full MCLBYTES and may become negative, which terminates the loop that fills mbufs for atomic protocols. Is this a problem?
- **16.3** Datagram and stream protocols have very different semantics. Divide the sosend and soreceive functions each into two functions, one to handle messages, and one to handle streams. Other than making the code clearer, what are the advantages of making this change?
- 16.4 For PR\_ATOMIC protocols, each write call specifies an implicit message boundary. The

R

socket layer delivers the message as a single unit to the protocol. The MSG\_EOR flag allows a process to specify explicit message boundaries. Why is the implicit technique insufficient?

- **16.5** What happens when sosend cannot immediately acquire a lock on the send buffer when the socket descriptor is marked as nonblocking and the process does not specify MSG\_DONTWAIT?
- **16.6** Under what circumstances would sb\_cc < sb\_hiwat yet sbspace would report no free space? Why should a process be blocked in this case?
- **16.7** Why isn't the length of a control message copied back to the process by recvit as is the name length?
- 16.8 Why does soreceive clear MSG\_EOR?
- 16.9 What might happen if the nselcoll code were removed from select and selwakeup?
- **16.10** Modify the select system call to return the time remaining in the timer when select returns.



# 17

# Socket Options

# 17.1 Introduction

We complete our discussion of the socket layer in this chapter by discussing several system calls that modify the behavior of sockets.

The setsockopt and getsockopt system calls were introduced in Section 8.8, where we described the options that provide access to IP features. In this chapter we show the implementation of these two system calls and the socket-level options that are controlled through them.

The ioctl function was introduced in Section 4.4, where we described the protocol-independent ioctl commands for network interface configuration. In Section 6.7 we described the IP specific ioctl commands used to assign network masks as well as unicast, broadcast, and destination addresses. In this chapter we describe the implementation of ioctl and the related features of the fcntl function.

Finally, we describe the getsockname and getpeername system calls, which return address information for sockets and connections.

Figure 17.1 shows the functions that implement the socket option system calls. The shaded functions are described in this chapter.

537



「東安山」の

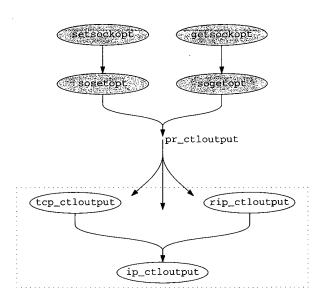


Figure 17.1 setsockopt and getsockopt system calls.

# 17.2 Code Introduction

The code in this chapter comes from the four files listed in Figure 17.2.

File	Description
kern/kern_descrip.c	fcntl system call
kern/uipc_syscalls.c	setsockopt, getsockopt, getsockname, and getpeername system calls
kern/uipc_socket.c	socket layer processing for setsockopt and getsockopt
kern/sys_socket.c	ioctl system call for sockets

Figure 17.2 Files discussed in this chapter.

# **Global Variables and Statistics**

No new global variables are introduced and no statistics are collected by the system calls we describe in this chapter.

Section 17.3

# 17.3 setsockopt System Call

Figure 8.29 listed the different protocol levels that can be accessed with this function (and with getsockopt). In this chapter we focus on the SOL\_SOCKET level options, which are listed in Figure 17.3.

optname	optval type	Variable	Description
SO_SNDBUF	int	so_snd.sb_hiwat	send buffer high-water mark
SO_RCVBUF	int	so_rcv.sb_hiwat	receive buffer high-water mark
SO_SNDLOWAT	int	so_snd.sb_lowat	send buffer low-water mark
SO_RCVLOWAT	int	so_rcv.sb_lowat	receive buffer low-water mark
SO_SNDTIMEO	struct timeval	so_snd.sb_timeo	send timeout
SO_RCVTIMEO	struct timeval	so_rcv.sb_timeo	receive timeout
SO_DEBUG	int	so_options	record debugging information for this socket
SO_REUSEADDR	int	so_options	socket can reuse a local address
SO_REUSEPORT	int	so_options	socket can reuse a local port
SO_KEEPALIVE	int	so_options	protocol probes idle connections
SO_DONTROUTE	int	so_options	bypass routing tables
SO_BROADCAST	int	so_options	socket allows broadcast messages
SO_USELOOPBACK	int	so_options	routing domain sockets only; sending
			process receives its own routing
			messages
SO_OOBINLINE	int	so_options	protocol queues out-of-band data inline
SO_LINGER	struct linger	so_linger	socket lingers on close
SO_ERROR	int	so_error	get error status and clear; getsockopt only
SO_TYPE	int	so_type	get socket type; getsockopt only
other			ENOPROTOOPT returned

Figure 17.3 setsockopt and getsockopt options.

The prototype for setsockopt is

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

Figure 17.4 shows the code for this system call.

565-597 getsock locates the file structure for the socket descriptor. If val is nonnull, valsize bytes of data are copied from the process into an mbuf allocated by m\_get. The data associated with an option can be no more than MLEN bytes in length, so if valsize is larger than MLEN, then EINVAL is returned. sosetopt is called and its value is returned.

:er 17

540 Socket Options

Chapter 17

			- uipc_syscalls.c
565 stru	ict sets	ockopt_args {	
566	int	S;	
567	int	level;	
568	int	name;	
569	caddr_t	val;	
570	int	valsize;	
571 };			
572 sets	sockopt (	p, uap, retval)	
573 stru			
574 stru	uct sets	ockopt_args *uap;	
575 int		val;	
576 {			
577		file *fp;	
578	struct	mbuf *m = NULL;	
579	int	error;	
580		ror = getsock(p->p_fd, uap->s, &fp))	
581		urn (error);	
582	if (uar	p->valsize > MLEN)	
583	ret	curn (EINVAL);	
584		p->val) {	
585		= m_get(M_WAIT, MT_SOOPTS);	
586	if	(m = NULL)	
587		return (ENOBUFS);	
588	if	(error = copyin(uap->val, mtod(m, caddr_t),	
589		(u_int) uap->valsize)) {	
590		(void) m_free(m);	
591		return (error);	
592	}		
593		<pre>&gt;m_len = uap-&gt;valsize;</pre>	
594	}		
595	return	<pre>(sosetopt((struct socket *) fp-&gt;f_data, uap-&gt;level,</pre>	
596		uap->name, m));	
597 }			— uipc_syscalls.c

Figure 17.4 setsockopt system call.

# sosetopt Function

This function processes all the socket-level options and passes any other options to the pr\_ctloutput function for the protocol associated with the socket. Figure 17.5 shows an overview of the function.

- <sup>752-764</sup> If the option is not for the socket level (SOL\_SOCKET), the PRCO\_SETOPT request is issued to the underlying protocol. Note that the protocol's pr\_ctloutput function is being called and not its pr\_usrreq function. Figure 17.6 shows which function is called for the Internet protocols.
- 765 The switch statement handles the socket-level options.
- 841-844 An unrecognized option causes ENOPROTOOPT to be returned after the mbuf holding the option is released.

**INTEL Ex.1013.566** 

Section 17.3 pter 17 /scalls.c 757 758 \_syscalls.c ns to the .5 shows

equest is nction is nction is

ouf hold-

uipc\_socket.c 752 sosetopt(so, level, optname, m0) 753 struct socket \*so; level, optname; 754 int 755 struct mbuf \*m0; 756 { error = 0;int struct mbuf \*m = m0; if (level != SOL\_SOCKET) { 759 if (so->so\_proto && so->so\_proto->pr\_ctloutput) 760 return ((\*so->so\_proto->pr\_ctloutput) 761 (PRCO\_SETOPT, so, level, optname, &m0)); 762 error = ENOPROTOOPT; 763 764 } else { switch (optname) { 765 /\* socket option processing \*/ 841 default: error = ENOPROTOOPT; 842 break; 843 844 } if (error == 0 && so->so\_proto && so->so\_proto->pr\_ctloutput) ( 845 (void) ((\*so->so\_proto->pr\_ctloutput) 846 (PRCO\_SETOPT, so, level, optname, &m0)); 847 /\* freed by protocol \*/ m = NULL; 848 849 } 850 } bad: 851 852 if (m) (void) m\_free(m); 853 return (error); 854 855 } - uipc\_socket.c

Figure 17.5 sosetopt function.

Protocol	pr_ctloutput Function	Reference
UDP TCP	ip_ctloutput tcp_ctloutput	Section 8.8 Section 30.6
ICMP IGMP raw IP	rip_ctloutput and ip_ctloutput	Section 8.8 and Section 32.8

Figure 17.6 pr\_ctloutput functions.

Unless an error occurs, control always falls through the switch, where the option 845-855 is passed to the associated protocol in case the protocol layer needs to respond to the request as well as the socket layer. None of the Internet protocols expect to process the socket-level options.

# INTEL Ex.1013.567

「「「「「「「「「「」」」」」」」「「「「「」」」」」」」

uipc socket.c

Notice that the return value from the call to the pr\_ctloutput function is explicitly discarded in case the option is not expected by the protocol. m is set to null to avoid the call to m\_free, since the protocol layer is responsible for releasing the mbuf.

Figure 17.7 shows the linger option and the options that set a single flag in the socket structure.

		uipc_socket.c
766	case SO_LINGER:	, _
767	if (m == NULL    m->m_len != sizeof(struct linger)) {	
768	error = EINVAL;	
769	goto bad;	
770	}	
771	so->so_linger = mtod(m, struct linger *)->l_linger;	
772	/* fall thru */	
773	case SO_DEBUG:	
774	case SO_KEEPALIVE:	
775	case SO_DONTROUTE:	
776	case SO_USELOOPBACK:	
777	case SO_BROADCAST:	
778	case SO_REUSEADDR:	
779	case SO_REUSEPORT:	
780	case SO_OOBINLINE:	
781	if (m == NULL    m->m_len < sizeof(int)) {	
782	error = EINVAL;	
783	goto bad;	
784	}	
785	if (*mtod(m, int *))	
786	so->so_options  = optname;	
787	else	
788	so->so_options &= ~optname;	
789	break;	

Figure 17.7 sosetopt function: linger and flag options.

766-772

The linger option expects the process to pass a linger structure:

struct linger {
 int l\_onoff; /\* option on/off \*/
 int l\_linger; /\* linger time in seconds \*/
};

After making sure the process has passed data and it is the size of a linger structure, the l\_linger member is copied into so\_linger. The option is enabled or disabled after the next set of case statements. so\_linger was described in Section 15.15 with the close system call.

773–789

These options are boolean flags set when the process passes a nonzero value and cleared when 0 is passed. The first check makes sure an integer-sized object (or larger) is present in the mbuf and then sets or clears the appropriate option.

setsockopt System Call 543 Section 17.3 apter 17 Figure 17.8 shows the socket buffer options. explico avoid uipc\_socket.c 790 case SO\_SNDBUF: 791 case SO\_RCVBUF: case SO\_SNDLOWAT: 792 3 in the case SO\_RCVLOWAT: 793 if (m == NULL || m->m\_len < sizeof(int)) { 794 error = EINVAL; 795 : socket.c goto bad; 796 797 } switch (optname) { 798 799 case SO\_SNDBUF: case SO\_RCVBUF: 800 801 if (sbreserve(optname == SO\_SNDBUF ? &so->so\_snd : &so->so\_rcv, 802  $(u long) * mtod(m, int *)) == 0) {$ 803 error = ENOBUFS; 804 805 goto bad; } 806 807 break; 808 case SO SNDLOWAT: so\_so\_snd.sb\_lowat = \*mtod(m, int \*); 809 810 break; case SO\_RCVLOWAT: 811 so->so\_rcv.sb\_lowat = \*mtod(m, int \*); 812 813 break; } 814 815 break; uipc\_socket.c Figure 17.8 sosetopt function: socket buffer options. >c\_socket.c This set of options changes the size of the send and receive buffers in a socket. The 790-815 first test makes sure the required integer has been provided for all four options. For SO\_SNDBUF and SO\_RCVBUF, sbreserve adjusts the high-water mark but does no buffer allocation. For SO\_SNDLOWAT and SO\_RCVLOWAT, the low-water marks are adjusted. Figure 17.9 shows the timeout options. The timeout value for SO\_SNDTIMEO and SO\_RCVTIMEO is specified by the process 816-824 in a timeval structure. If the right amount of data is not available, EINVAL is er strucreturned. d or dis-The time interval stored in the timeval structure must be small enough so that ion 15.15 825-830 when it is represented as clock ticks, it fits within a short integer, since sb\_timeo is a short integer. alue and The code on line 826 is incorrect. The time interval cannot be represented as a short or larger) integer if:

**INTEL Ex.1013.569** 

2

816	case SO_SNDTIMEO:	—— uipc_socket.c
817	case SO_RCVTIMEO:	
	_	
818	{ struct timeval *tv;	
819		
820	short val;	
821	if (m == NULL    m->m_len < sizeof(*tv)) {	
822	error = EINVAL;	
823	goto bad;	
824	}	
825	<pre>tv = mtod(m, struct timeval *);</pre>	
826	if (tv->tv_sec > SHRT_MAX / hz - hz) {	
827	error = EDOM;	
828	goto bad;	
829	}	
830	<pre>val = tv-&gt;tv_sec * hz + tv-&gt;tv_usec / tick;</pre>	
831	switch (optname) {	
832	case SO_SNDTIMEO:	
833	<pre>so-&gt;so_snd.sb_timeo = val;</pre>	
834	break;	
835	case SO_RCVTIMEO:	
836	<pre>so-&gt;so_rcv.sb_timeo = val;</pre>	
837	break;	
838	}	
839	break;	
840	}	——— uipc_socket.

Figure 17.9 sosetopt function: timeout options.

$$tv\_sec \times hz + \frac{tv\_usec}{tick} > SHRT\_MAX$$

where

tick = 
$$\frac{1,000,000}{hz}$$
 and SHRT\_MAX = 32767

So EDOM should be returned if

$$tv\_sec > \frac{SHRT\_MAX}{hz} - \frac{tv\_usec}{tick \times hz} = \frac{SHRT\_MAX}{hz} - \frac{tv\_usec}{1,000,000}$$

The last term in this equation is not hz as specified in the code. The correct test is

if (tv->tv\_sec\*hz + tv->tv\_usec/tick > SHRT\_MAX)

but see Exercise 17.3 for more discussion.

<sup>831–840</sup> The converted time, val, is saved in the send or receive buffer as requested. sb\_timeo limits the amount of time a process will wait for data in the receive buffer or space in the send buffer. See Sections 16.7 and 16.12 for details.

The timeout values are passed as the last argument to tsleep, which expects an integer, so the process is limited to 65535 ticks. At 100 Hz, this less than 11 minutes.

getsockopt System Call 545

### Section 17.4

С

l. r

٦e

# 17.4 getsockopt System Call

getsockopt returns socket and protocol options as requested. The prototype for this system call is

int getsockopt(int s, int level, int name, caddr\_t val, int \*valsize);

The code is shown in Figure 17.10.

```
uipc_syscalls.c
598 struct getsockopt_args {
599
       int
                s;
600
                level;
        int
601
        int
                name:
       caddr_t val;
602
603
               *avalsize;
       int
604 };
605 getsockopt(p, uap, retval)
606 struct proc *p;
607 struct getsockopt_args *uap;
608 int
           *retval;
609 {
        struct file *fp;
610
611
        struct mbuf *m = NULL;
               valsize, error;
612
        int
        if (error = getsock(p->p_fd, uap->s, &fp))
613
           return (error);
614
615
        if (uap->val) {
            if (error = copyin((caddr_t) uap->avalsize, (caddr_t) & valsize,
615
617
                               sizeof(valsize)))
                return (error);
618
619
        } else
62.0
            valsize = 0;
        if ((error = sogetopt((struct socket *) fp->f_data, uap->level,
621
                   uap->name, &m)) == 0 && uap->val && valsize && m != NULL) {
622
            if (valsize > m->m_len)
623
624
                valsize = m->m_len;
            error = copyout(mtod(m, caddr_t), uap->val, (u_int) valsize);
625
            if (error == 0)
626
627
                error = copyout((caddr_t) & valsize,
                                 (caddr_t) uap->avalsize, sizeof(valsize));
628
629
        }
        if (m != NULL)
630
631
            (void) m_free(m);
632
        return (error);
633 }
                                                                       - uipc_syscalls.c
```

Figure 17.10 getsockopt system call.

<sup>598-633</sup> The code should look pretty familiar by now. getsock locates the socket, the size of the option buffer is copied into the kernel, and sogetopt is called to get the value of the requested option. The data returned by sogetopt is copied out to the buffer in the process along with the possibly new length of the buffer. It is possible that the data will

# INTEL Ex.1013.571

be silently truncated if the process did not provide a large enough buffer. As usual, the mbuf holding the option data is released before the function returns.

### sogetopt Function

As with sosetopt, the sogetopt function handles the socket-level options and passes any other options to the protocol associated with the socket. The beginning and end of the function are shown in Figure 17.11.

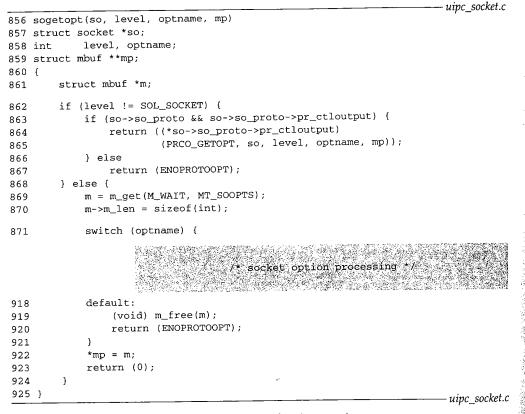


Figure 17.11 sogetopt function: overview.

As with sosetopt, options that do not pertain to the socket level are immediately passed to the protocol level through the PRCO\_GETOPT protocol request. The protocol returns the requested option in the mbuf pointed to by \*mp.

For socket-level options, a standard mbuf is allocated to hold the option value, which is normally an integer, so m\_len is set to the size of an integer. The appropriate option is copied into the mbuf by the code in the switch statement.

918-925 If the default case is taken by the switch, the mbuf is released and ENOPROTOOPT returned. Otherwise, after the switch statement, the pointer to the

547 getsockopt System Call

Section 17.4

In Figure 17.12 the linger option and the options that are implemented as boolean flags are processed. case SO\_LINGER: 872 m->m\_len = sizeof(struct linger); 873 mtod(m, struct linger \*)->l\_onoff = 874 so->so\_options & SO\_LINGER; 875 876 break; 877 case SO\_USELOOPBACK: 878 case SO\_DONTROUTE: 879 case SO\_DEBUG: 880 881 case SO\_KEEPALIVE: 882 case SO\_REUSEADDR: case SO\_REUSEPORT: 883 884 case SO\_BROADCAST: case SO\_OOBINLINE: 885 \*mtod(m, int \*) = so->so\_options & optname; 886 break; 887 Figure 17.12 sogetopt function: SO\_LINGER and boolean options. The SO\_LINGER option requires two copies, one for the flag into l\_onoff and a 872-877 second for the linger time into l\_linger. 878-887 off. Notice that the return value is not necessarily 1 when the flag is on. into the mbuf. case SO\_TYPE: 888 \*mtod(m, int \*) = so->so\_type; 889 break; \_socket.c 890 case SO\_ERROR: 891 \*mtod(m, int \*) = so->so\_error; 892 so->so\_error = 0; 893 diately 894 break; rotocol case SO\_SNDBUF: 895 \*mtod(m, int \*) = so->so\_snd.sb\_hiwat; value, 896 break; 897 opriate case SO\_RCVBUF: 898 \*mtod(m, int \*) = so->so\_rcv.sb\_hiwat; d and 899 to the 900 break:

mbuf is saved in \*mp. When this function returns, getsockopt copies the option from the mbuf to the process and releases the mbuf.

uipc\_socket.c mtod(m, struct linger \*)->l\_linger = so->so\_linger; - uipc\_socket.c

:er 17

l, the

asses

nd of

ocket.c

The remaining options are implemented as boolean flags. so\_options is masked with optname, which results in a nonzero value if the option is on and 0 if the option is

In the next part of sogetopt (Figure 17.13), the integer-valued options are copied

uipc\_socket.c

64) 5 (1

Ś

CAN'S ST

Socket Options 548

901 902 903	<pre>case S0_SNDLOWAT:  *mtod(m, int *) = so-&gt;so_snd.sb_lowat;  break;</pre>	
904 905 906	<pre>case SO_RCVLOWAT: *mtod(m, int *) = so-&gt;so_rcv.sb_lowat; break; uipc_socket.c</pre>	

Figure 17.13 sogetopt function: integer valued options.

Each option is copied as an integer into the mbuf. Notice that some of the options are stored as shorts in the kernel (e.g., the high-water and low-water marks) but 888-906 returned as integers. Also, so\_error is cleared once the value is copied into the mbuf. This is the only time that a call to getsockopt changes the state of the socket.

The fourth and last part of sogetopt is shown in Figure 17.14, where the SO\_SNDTIMEO and SO\_RCVTIMEO options are handled.

50_5MD	uipc_socket.c
907	case SO_SNDTIMEO:
908	case SO_RCVTIMEO:
909 910	{     int val = (optname == SO_SNDTIMEO ?
911	<pre>int val = (optimize = c; so-&gt;so_snd.sb_timeo : so-&gt;so_rcv.sb_timeo);</pre>
912	m->m_len = sizeof(struct timeval); mtod(m, struct timeval *)->tv_sec = val / hz;
913	mtod(m, struct timeval *)->tv_usec =
914	
915	(val % hz) / tick;
916	break;
917	} uipc_socket.c

Figure 17.14 sogetopt function: timeout options.

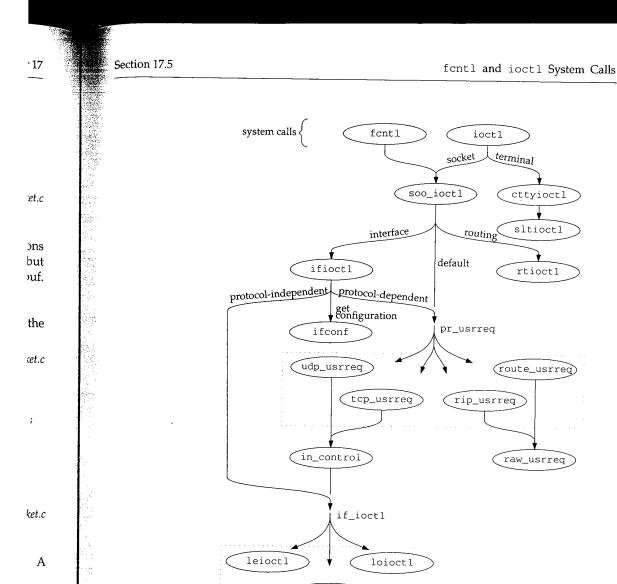
The sb\_timeo value from the send or receive buffer is copied into val. A 907–917 timeval structure is constructed in the mbuf based on the clock ticks in val.

There is a bug in the calculation of tv\_usec. The expression should be "(val % hz) \* tick".

#### fcntl and ioctl System Calls 17.5

Due more to history than intent, several features of the sockets API can be accessed from either ioctl or fcntl. We have already discussed many of the ioctl com-mands and have mentioned fcntl several times.

Figure 17.15 highlights the functions described in this chapter.



# Figure 17.15 fcntl and ioctl functions.

The prototypes for ioctl and fcntlare:

be

sed

om-

int ioctl(int fd, unsigned long result, char \*argp);

slioctl

int fcntl(int fd, int cmd, ... /\* int arg \*/);

Figure 17.16 summarizes the features of these two system calls as they relate to sockets. We show the traditional constants in Figure 17.16, since they appear in the code. For Posix compatibility, O\_NONBLOCK can be used instead of FNONBLOCK, and O\_ASYNC can be used instead of FASYNC.

549

i.

Description	fcntl	ioctl
enable or disable nonblocking semantics by	FNONBLOCK file status flag	FIONBIO command
turning SS_NBIO on or off in so_state enable or disable asynchronous notification	FASYNC file status flag	FIOASYNC command
by turning SB_ASYNC on or off in		
sb_flags set or get so_pgid, which is the target	F_SETOWN OF F_GETOWN	SIOCSPGRP or SIOCGPGRP
process or process group for SIGIO and SIGURG signals		commands
get number of bytes in receive buffer; return		FIONREAD
so_rcv.sb_cc return OOB synchronization mark; the SS_RCVATMARK flag in so_state		SIOCATMARK

Figure 17.16 fcntl and ioctl commands.

# fcntl Code

Figure 17.17 shows an overview of the fcntl function.

```
— kern_descrip.c
133 struct fcntl_args {
              fd;
134
       int
135
        int
                cmd;
        int
                arg;
136
137 };
138 /* ARGSUSED */
139 fcntl(p, uap, retval)
140 struct proc *p;
141 struct fcntl_args *uap;
           *retval;
142 int
143 {
        struct filedesc *fdp = p->p_fd;
144
        struct file *fp;
145
        struct vnode *vp;
146
              i, tmp, error, flg = F_POSIX;
147
        int
        struct flock fl;
148
        u_int newmin;
149
        if ((unsigned) uap->fd >= fdp->fd_nfiles ||
150
             (fp = fdp->fd_ofiles[uap->fd]) == NULL)
151
             return (EBADF);
 152
         switch (uap->cmd) {
 153
                               /* command processing */
         default:
 253
             return (EINVAL);
 254
 255
         }
         /* NOTREACHED */
 256
                                                                       - kern_descrip.c
 257 }
```

Figure 17.17 fcntl system call: overview.

r 17

·scrip.c

After verifying that the descriptor refers to an open file, the switch statement processes the requested command.

253-257 If the command is not recognized, fcntl returns EINVAL.

Figure 17.18 shows only the cases from fcntl that are relevant to sockets.

1979	6	kern_descrip.c
	168	case F_GETFL:
	169	<pre>*retval = OFLAGS(fp-&gt;f_flag);</pre>
19 	170	return (0);
	171	case F_SETFL:
	172	<pre>fp-&gt;f_flag &amp;= FCNTLFLAGS;</pre>
42.4 1.4 1.4 2.4 2.4	173	fp->f_flag  = FFLAGS(uap->arg) & FCNTLFLAGS;
	174	<pre>tmp = fp-&gt;f_flag &amp; FNONBLOCK; tmp = fp-&gt;f_flag &amp; FNONBLOCK;</pre>
	175	<pre>error = (*fp-&gt;f_ops-&gt;fo_ioctl) (fp, FIONBIO, (caddr_t) &amp; tmp, p);</pre>
and	176	if (error)
	177	return (error);
	178	<pre>tmp = fp-&gt;f_flag &amp; FASYNC;</pre>
25	179	error = (*fp->f_ops->fo_ioctl) (fp, FIOASYNC, (caddr_t) & tmp, p);
/ 8 · . 1997 :	180	if (!error)
	181	return (0);
iller Filter	182	fp->f_flag &= ~FNONBLOCK;
	183	tmp = 0;
	184	<pre>(mp = 0; (void) (*fp-&gt;f_ops-&gt;fo_ioctl) (fp, FIONBIO, (caddr_t) &amp; tmp, p);</pre>
20. 22 24 -	185	return (error);
	186	case F_GETOWN:
	187	if (fp->f_type == DTYPE_SOCKET) {
A.	188	<pre>*retval = ((struct socket *) fp-&gt;f_data)-&gt;so_pgid;</pre>
Ay .	189	return (0);
ng kanalan Na kanalan	190	}
	191	error = (*fp->f_ops->fo_ioctl)
e e de la se Esta de la secta de la sect Esta de la secta	192	(fp, (int) TIOCGPGRP, (caddr_t) retval, p);
1.500 C	193	<pre>*retval = -*retval;</pre>
	194	return (error);
an ta	195	case F_SETOWN:
	196	if (fp->f type == DTYPE_SOCKET) {
	197	((struct socket *) fp->f_data)->so_pgid = uap->arg;
	198	return (0);
	190	
1919 	200	if (uap->arg <= 0) {
	200	uap->arg = -uap->arg;
	201	} else {
	202	<pre>struct proc *p1 = pfind(uap-&gt;arg);</pre>
		if (p1 == 0)
d.	204	return (ESRCH);
े । सन्दर्भ	205	uap->arg = p1->p_pgrp->pg_id;
	206	
i segen daga Kabupatèn Kabupatèn Kabupatèn Kabupatèn	207	} return ((*fp->f_ops->fo_ioctl)
	208	(1, 1) $(1, 1)$ $(1, 2)$ $($
	209	(fp, (int) TIOCSPGRP, (Caddingt) a dap satisfy prover kern_descri



Figure 17.18 fcntl system call: socket processing.

- <sup>168–185</sup> F\_GETFL returns the current file status flags associated with the descriptor and F\_SETFL sets the flags. The new settings for FNONBLOCK and FASYNC are passed to the associated socket by calling fo\_ioctl, which for sockets is the soo\_ioctl function described with Figure 17.20. The third call to fo\_ioctl is made only if the second call fails. It clears the FNONBLOCK flag, but should instead restore the flag to its original setting.
- <sup>186-209</sup> F\_GETOWN returns so\_pgid, the process or process group associated with the socket. For a descriptor other than a socket, the TIOCGPGRP ioctl command is passed to the associated fo\_ioctl function. F\_SETOWN assigns a new value to so\_pgid.

For a descriptor other than a socket, the process group is checked in this function, but for sockets, the value is checked just before a signal is sent in sohasoutofband and in sowakeup.

### ioct1 Code

We skip the ioctl system call itself and start with soo\_ioctl in Figure 17.20, since most of the code in ioctl duplicates the code we described with Figure 17.17. We've already shown that this function sends routing commands to rtioctl, interface commands to ifioctl, and any remaining commands to the pr\_usrreq function of the underlying protocol.

- A few commands are handled by soo\_ioctl directly. FIONBIO turns on nonblocking semantics if \*data is nonzero, and turns them off otherwise. As we have seen, this flag affects the accept, connect, and close system calls as well as the various read and write system calls.
- 69-79 FIOASYNC enables or disables asynchronous I/O notification. Whenever there is activity on a socket, sowakeup gets called and if SS\_ASYNC is set, the SIGIO signal is sent to the process or process group.
- <sup>80–88</sup> FIONREAD returns the number of bytes available in the receive buffer. SIOCSPGRP sets the process group associated with the socket, and SIOCGPGRP gets it. so\_pgid is used as a target for the SIGIO signal as we just described and for the SIGURG signal when out-of-band data arrives for a socket. The signal is sent when the protocol layer calls the sohasoutofband function.
- <sup>89–92</sup> SIOCATMARK returns true if the socket is at the out-of-band synchronization mark, false otherwise.
  - ioctl commands, the FIOxxx and SIOxxx constants, have an internal structure illustrated in Figure 17.19.

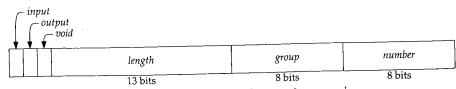


Figure 17.19 The structure of an ioctl command.

Section 17.5

17

١d

he

эn

all

et-

he

ed

on,

nd

nce

've

)m-

ion-

ave

·ari-

e is

al is

GRP

\_d is

gnal

ayer

hark,

cture

the  $\cdot$ 

sys\_socket.c 55 soo\_ioctl(fp, cmd, data, p) 56 struct file \*fp; 57 int cmd; 58 caddr\_t data; 59 struct proc \*p; 60 { struct socket \*so = (struct socket \*) fp->f\_data; 61 62 switch (cmd) { 63 case FIONBIO: 64 if (\*(int \*) data) 65 so->so\_state |= SS\_NBIO; 66 else so->so\_state &= ~SS\_NBIO; 67 68 return (0); case FIOASYNC: 69 if (\*(int \*) data) { 70 so->so\_state |= SS\_ASYNC; 71 so->so\_rcv.sb\_flags |= SB\_ASYNC; 72 73 so->so\_snd.sb\_flags |= SB\_ASYNC; 74} else { so->so\_state &= ~SS\_ASYNC; 75 so->so\_rcv.sb\_flags &= ~SB\_ASYNC; 76 so->so\_snd.sb\_flags &= ~SB\_ASYNC; 77 78 } return (0); 79 80 case FIONREAD: \*(int \*) data = so->so\_rcv.sb\_cc; 81 82 return (0); 83 case SIOCSPGRP: so\_pgid = \*(int \*) data; 84 85 return (0); 86 case SIOCGPGRP: \*(int \*) data = so->so\_pgid; 87 88 return (0); 89 case SIOCATMARK: \*(int \*) data = (so->so\_state & SS\_RCVATMARK) != 0; 90 return (0); 91 92 } 93 /\* \* Interface/routing/protocol specific ioctls: 94 95 \* interface and routing ioctls should have a \* different entry since a socket's unnecessary 96 97 \*/ if (IOCGROUP(cmd) == 'i') 98 return (ifioctl(so, cmd, data, p)); 99 if (IOCGROUP(cmd) == 'r') 100 return (rtioctl(cmd, data, p)); 101 return ((\*so->so\_proto->pr\_usrreq) (so, PRU\_CONTROL, 102 (struct mbuf \*) cmd, (struct mbuf \*) data, (struct mbuf \*) 0)); 103 104 } sys\_socket.c

Figure 17.20 soo\_ioctl function.

な理論を確認するになってい

If the third argument to ioctl is used as input, *input* is set. If the argument is used as output, *output* is set. If the argument is unused, *void* is set. *length* is the size of the argument in bytes. Related commands are in the same group but each command has its own *number* within the group. The macros in Figure 17.21 extract the components of an ioctl command.

Macro	Description
IOCPARM_LEN(cmd)	the <i>length</i> from <i>cmd</i>
IOCBASECMD(cmd)	the command with <i>length</i> set to 0
IOCGROUP(cmd)	the <i>group</i> from <i>cmd</i>

Figure 17.21	ioctl command macros.
--------------	-----------------------

93-104 The macro IOCGROUP extracts the 8-bit group from the command. Interface commands are handled by ifioctl. Routing commands are processed by rtioctl. All other commands are passed to the socket's protocol through the PRU\_CONTROL request.

As we describe in Chapter 19, Net/2 introduced a new interface to the routing tables in which messages are passed to the routing subsystem through a socket created in the PF\_ROUTE domain. This method replaces the ioctl method shown here. rtioctl always returns ENOTSUPP in kernels that do not have compatibility code compiled in.

# 17.6 getsockname System Call

The prototype for this system call is:

int getsockname(int fd, caddr\_t asa, int \*alen);

getsockname retrieves the local address bound to the socket *fd* and places it in the buffer pointed to by *asa*. This is useful when the kernel has selected an address during an implicit bind or when the process specified a wildcard address (Section 22.5) during an explicit call to bind. The getsockname system call is shown in Figure 17.22.

682-715

getsock locates the file structure for the descriptor. The size of the buffer specified by the process is copied from the process into len. This is the first call to m\_getclr that we've seen—it allocates a standard mbuf and clears it with bzero. The protocol processing layer is responsible for returning the local address in m when the PRU\_SOCKADDR request is issued.

If the address is larger than the buffer specified by the process, it is silently truncated when it is copied out to the process. \*alen is updated to the number of bytes copied out to the process. Finally, the mbuf is released and getsockname returns.

# 17.7 getpeername System Call

The prototype for this system call is:

int getpeername(int fd, caddr\_t asa, int \*alen);

Section 17.7

.7

٠d

۱e

ts

in

m-

All

est.

uich

JTE

irns

the

ring

ring

beci-

l to The

1 the

run-

ytes

uipc\_syscalls.c

682 struct getsockname\_args { 683 int fdes; 684 caddr\_t asa; 685 int \*alen; 686 }; 687 getsockname(p, uap, retval) 688 struct proc \*p; 689 struct getsockname\_args \*uap; \*retval; 690 int 691 { 692 struct file \*fp; struct socket \*so; 693 struct mbuf \*m; 694 695 int len, error; if (error = getsock(p->p\_fd, uap->fdes, &fp)) 696 697 return (error); if (error = copyin((caddr\_t) uap->alen, (caddr\_t) & len, sizeof(len))) 698 699 return (error); so = (struct socket \*) fp->f\_data; 700 m = m\_getclr(M\_WAIT, MT\_SONAME); 701 if (m == NULL) 702 return (ENOBUFS); 703 if (error = (\*so->so\_proto->pr\_usrreq) (so, PRU\_SOCKADDR, 0, m, 0)) 704 705 goto bad; 706 if (len > m->m\_len) len = m->m\_len; 707 error = copyout(mtod(m, caddr\_t), (caddr\_t) uap->asa, (u\_int) len); 708 if (error == 0) 709 error = copyout((caddr\_t) & len, (caddr\_t) uap->alen, 710 711 sizeof(len)); 712 bad: m\_freem(m); 713 714 return (error); 715 } uipc\_syscalls.c

Figure 17.22 getsockname system call.

The getpeername system call returns the address of the remote end of the connection associated with the specified socket. This function is often called when a server is invoked through a fork and exec by the process that calls accept (i.e., any server started by inetd). The server doesn't have access to the peer address returned by accept and must use getpeername. The returned address is often checked against an access list for the application, and the connection is closed if the address is not on the list.

Some protocols, such as TP4, utilize this function to determine if an incoming connection should be rejected or confirmed. In TP4, the connection associated with a socket returned by accept is not yet complete and must be confirmed before the connection completes. Based on the address returned by getpeername, the server can close the connection or implicitly confirm the connection by sending or receiving data. This 1

and South Strategies , 1 (2020) 23.

feature is irrelevant for TCP, since TCP doesn't make a connection available to accept until the three-way handshake is complete. Figure 17.23 shows the getpeername function.

uipc_syscalls.	c
719 struct getpeername_args {	Ű
720 int fdes;	
721 caddr_t asa;	
722 int *alen;	
723 };	
724 getpeername(p, uap, retval)	
725 struct proc *p;	
726 struct getpeername_args *uap;	
727 int *retval;	
728 (	
729 struct file *fp;	
730 struct socket *so;	
731 struct mbuf *m;	
732 int len, error;	
733 if (error = getsock( $p$ -> $p_fd$ , uap->fdes, &fp))	
734 return (error);	
<pre>735 so = (struct socket *) fp-&gt;f_data;</pre>	
<pre>if ((so-&gt;so_state &amp; (SS_ISCONNECTED   SS_ISCONFIRMING)) == 0)</pre>	
737 return (ENOTCONN);	
<pre>738 if (error = copyin((caddr_t) uap-&gt;alen, (caddr_t) &amp; len, sizeof(len)))</pre>	
739 return (error);	
<pre>740 m = m_getclr(M_WAIT, MT_SONAME);</pre>	
741 if $(m == NULL)$	
742 return (ENOBUFS);	
<pre>743 if (error = (*so-&gt;so_proto-&gt;pr_usrreq) (so, PRU_PEERADDR, 0, m, 0))</pre>	
744 goto bad;	
745 if (len > m->m_len)	
746 len = $m \rightarrow m_len;$	
<pre>747 if (error = copyout(mtod(m, caddr_t), (caddr_t) uap-&gt;asa, (u_int) len)) 748</pre>	
748 goto bad;	
<pre>749 error = copyout((caddr_t) &amp; len, (caddr_t) uap-&gt;alen, sizeof(len));</pre>	
750 bad:	
751 m_freem(m);	
752 return (error); 753 }	
uipc syscalls.c	

Figure 17.23 getpeername system call.

719-753

The code here is almost identical to the getsockname code. getsock locates the socket and ENOTCONN is returned if the socket is not yet connected to a peer or if the connection is not in a confirmation state (e.g., TP4). If it is connected, the size of the buffer is copied in from the process and an mbuf is allocated to hold the address. The PRU\_PEERADDR request is issued to get the remote address from the protocol layer. The address and the length of the address are copied from the kernel mbuf to the buffer in the process. The mbuf is released and the function returns.

С

i.C

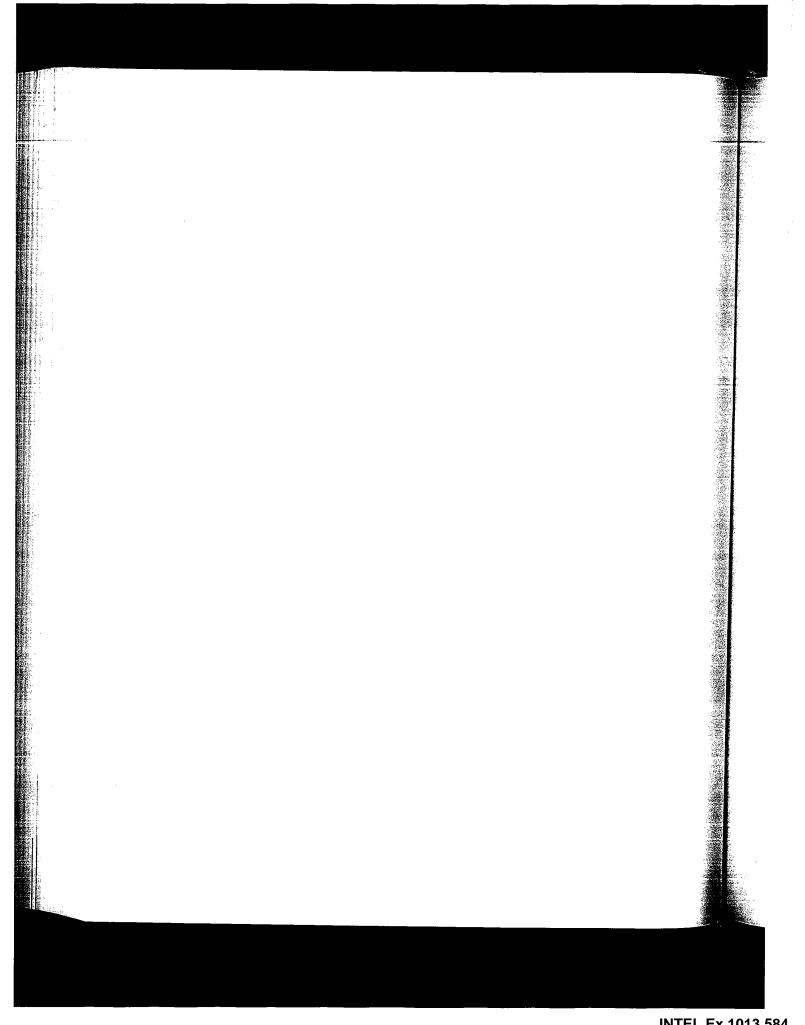
າຍ າຍ າຍ າຍ າຍ

# 17.8 Summary

In this chapter we discussed the six functions that modify the semantics of a socket. Socket options are processed by setsockopt and getsockopt. Additional options, some of which are not unique to sockets, are handled by fcntl and ioctl. Finally, connection information is available through getsockname and getpeername.

# **Exercises**

- 17.1 Why do you think options are limited to the size of a standard mbuf (MHLEN, 128 bytes)?
- 17.2 Why does the code at the end of Figure 17.7 work for the SO\_LINGER option?
- 17.3 There is a problem with the suggested code used to test the timeval structure in Figure 17.9 since tv->tv\_sec \* hz may cause an overflow. Suggest a change to the code to solve this problem.



INTEL Ex.1013.584

# 18

# Radix Tree Routing Tables

# 18.1 Introduction

The routing performed by IP, when it searches the routing table and decides which interface to send a packet out on, is a *routing mechanism*. This differs from a *routing policy*, which is a set of rules that decides which routes go into the routing table. The Net/3 kernel implements the routing mechanism while a routing daemon, typically routed or gated, implements the routing policy. The structure of the routing table must recognize that the packet forwarding occurs frequently—hundreds or thousands of times a second on a busy system—while routing policy changes are less frequent.

Routing is a detailed issue and we divide our discussion into three chapters.

- This chapter looks at the structure of the radix tree routing tables used by the Net/3 packet forwarding code. The tables are consulted by IP every time a packet is sent (since IP must determine which local interface receives the packet) and every time a packet is forwarded.
- Chapter 19 looks at the functions that interface between the kernel and the radix tree functions, and also at the routing messages that are exchanged between the kernel and routing processes—normally the routing daemons that implement the routing policy. These messages allow a process to modify the kernel's routing table (add a route, delete a route, etc.) and let the kernel notify the daemons when an asynchronous event occurs that might affect the routing policy (a redirect is received, an interface goes down, and so on).
- Chapter 20 presents the routing sockets that are used to exchange routing messages between the kernel and a process.

559

# 18.2 Routing Table Structure

Before looking at the internal structure of the Net/3 routing table, we need to understand the type of information contained in the table. Figure 18.1 is the bottom half of Figure 1.17: the four systems on the author's Ethernet.

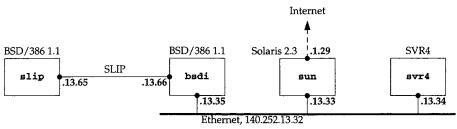


Figure 18.1 Subnet used for routing table example.

Figure 18.2 shows the routing table for bsdi in Figure 18.1.

```
bsdi $ netstat -rn
Routing tables
```

Internet:					
Destination	Gateway	Flags	Refs	Use	Interface
default	140.252.13.33	UG S	0	3	le0
127	127.0.0.1	UG S R	0	2	100
127.0.0.1	127.0.0.1	υн	1	55	100
128.32.33.5	140.252.13.33	UGHS	2	16	le0
140.252.13.32	link#1	υC	0	0	le0
140.252.13.33	8:0:20:3:f6:42	UHL	11	55146	1e0
140.252.13.34	0:0:c0:c2:9b:26	UHL	0	3	1e0
140.252.13.35	0:0:c0:6f:2d:40	UH L	1	12	100
140.252.13.65	140.252.13.66	υн	0	41	s10
224	link#1	UC	0	0	le0
224.0.0.1	link#1	UHL	0	5	le0

Figure 18.2 Routing table on the host bsdi.

We have modified the "Flags" column from the normal netstat output, making it easier to see which flags are set for the various entries.

The routes in this table were entered as follows. Steps 1, 3, 5, 8, and 9 are performed at system initialization when the /etc/netstart shell script is executed.

- 1. A default route is added by the route command to the host sun (140.252.13.33), which contains a PPP link to the Internet.
- 2. The entry for network 127 is typically created by a routing daemon such as gated, or it can be entered with the route command in the /etc/netstart file. This entry causes all packets sent to this network, other than references to the host 127.0.0.1 (which are covered by the more specific route entered in the next step), to be rejected by the loopback driver (Figure 5.27).

Chapter 18

а. "A

「「ないない」

「「「「「「「」」」」

- 3. The entry for the loopback interface (127.0.0.1) is configured by ifconfig.
- 4. The entry for vangogh.cs.berkeley.edu (128.32.33.5) was created by hand using the route command. It specifies the same router as the default route (140.252.13.33), but having a host-specific route, instead of using the default route for this host, allows routing metrics to be stored in this entry. These metrics can optionally be set by the administrator, are used by TCP each time a connection is established to the destination host, and are updated by TCP when the connection is closed. We describe these metrics in more detail with Figure 27.3.
- 5. The interface le0 is initialized using the ifconfig command. This causes the entry for network 140.252.13.32 to be entered into the routing table.
- 6. The entries for the other two hosts on the Ethernet, sun (140.252.13.33) and svr4 (140.252.13.34), were created by ARP, as we describe in Chapter 21. These are temporary entries that are removed if they are not used for a certain period of time.
- 7. The entry for the local host, 140.252.13.35, is created the first time the host's own IP address is referenced. The interface is the loopback, meaning any IP datagrams sent to the host's own IP address are looped back internally. The automatic creation of this entry is new with 4.4BSD, as we describe in Section 21.13.
- 8. The entry for the host 140.252.13.65 is created when the SLIP interface is configured by ifconfig.
- 9. The route command adds the route to network 224 through the Ethernet interface.
- 10. The entry for the multicast group 224.0.0.1 (the all-hosts group) was created by running the Ping program, pinging the address 224.0.0.1. This is also a temporary entry that is removed if not used for a certain period of time.

The "Flags" column in Figure 18.2 needs a brief explanation. Figure 18.25 provides a list of all the possible flags.

- U The route is up.
- G The route is to a gateway (router). This is called an *indirect route*. If this flag is not set, the destination is directly connected; this is called a *direct route*.
- H The route is to a host, that is, the destination is a complete host address. If this flag is *not* set, the route is to a network, and the destination is a network address: a network ID, or a combination of a network ID and a subnet ID. The netstat command doesn't show it, but each network route also contains a network mask. A host route has an implied mask of all one bits.
- S The route is static. The three entries created by the route command in Figure 18.2 are static.

underhalf of

g it eas-

rformed

ost sun

such as tstart ences to ed in the

a,

- C The route is cloned to create new routes. Two entries in this routing table have this flag set: (1) the route for the local Ethernet (140.252.13.32), which is cloned by ARP to create the host-specific routes of other hosts on the Ethernet, and (2) the route for multicast groups (224), which is cloned to create specific multicast group routes such as 224.0.0.1
- L The route contains a link-layer address. The host routes that ARP clones from the Ethernet network routes all have the link flag set. This applies to unicast and multicast addresses.
- R The loopback driver (the normal interface for routes with this flag) rejects all datagrams that use this route.

The ability to enter a route with the "reject" flag was provided in Net/2. It provides a simple way of preventing datagrams destined to network 127 from appearing outside the host. See also Exercise 6.6.

Before 4.3BSD Reno, two distinct routing tables were maintained by the kernel for IP addresses: one for host routes and one for network routes. A given route was entered into one table or the other, based on the type of route. The default route was stored in the network routing table with a destination address of 0.0.0.0. There was an implied hierarchy: a search was made for a host route first, and if not found a search was made for a network route, and if still not found, a search was made for a default route. Only if all three searches failed was the destination unreachable. Section 11.5 of [Leffler et al. 1989] describes the hash table with linked lists used for the host and network routing tables in Net/1.

Major changes took place in the internal representation of the routing table with 4.3BSD Reno [Sklower 1991]. These changes allow the same routing table functions to access a routing table for other protocol suites, notably the OSI protocols, which use variable-length addresses, unlike the fixed-length 32-bit Internet addresses. The internal structure was also changed, to provide faster lookups.

The Net/3 routing table uses a Patricia tree structure [Sedgewick 1990] to represent both host addresses and network addresses. (Patricia stands for "Practical Algorithm to Retrieve Information Coded in Alphanumeric.") The address being searched for and the addresses in the tree are considered as sequences of bits. This allows the same functions to maintain and search one tree containing fixed-length 32-bit Internet addresses, another tree containing fixed-length 48-bit XNS addresses, and another tree containing variable-length OSI addresses.

The idea of using Patricia trees for the routing table is attributed to Van Jacobson in [Sklower 1991]. These are actually binary radix tries with one-way branching removed.

An example is the easiest way to describe the algorithm. The goal of routing lookup is to find the most specific address that matches the given destination: the search key. The term *most specific* implies that a host address is preferred over a network address, which is preferred over a default address.

Each entry has an associated network mask, although no mask is stored with a host route; instead host routes have an implied mask of all one bits. An entry in the routing table matches a search key if the search key logically ANDed with the network mask of

Routing Table Structure 563

Section 18.2

the entry equals the entry itself. A given search key might match multiple entries in the routing table, so with a single table for both network route and host routes, the table must be organized so that more-specific routes are considered before less-specific routes.

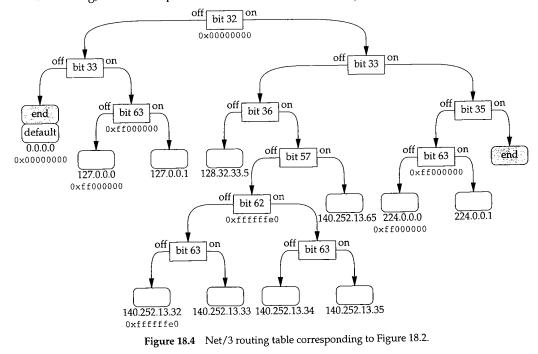
Consider the examples in Figure 18.3. The two search keys are 127.0.0.1 and 127.0.0.2, which we show in hexadecimal since the logical ANDing is easier to illustrate. The two routing table entries are the host entry for 127.0.0.1 (with an implied mask of  $0 \times ffffffff$ ) and the network entry for 127.0.0.0 (with a mask of  $0 \times ff000000$ ).

		search key = 127.0.0.1		search key = 127.0.0.2	
		host route	net route	host route	net route
1	search key	7f000001	7£000001	7£000002	7£000002
2	routing table key	7£000001	7£000000	7£000001	7£000000
3	routing table mask	fffffff	ff000000	fffffff	ff000000
4	logical AND of 1 and 3	7£000001	7£000000	7f000002	7£000000
	2 and 4 equal?	yes	yes	no	yes

Figure 18.3 Example routing table lookups for the two search keys 127.0.0.1 and 127.0.0.2.

Since the search key 127.0.0.1 matches both routing table entries, the routing table must be organized so that the more-specific entry (127.0.0.1) is tried first.

Figure 18.4 shows the internal representation of the Net/3 routing table corresponding to Figure 18.2. This table was built from the output of the netstat command with the -A flag, which dumps the tree structure of the routing tables.



er 18

lave

med

1(2)

icast

um to

1.5 of

l net-

with

ns to

r and funcesses, ining

klower

okup 1 key.

dress, a host



The two shaded boxes labeled "end" are leaves with special flags denoting the end of the tree. The left one has a key of all zero bits and the right one has a key of all one bits. The two boxes stacked together at the left, labeled "end" and "default," are a special representation used for duplicate keys, which we describe in Section 18.9.

The square-cornered boxes are called *internal nodes* or just *nodes*, and the boxes with rounded corners are called *leaves*. Each internal node corresponds to a bit to test in the search key, and a branch is made to the left or the right. Each leaf corresponds to either a host address or a network address. If there is a hexadecimal number beneath a leaf, that leaf is a network address and the number specifies the network mask for the leaf. The absence of a hexadecimal mask beneath a leaf node implies that the leaf is a host address with an implied mask of 0xfffffff.

Some of the internal nodes also contain network masks, and we'll see how these are used in backtracking. Not shown in this figure is that every node also contains a pointer to its parent, to facilitate backtracking, deletion, and nonrecursive walks of the tree.

The bit comparisons are performed on socket address structures, so the bit positions given in Figure 18.4 are from the start of the socket address structure. Figure 18.5 shows the bit positions for a sockaddr\_in structure.

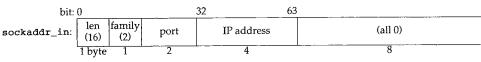


Figure 18.5 Bit offsets in Internet socket address structure.

The highest-order bit of the IP address is at bit position 32 and the lowest-order bit is at bit position 63. We also show the length as 16 and the address family as 2 (AF\_INET), as we'll encounter these two values throughout our examples.

To work through the examples we also need to show the bit representations of the various IP addresses in the tree. These are shown in Figure 18.6 along with some other IP addresses that are used in the examples that follow. The bit positions used in Figure 18.4 as branching points are shown in a bolder font.

We now provide some specific examples of how the routing table searches are performed.

## Example—Host Match

Assume the host address 127.0.0.1 is the search key—the destination address being looked up. Bit 32 is off, so the left branch is made from the top of the tree. Bit 33 is on, so the right branch is made from the next node. Bit 63 is on, so the right branch is made from the next node. This next node is a leaf, so the search key (127.0.0.1) is compared to the address in the leaf (127.0.0.1). They match exactly so this routing table entry is returned by the lookup function.

dotted-decimal 32-bit IP address (bits 32-63) **5**55 **66** 333 bit: **7**89 789 10.1.2.3 112.0.0.1 127.0.0.0 127.0.0.1 127.0.0.3 128.32.33.5 128.32.33.6 140.252.13.32 140.252.13.33 140.252.13.34 140.252.13.35 140.252.13.65 224.0.0.0 224.0.0.1 

Figure 18.6 Bit representations of the IP addresses in Figures 18.2 and 18.4.

#### Example—Host Match

Next assume the search key is the address 140.252.13.35. Bit 32 is on, so the right branch is made from the top of the tree. Bit 33 is off, bit 36 is on, bit 57 is off, bit 62 is on, and bit 63 is on, so the search ends at the leaf on the bottom labeled 140.252.13.35. The search key matches the routing table key exactly.

# Example—Network Match

The search key is 127.0.0.2. Bit 32 is off, bit 33 is on, and bit 63 is off so the search ends up at the leaf labeled 127.0.0.0. The search key and the routing table key don't match exactly, so a network match is tried. The search key is logically ANDed with the network mask (0xff000000) and since the result equals the routing table key, this entry is considered a match.

#### Example—Default Match

The search key is 10.1.2.3. Bit 32 is off and bit 33 is off, so the search ends up at the leaf with the duplicate keys labeled "end" and "default." The routing table key that is duplicated in these two leaves is 0.0.0. The search key and the routing table key don't match exactly, so a network match is tried. This match is tried for all duplicate keys that have a network mask. The first key (the end marker) doesn't have a network mask, so it is skipped. The next key (the default entry) has a mask of 0x00000000. The search key is logically ANDed with this mask and since the result equals the routing table key (0), this entry is considered a match. The default route is used.

ows

r 18

end

one

pe-

vith

the

:her

eaf,

.eaf.

lost

: are

is at IET),

f the other Fig-

per-

being s on, nade ed to ry is

# Example—Network Match with Backtracking

The search key is 127.0.0.3. Bit 32 is off, bit 33 is on, and bit 63 is on, so the search ends up at the leaf labeled 127.0.0.1. The search key and the routing table key don't match exactly. A network match cannot be attempted since this leaf does not have a network mask. Backtracking now takes place.

The backtracking algorithm is to move up the tree, one level at a time. If an internal node is encountered that contains a mask, the search key is logically ANDed with the mask and another search is made of the subtree starting at the node with the mask, looking for a match with the ANDed key. If a match isn't found, the backtrack keeps moving up the tree, until the top is reached.

In this example the search moves up one level to the node for bit 63 and this node contains a mask. The search key is logically ANDed with the mask (0xff000000), giving a new search key of 127.0.0.0. Another search is made starting at this node for 127.0.0.0. Bit 63 is off, so the left branch is taken to the leaf labeled 127.0.0.0. The new search key is compared to the routing table key and since they're equal, this leaf is the match.

# Example—Backtracking Multiple Levels

The search key is 112.0.0.1. Bit 32 is off, bit 33 is on, and bit 63 is on, so the search ends up at the leaf labeled 127.0.0.1. The keys are not equal and the routing table entry does not have a network mask, so backtracking takes place

The search moves up one level to the node for bit 63, which contains a mask. The search key is logically ANDed with the mask of  $0 \times ff000000$  and another search is made starting at that node. Bit 63 is off in the new search key, so the left branch is made to the leaf labeled 127.0.0.0. A comparison is made but the ANDed search key (112.0.0.0) doesn't equal the search key in the table.

Backtracking continues up one level from the bit-63 node to the bit-33 node. But this node does not have a mask, so the backtracking continues upward. The next level is the top of the tree (bit 32) and it has a mask. The search key (112.0.0.1) is logically ANDed with the mask ( $0 \times 00000000$ ) and a new search started from that point. Bit 32 is off in the new search key, as is bit 33, so the search ends up at the leaf labeled "end" and "default." The list of duplicate keys is traversed and the default key matches the new search key, so the default route is used.

As we can see in this example, if a default route is present in the routing table, when the backtrack ends up at the top node in the tree, its mask is all zero bits, which causes the search to proceed to the leftmost leaf in the tree for a match with the default.

# Example—Host Match with Backtracking and Cloning

The search key is 224.0.0.5. Bit 32 is on, bit 33 is on, bit 35 is off, and bit 63 is on, so the search ends up at the leaf labeled 224.0.0.1. This routing table key does not equal the search key, and the routing table entry does not contain a network mask, so backtracking takes place.

The backtrack moves one level up to the node that tests bit 63. This node contains the mask 0xff000000, so the search key ANDed with the mask yields a new search key of 224.0.0.0. Another search is made, starting at this node. Since bit 63 is off in the ANDed key, the left branch is taken to the leaf labeled 224.0.0.0. This routing table key matches the ANDed search key, so this entry is a match.

This route has the "clone" flag set (Figure 18.2), so a new leaf is created for the address 224.0.0.5. The new routing table entry is

Destination	Gateway	Flags	Refs	Use	Interface
224.0.0.5	link#1	UHL	0	0	1e0

and Figure 18.7 shows the new arrangement of the right side of the routing table tree from Figure 18.4, starting with the node for bit 35. Notice that whenever a new leaf is added to the tree, two nodes are needed: one for the leaf and one for the internal node specifying the bit to test.

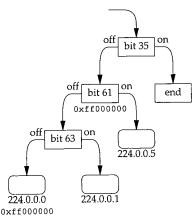


Figure 18.7 Modification of Figure 18.4 after inserting entry for 224.0.0.5.

This newly created entry is the one returned to the caller who was searching for 224.0.0.5.

### The Big Picture

Figure 18.8 shows a bigger picture of all the data structures involved. The bottom portion of this figure is from Figure 3.32.

There are numerous points about this figure that we'll note now and describe in detail later in this chapter.

 rt\_tables is an array of pointers to radix\_node\_head structures. There is one entry in the array for each address family. rt\_tables[AF\_INET] points to the top of the Internet routing table tree.

de vor w he ds es he is de

Suf

vel llv

32

d″ :he

ien

ses

18

ds

ch

rk

ıal he ⊧k, ps

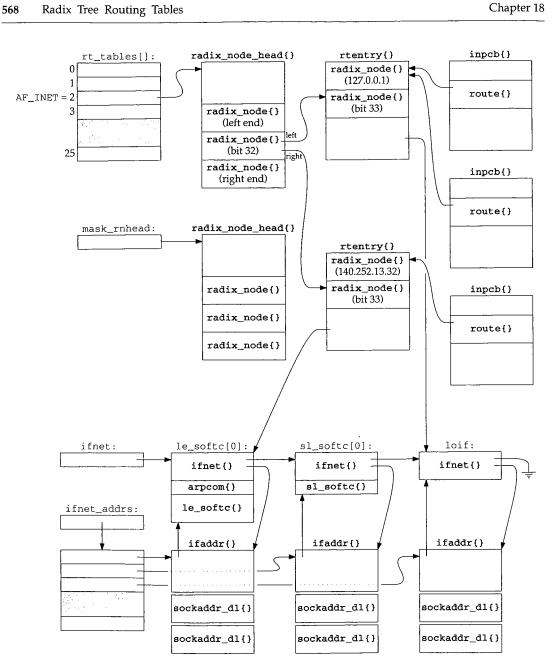


Figure 18.8 Data structures involved with routing tables.

X

Routing Sockets 569

Section 18.3

- The radix\_node\_head structure contains three radix\_node structures. These structures are built when the tree is initialized and the middle of the three is the top of the tree. This corresponds to the top box in Figure 18.4, labeled "bit 32." The first of the three radix\_node structures is the leftmost leaf in Figure 18.4 (the shared duplicate with the default route) and the third of the three is the rightmost leaf. An empty routing table consists of just these three radix\_node structures; we'll see how it is constructed by the rn\_inithead function.
- The global mask\_rnhead also points to a radix\_node\_head structure. This is the head of a separate tree of all the masks. Notice in Figure 18.4 that of the eight masks shown, one is duplicated four times and two are duplicated once. By keeping a separate tree for the masks, only one copy of each unique mask is maintained.
- The routing table tree is built from rtentry structures, and we show two of these in Figure 18.8. Each rtentry structure contains two radix\_node structures, because each time a new entry is inserted into the tree, two nodes are required: an internal node corresponding to a bit to be tested, and a leaf node corresponding to a host route or a network route. In each rtentry structure we also show which bit test the internal node corresponds to and the address contained in the leaf node.

The remainder of the rtentry structure is the focal point of information for this route. We show only a single pointer from this structure to the corresponding ifnet structure for the route, but this structure also contains a pointer to the ifaddr structure, the flags for the route, a pointer to another rtentry structure if this entry is an indirect route, the metrics for the route, and so on.

• Protocol control blocks (Chapter 22), of which one exists for each UDP and TCP socket (Figure 22.1), contain a route structure that points to an rtentry structure. The UDP and TCP output functions both pass a pointer to the route structure in a PCB as the third argument to ip\_output, each time an IP datagram is sent. PCBs that use the same route point to the same routing table entry.

## 18.3 Routing Sockets

When the routing table changes were made with 4.3BSD Reno, the interaction of processes with the routing subsystem also changed—the concept of routing sockets was introduced. Prior to 4.3BSD Reno, fixed-length ioctls were issued by a process (such as the route command) to modify the routing table. 4.3BSD Reno changed this to a more generalized message-passing scheme using the new PF\_ROUTE domain. A process creates a raw socket in the PF\_ROUTE domain and can send routing messages to the kernel, and receives routing messages from the kernel (e.g., redirects and other asynchronous notifications from the kernel).

Figure 18.9 shows the 12 different types of routing messages. The message type is the rtm\_type field in the rt\_msghdr structure, which we describe in Figure 19.16. Only five of the messages can be issued by a process (a write to a routing socket), but all 12 can be received by a process.

We'll defer our discussion of these routing messages until Chapter 19.

Chapter 18

rtm_type	To kernel?	From kernel?	Description	Structure type
RTM ADD	•	•	add route	rt_msghdr
RTM_HDD RTM_CHANGE	•	•	change gateway, metrics, or flags	rt_msghdr
RTM_DELADDR		•	address being removed from interface	ifa_msghdr
RTM DELETE	•	•	delete route	rt_msghdr
RTM GET	•	•	report metrics and other route information	rt_msghdr
RTM IFINFO		•	interface going up, down, etc.	if_msghdr
RTM_LOCK	•	•	lock specified metrics	rt_msghdr
RTM_LOSING		•	kernel suspects route is failing	rt_msghdr
RTM MISS		•	lookup failed on this address	rt_msghdr
RTM_NEWADDR		•	address being added to interface	ifa_msghdr
RTM_REDIRECT		•	kernel told to use different route	rt_msghdr
RTM_RESOLVE		•	request to resolve destination to link-layer address	rt_msghdr

Figure 18.9	Types of messages exchanged	l across a routing socket.
-------------	-----------------------------	----------------------------

## 18.4 Code Introduction

Three headers and five C files define the various structures and functions used for routing. These are summarized in Figure 18.10.

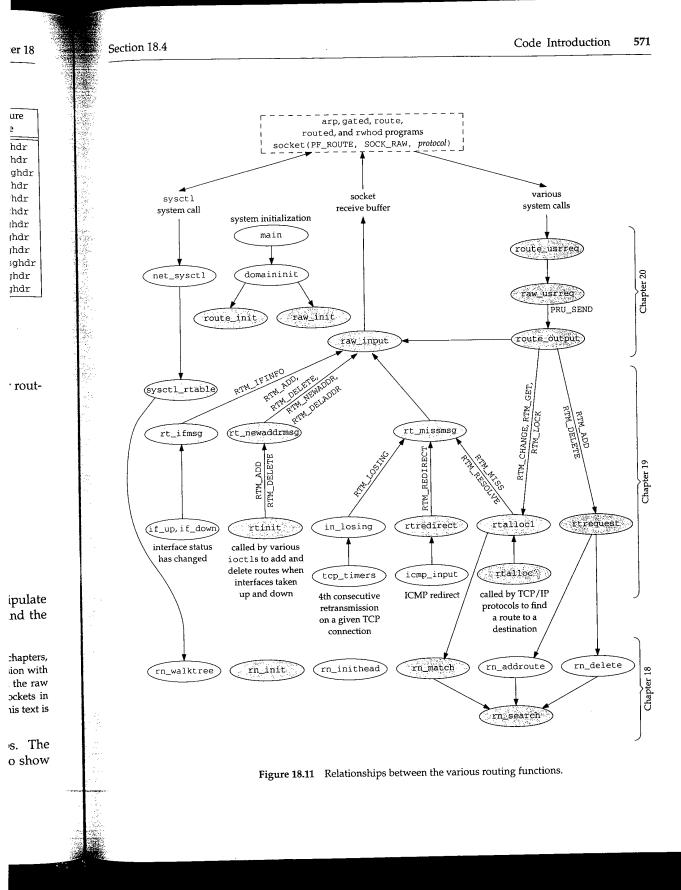
File	Description
net/radix.h net/raw_cb.h net/route.h	radix node definitions routing control block definitions routing structures
<pre>net/radix.c net/raw_cb.c net/raw_usrreq.c net/route.c net/rtsock.c</pre>	radix node (Patricia tree) functions routing control block functions routing control block functions routing functions routing socket functions

Figure 18.10 Files discussed in this chapter.

In general, the prefix rn\_ denotes the radix node functions that search and manipulate the Patricia trees, the raw\_ prefix denotes the routing control block functions, and the three prefixes route\_, rt\_, and rt denote the general routing functions.

We use the term *routing control blocks* instead of *raw control blocks* in all the routing chapters, even though the files and functions begin with the prefix raw. This is to avoid confusion with the raw IP control blocks and functions, which we discuss in Chapter 32. Although the raw control blocks and their associated functions are used for more than just routing sockets in Net/3 (one of the raw OSI protocols uses these structures and functions), our use in this text is only with routing sockets in the PF\_ROUTE domain.

Figure 18.11 shows the primary routing functions and their relationships. The shaded ellipses are the ones we cover in this chapter and the next two. We also show where each of the 12 routing message types are generated.



**INTEL Ex.1013.597** 

rtalloc is the function called by the Internet protocols to look up routes to destinations. We've already encountered rtalloc in the ip\_rtaddr, ip\_forward, ip\_output, and ip\_setmoptions functions. We'll also encounter it later in the in\_pcbconnect and tcp\_mss functions.

We also show in Figure 18.11 that five programs typically create sockets in the routing domain:

- arp manipulates the ARP cache, which is stored in the IP routing table in Net/3 (Chapter 21),
- gated and routed are routing daemons that communicate with other routers and manipulate the kernel's routing table as the routing environment changes (routers and links go up or down),
- route is a program typically executed by start-up scripts or by the system administrator to add or delete routes, and
- rwhod issues a routing sysctl on start-up to determine the attached interfaces.

Naturally, any process (with superuser privilege) can open a routing socket to send and receive messages to and from the routing subsystem; we show only the common system programs in Figure 18.11.

# **Global Variables**

The global variables introduced in the three routing chapters are shown in Figure 18.12.

Variable	Datatype	Description
rt_tables mask_rnhead rn_mkfreelist	struct radix_node_head * [] struct radix_node_head * struct radix_mask *	array of pointers to heads of routing tables pointer to head of mask table head of linked list of available radix_mask structures
max_keylen rn_zeros rn_ones maskedKey	int char * char * char *	longest routing table key, in bytes array of all zero bits, of length max_keylen array of all one bits, of length max_keylen array for masked search key, of length max_keylen
rtstat	struct rtstat	routing statistics (Figure 18.13) #routes not in table but not freed
rttrash rawcb raw_recvspace raw_sendspace	struct rawcb u_long u long	head of doubly linked list of routing control blocks default size of routing socket receive buffer, 8192 bytes default size of routing socket send buffer, 8192 bytes
raw_sendspace route_cb route_dst route_src route_proto	struct route_cb struct sockaddr struct sockaddr struct sockaddr	#routing socket listeners, per protocol, and total temporary for destination of routing message temporary for source of routing message temporary for protocol of routing message

Figure 18.12 Global variables in the three routing chapters.

## Statistics

Some routing statistics are maintained in the global structure rtstat, described in Figure 18.13.

rtstat member	Description	Used by SNMP
rts_badredirect rts_dynamic rts_newgateway rts_unreach rts_wildcard	<ul> <li>#invalid redirect calls</li> <li>#routes created by redirects</li> <li>#routes modified by redirects</li> <li>#lookups that failed</li> <li>#lookups matched by wildcard (never used)</li> </ul>	

Figure 18.13 Routing statistics maintained in the rtstat structure.

We'll see where these counters are incremented as we proceed through the code. None are used by SNMP.

Figure 18.14 shows some sample output of these statistics from the netstat -rs command, which displays this structure.

netstat -rs output	rtstat member
1029 bad routing redirects	rts_badredirect
0 dynamically created routes	rts_dynamic
0 new gateways due to redirects	rts_newgateway
0 destinations found unreachable	rts_unreach
0 uses of a wildcard route	rts_wildcard

Figure 18.14 Sample routing statistics.

## **SNMP Variables**

Figure 18.15 shows the IP routing table, named ipRouteTable, and the kernel variables that supply the corresponding value.

For ipRouteType, if the RTF\_GATEWAY flag is set in rt\_flags, the route is remote (4); otherwise the route is direct (3). For ipRouteProto, if either the RTF\_DYNAMIC or RTF\_MODIFIED flag is set, the route was created or modified by ICMP (4), otherwise the value is other (1). Finally, if the rt\_mask pointer is null, the returned mask is all one bits (i.e., a host route).

## 18.5 Radix Node Data Structures

In Figure 18.8 we see that the head of each routing table is a radix\_node\_head and all the nodes in the routing tree, both the internal nodes and the leaves, are radix\_node structures. The radix\_node\_head structure is shown in Figure 18.16.

Chapter 18

IP routing table, index = < ipRouteDest >				
SNMP variable	Variable	Description		
ipRouteDest	rt_key	Destination IP address. A value of 0.0.0.0 indicates a default entry.		
ipRouteIfIndex ipRouteMetric1	rt_ifp.if_index -1	Interface number: ifIndex. Primary routing metric. The meaning of the metric depends on the routing protocol (ipRouteProto). A value of -1 means it is not used.		
ipRouteMetric2 ipRouteMetric3 ipRouteMetric4 ipRouteNextHop ipRouteType	-1 -1 rt_gateway (see text)	Alternative routing metric. Alternative routing metric. Alternative routing metric. IP address of next-hop router. Route type: 1 = other, 2 = invalidated route, 3 = direct, 4 = indirect.		
ipRouteProto	(see text)	Routing protocol: 1 = other, 4 = ICMP redirect, 8 = RIP, 13 = OSPF, 14 = BGP, and others.		
ipRouteAge	(not implemented)	Number of seconds since route was last updated or determined to be correct.		
ipRouteMask	rt_mask	Mask to be logically ANDed with destination IP address before being compared with ipRouteDest.		
ipRouteMetric5 ipRouteInfo	-1 NULL	Alternative routing metric. Reference to MIB definitions specific to this particular routing protocol.		

Figure 18.15 IP routing table: ipRouteTable.

		radix.h
91 stru	uct radi	x_node_head (
92	atruct	radix node *rnh treetop;
93	int	<pre>rnh_addrsize; /* (not currently used) */</pre>
94	int	<pre>rnh_pktsize; /* (not currently used) */ rnh_pktsize; /* (not currently used) */</pre>
95	struct	radix_node *(*rnh_addaddr) /* add based on sockaddr */
96		(void *v, void *mask,
97		struct radix_node_head * head, struct radix_node nodes[]);
98	struct	radix_node *(*rnh_addpkt) /* add based on packet hdr */
99		(void *v, void *mask,
100		struct radix_node_head * head, struct radix_node nodes[]);
101	struct	radix_node *(*rnh_deladdr) /* remove based on sockaddr */
102		(void *v, void *mask, struct radix_node_head * head);
103	struct	(Void *V, Void mask, bitset /* remove based on packet hdr */ radix_node *(*rnh_delpkt) /* remove based on packet hdr */
104		<pre>radix_node ('fim_dorphi's truct radix_node_head * head); (void *v, void *mask, struct radix_node_head * head); radix_node *(*rnh_matchaddr) /* locate based on sockaddr */</pre>
105	struct	radix_node *(*rnh_matchadur) / focute subset
106		<pre>(void *v, struct radix_node_head * head); radix_node *(*rnh_matchpkt) /* locate based on packet hdr */</pre>
107	struct	radix_node *(*rnh_matchpkt) /* Tocate back an part
108		<pre>(void *v, struct radix_node_head * head);   (void *v, struct radix_node_head * head);</pre>
109	int	(*rnh_walktree) /* traverse tree */ (struct radix_node_head * head, int (*f) (), void *w);
110		(struct radix_node_nead ~ nead, int ( 1, (), topic a,
111	atruct	radix_node rnh_nodes[3]; /* top and end nodes */
111 112 };	struct	radix_noue inn_nous(i), radix.

Figure 18.16 radix\_node\_head structure: the top of each routing tree.

92

ı

rnh\_treetop points to the top radix\_node structure for the routing tree. Notice that three of these structures are allocated at the end of the radix\_node\_head, and the middle one of these is initialized as the top of the tree (Figure 18.8).

93-94 rnh\_addrsize and rnh\_pktsize are not currently used.

rnh\_addrsize is to facilitate porting the routing table code to systems that don't have a length byte in the socket address structure. rnh\_pktsize is to allow using the radix node machinery to examine addresses in packet headers without having to copy the address into a socket address structure.

<sup>95-110</sup> The seven function pointers, rnh\_addaddr through rnh\_walktree, point to functions that are called to operate on the tree. Only four of these pointers are initialized by rn\_inithead and the other three are never used by Net/3, as shown in Figure 18.17.

Member	Initialized to (by rn_inithead)
rnh_addaddr	rn_addroute
rnh_addpkt	NULL
rnh_deladdr	rn_delete
rnh_delpkt	NULL
rnh_matchaddr	rn_match
rnh_matchpkt	NULL
rnh_walktree	rn_walktree

Figure 18.17 The seven function pointers in the radix\_node\_head structure.

<sup>111-112</sup> Figure 18.18 shows the radix\_node structure that forms the nodes of the tree. In Figure 18.8 we see that three of these are allocated in the radix\_node\_head and two are allocated in each rtentry structure.

```
- radix.h
40 struct radix_node {
    struct radix_mask *rn_mklist; /* list of masks contained in subtree */
41
       struct radix_node *rn_p; /* parent pointer */
42
      short rn_b;
                                   /* bit offset; -1-index(netmask) */
43
                                   /* node: mask for bit test */
44
      char
               rn_bmask;
       u_char rn_flags;
                                  /* Figure 18.20 */
45
46
      union {
                                   /* leaf only data: rn_b < 0 */</pre>
47
           struct {
               caddr_t rn_Key; /* object of search */
caddr_t rn_Mask; /* netmask, if present */
48
               caddr_t rn_Mask;
49
50
               struct radix_node *rn_Dupedkey;
51
           } rn_leaf;
                                    /* node only data: rn_b >= 0 */
           struct {
52
                      rn_Off;
                                  /* where to start compare */
53
               int
               struct radix_node *rn_L; /* left pointer */
54
               struct radix_node *rn_R;
55
                                          /* right pointer */
56
           } rn_node;
57
       } rn_u;
58 };
59 #define rn_dupedkey rn_u.rn_leaf.rn_Dupedkey
                      rn_u.rn_leaf.rn_Key
60 #define rn_key
```

**INTEL Ex.1013.601** 

Chapter 18

Ş

2

61 #define rn 62 #define rn 63 #define rn	off rn_u.ri	n_leaf.rn_Mask n_node.rn_Off n_node.rn_L
64 #define rn	·	n_node.rn_R radix.h

Figure 18.18 radix\_node structure: the nodes of the routing tree.

<sup>41–45</sup> The first five members are common to both internal nodes and leaves, followed by a union defining three members if the node is a leaf, or a different three members if the node is internal. As is common throughout the Net/3 code, a set of #define statements provide shorthand names for the members in the union.

41-42 rn\_mklist is the head of a linked list of masks for this node. We describe this field in Section 18.9. rn\_p points to the parent node.

43

If rn\_b is greater than or equal to 0, the node is an internal node, else the node is a leaf. For the internal nodes, rn\_b is the bit number to test: for example, its value is 32 in the top node of the tree in Figure 18.4. For leaves, rn\_b is negative and its value is -1 minus the *index of the network mask*. This index is the first bit number where a 0 occurs. Figure 18.19 shows the indexes of the masks from Figure 18.4.

·			32-bi	t IP mas	k (bits 32	2-63)			index	rn_b
	3333	3333	4444	4444	4455	5555	5555	6666		
	2345	6789	0123	4567	8901	2345	6789	0123		
00000000:	0000	0000	0000	0000	0000	0000	0000	0000	0	-1
ff000000:	1111	1111	0000	0000	0000	0000	0000	0000	40	-41
ffffffe0:	1111	1111	1111	1111	1111	1111	1110	0000	59	60

## Figure 18.19 Example of mask indexes.

As we can see, the index of the all-zero mask is handled specially: its index is 0, not 32. rn\_bmask is a 1-byte mask used with the internal nodes to test whether the corresponding bit is on or off. Its value is 0 in leaves. We'll see how this member is used

45

44

Figure 18.20 shows the three values for the rn\_flags member.

with the rn\_off member shortly.

Constant	Description
RNF_ACTIVE	this node is alive (for rtfree)
RNF_NORMAL	leaf contains normal route (not currently used)
RNF_ROOT	node is in the radix_node_head structure

Figure 18.20 rn\_flags values.

The RNF\_ROOT flag is set only for the three radix nodes in the radix\_node\_head structure: the top of the tree and the left and right end nodes. These three nodes can never be deleted from the routing tree.



<sup>48–49</sup> For a leaf, rn\_key points to the socket address structure and rn\_mask points to a socket address structure containing the mask. If rn\_mask is null, the implied mask is all one bits (i.e., this route is to a host, not to a network).

Figure 18.21 shows an example corresponding to the leaf for 140.252.13.32 in Figure 18.4.

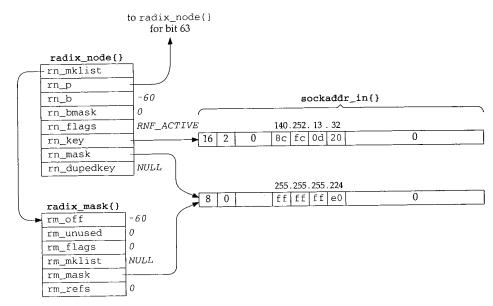


Figure 18.21 radix\_node structure corresponding to leaf for 140.252.13.32 in Figure 18.4.

This example also shows a radix\_mask structure, which we describe in Figure 18.22. We draw this latter structure with a smaller width, to help distinguish it as a different structure from the radix\_node; we'll encounter both structures in many of the figures that follow. We describe the reason for the radix\_mask structure in Section 18.9.

The rn\_b of -60 corresponds to an index of 59. rn\_key points to a sockaddr\_in, with a length of 16 and an address family of 2 (AF\_INET). The mask structure pointed to by rn\_mask and rm\_mask has a length of 8 and a family of 0 (this family is AF\_UNSPEC, but it is never even looked at).

<sup>50-51</sup> The rn\_dupedkey pointer is used when there are multiple leaves with the same key. We describe these in Section 18.9.

52-58 We describe rn\_off in Section 18.8. rn\_l and rn\_r are the left and right pointers for the internal node.

Figure 18.22 shows the radix\_mask structure.

er 18

dix.h

by a f the

tate-

field

e is a

is 32

is –1

curs.

: 32.

corre-

used

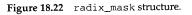
\_head

es can

Section 18.5

ab

Chapter 18 Radix Tree Routing Tables radix h 76 extern struct radix\_mask { /\* bit offset; -1-index(netmask) \*/ short rm\_b; 77 /\* cf. rn\_bmask \*/ rm\_unused; /\* cf. rn\_flags \*/ 78 char u\_char rm\_flags; 79 struct radix\_mask \*rm\_mklist; /\* more masks to try \*/ 80 /\* the mask \*/ caddr\_t rm\_mask; 81 /\* # of references to this struct \*/ rm\_refs; int 82 \*rn\_mkfreelist; 83 } radix.h



<sup>76-83</sup> Each of these structures contains a pointer to a mask: rm\_mask, which is really a pointer to a socket address structure containing the mask. Each radix\_node structure points to a linked list of radix\_mask structures, allowing multiple masks per node: rn\_mklist points to the first, and then each rm\_mklist points to the next. This structure definition also declares the global rn\_mkfreelist, which is the head of a linked list of available structures.

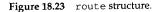
## **18.6 Routing Structures**

578

The focal points of access to the kernel's routing information are

- 1. the rtalloc function, which searches for a route to a destination,
- 2. the route structure that is filled in by this function, and
- 3. the rtentry structure that is pointed to by the route structure.

Figure 18.8 showed that the protocol control blocks (PCBs) used by UDP and TCP (Chapter 22) contain a route structure, which we show in Figure 18.23.



ro\_dst is declared as a generic socket address structure, but for the Internet protocols it is a sockaddr\_in. Notice that unlike most references to this type of structure, ro dst is the structure itself, not a pointer to one.

At this point it is worth reviewing Figure 8.24, which shows the use of these routes every time an IP datagram is output.

• If the caller passes a pointer to a route structure, that structure is used. Otherwise a local route structure is used and it is set to 0, setting ro\_rt to a null pointer. UDP and TCP pass a pointer to the route structure in their PCB to ip\_output.

route.h

pter 18 radix.h •radix.h ip\_output returns an error. eally a ucture ' node: 3 structhe destination address of the IP datagram. linked Figure 18.24 shows the rtentry structure. 83 struct rtentry { 84 85 struct sockaddr \*rt\_gateway; 86 short rt flags; 87 short rt\_refcnt; 88 u\_long rt\_use; 89 struct ifnet \*rt ifp: struct ifaddr \*rt\_ifa; 90 struct sockaddr \*rt\_genmask; 91 92 caddr\_t rt\_llinfo; 93 ıd TCP 94 95 }; route.h 96 #define rt\_key(r) 97 #define rt\_mask(r) \* / – route.h 83-84 rotocols ructure, 86 e routes Othero a null mally stored in the routing table entry. PCB to 85

If the route structure points to an rtentry structure (the ro\_rt pointer is nonnull), and if the referenced interface is still up, and if the destination address in the route structure equals the destination address of the IP datagram, that route is used. Otherwise the socket address structure ro\_dst is filled in with the destination IP address and rtalloc is called to locate a route to that destination. For a TCP connection the destination address of the datagram never changes from the destination address of the route, but a UDP application can send a datagram to a different destination with each sendto.

- If rtalloc returns a null pointer in ro\_rt, a route was not found and
- If the RTF\_GATEWAY flag is set in the rtentry structure, the route is indirect (the G flag in Figure 18.2). The destination address (dst) for the interface output function becomes the IP address of the gateway, the rt\_gateway member, not

route.h struct radix\_node rt\_nodes[2]; /\* a leaf and an internal node \*/ /\* value associated with rn\_key \*/ /\* Figure 18.25 \*/ /\* #held references \*/ /\* raw #packets sent \*/ /\* interface to use \*/ /\* interface address to use \*/ /\* for generation of cloned routes \*/ /\* pointer to link level info cache \*/ struct rt\_metrics rt\_rmx; /\* metrics: Figure 18.26 \*/ struct rtentry \*rt\_gwroute; /\* implied entry for gatewayed routes \*/ ((struct sockaddr \*)((r)->rt\_nodes->rn\_key)) ((struct sockaddr \*)((r)->rt\_nodes->rn\_mask)) route.h

Figure 18.24 rtentry structure.

Two radix node structures are contained within this structure. As we noted in the example with Figure 18.7, each time a new leaf is added to the routing tree a new internal node is also added. rt\_nodes[0] contains the leaf entry and rt\_nodes[1] contains the internal node. The two #define statements at the end of Figure 18.24 provide a shorthand access to the key and mask of this leaf node.

Figure 18.25 shows the various constants stored in rt\_flags and the corresponding character output by netstat in the "Flags" column (Figure 18.2).

The RTF\_BLACKHOLE flag is not output by netstat and the two with lowercase flag characters, RTF\_DONE and RTF\_MASK, are used in routing messages and not nor-

If the RTF\_GATEWAY flag is set, rt\_gateway contains a pointer to a socket address structure containing the address (e.g., the IP address) of that gateway. Also,

87

88

Chapter 18

Constant	netstat flag	Description	
RTF_BLACKHOLE		discard packets without error (loopback driver: Figure 5.27)	
RTF_CLONING	С	generate new routes on use (used by ARP)	
RTF_DONE	d	kernel confirmation that message from process was completed	
RTF_DYNAMIC	D	created dynamically (by redirect)	
RTF_GATEWAY	G	destination is a gateway (indirect route)	
RTF_HOST	н	host entry (else network entry)	
RTF_LLINFO	L	set by ARP when rt_llinfo pointer valid	
RTF_MASK	m	subnet mask present (not used)	
RTF_MODIFIED	M	modified dynamically (by redirect)	
RTF_PROTO1	1	protocol-specific routing flag	
RTF_PROTO2	2	protocol-specific routing flag (ARP uses)	
RTF_REJECT	R	discard packets with error (loopback driver: Figure 5.27)	
RTF_STATIC	S	manually added entry (route program)	
RTF_UP	U	route usable	
RTF_XRESOLVE	X	external daemon resolves name (used with X.25)	

Figure 18.25 rt\_flags values.

rt\_gwroute points to the rtentry for that gateway. This latter pointer was used in ether\_output (Figure 4.15).

- rt\_refcnt counts the "held" references to this structure. We describe this counter at the end of Section 19.3. This counter is output as the "Refs" column in Figure 18.2.
- rt\_use is initialized to 0 when the structure is allocated; we saw it incremented in Figure 8.24 each time an IP datagram was output using the route. This counter is also the value printed in the "Use" column in Figure 18.2.

89-90 rt\_ifp and rt\_ifa point to the interface structure and the interface address structure, respectively. Recall from Figure 6.5 that a given interface can have multiple addresses, so minimally the rt\_ifa is required.

<sup>92</sup> The rt\_llinfo pointer allows link-layer protocols to store pointers to their protocol-specific structures in the routing table entry. This pointer is normally used with the RTF\_LLINFO flag. Figure 21.1 shows how ARP uses this pointer.

			Toute:
54 s	truct rt_m	etrics {	
55	u_long	rmx_locks;	/* bitmask for values kernel leaves alone */
56	u_long	rmx_mtu;	/* MTU for this path */
57	u_long	rmx_hopcount;	/* max hops expected */
58	u_long	rmx_expire;	/* lifetime for route, e.g. redirect */
59	u long	rmx recvpipe;	/* inbound delay-bandwith product */
60	u long	rmx_sendpipe;	/* outbound delay-bandwith product */
61	u long	rmx_ssthresh;	/* outbound gateway buffer limit */
62	u long	rmx rtt;	/* estimated round trip time */
63	u long	rmx rttvar;	/* estimated RTT variance */
64	u_long	rmx_pksent;	<pre>/* #packets sent using this route */</pre>
65 }	;		route h

#### Figure 18.26 rt\_metrics structure.

INTEL Ex.1013.606

route h

93

Figure 18.26 shows the rt\_metrics structure, which is contained within the rtentry structure. Figure 27.3 shows that TCP uses six members in this structure.

54-65 rmx\_locks is a bitmask telling the kernel which of the eight metrics that follow must not be modified. The values for this bitmask are shown in Figure 20.13.

rmx\_expire is used by ARP (Chapter 21) as a timer for each ARP entry. Contrary to the comment with rmx\_expire, it is not used for redirects.

Figure 18.28 summarizes the structures that we've described, their relationships, and the various types of socket address structures they reference. The rtentry that we show is for the route to 128.32.33.5 in Figure 18.2. The other radix\_node contained in the rtentry is for the bit 36 test right above this node in Figure 18.4. The two sockaddr\_dl structures pointed to by the first ifaddr were shown in Figure 3.38. Also note from Figure 6.5 that the ifnet structure is contained within an le\_softc structure, and the second ifaddr structure is contained within an in\_ifaddr structure.

# 18.7 Initialization: route\_init and rtable\_init Functions

The initialization of the routing tables is somewhat obscure and takes us back to the domain structures in Chapter 7. Before outlining the function calls, Figure 18.27 shows the relevant fields from the domain structure (Figure 7.5) for various protocol families.

Member	OSI value	Internet value	Routing value	Unix value	XNS value	Comment
dom_family	AF_ISO	AF_INET	PF_ROUTE	AF_UNIX	AF_NS	
dom_init	0	0	route_init	0	0	
dom_rtattach	rn_inithead	rn_inithead	0	0	rn_inithead	
dom_rtoffset	48	32	0	0	16	in bits
dom_maxrtkey	32	16	0	0	16	in bytes

Figure 18.27 Members of domain structure relevant to routing.

The PF\_ROUTE domain is the only one with an initialization function. Also, only the domains that require a routing table have a dom\_rtattach function, and it is always rn\_inithead. The routing domain and the Unix domain protocols do not require a routing table.

The dom\_rtoffset member is the offset, in bits, (from the beginning of the domain's socket address structure) of the first bit to be examined for routing. The size of this structure in bytes is given by dom\_maxrtkey. We saw earlier in this chapter that the offset of the IP address in the sockaddr\_in structure is 32 bits. The dom\_maxrtkey member is the size in bytes of the protocol's socket address structure: 16 for sockaddr\_in.

INTEL Ex.1013.607

Figure 18.29 outlines the steps involved in initializing the routing tables.

18

1 in

nter

d in also

ruciple

otothe

ute.h

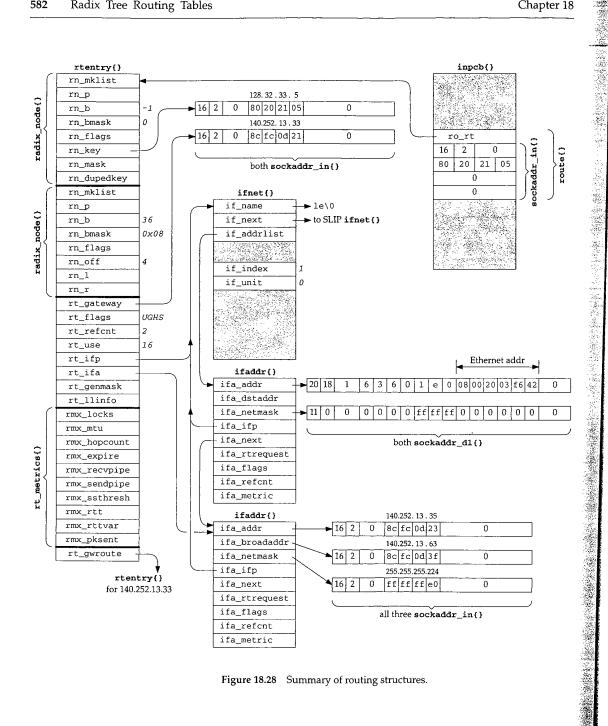
e \*/

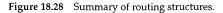
pute.h

and a start of the s

ł

Chapter 18





主体が多くのない

```
main()
                    /* kernel initialization */
  {
       . . .
      ifinit();
      -domaininit();
 domaininit()
                    /* Figure 7.15 */
 {
      . . .
      ADDDOMAIN(unix);
      ADDDOMAIN(route);
      ADDDOMAIN(inet);
      ADDDOMAIN(osi);
      • • •
      for ( dp = all domains ) {
           (*dp->dom_init)();
           for ( pr = all protocols for this domain )
                -(*pr->pr_init)();
 }
 raw_init()
                   /* pr_init() function for SOCK_RAW/PF_ROUTE protocol */
 ł
     initialize head of routing protocol control blocks;
route_init()
                   /* dom_init() function for PF_ROUTE domain */
     -rn_init();
    --rtable_init();
3
rn_init()
     for ( dp = all domains )
          if (dp->dom_maxrtkey > max_keylen)
               max_keylen = dp->dom_maxrtkey;
     allocate and initialize rn_zeros, rn_ones, masked_key;
     -rn_inithead(&mask_rnhead); /* allocate and init tree for masks */
rtable_init()
{
     for (dp = all domains)
          (*dp->dom_rtattach)(&rt_tables[dp->dom_family]);
rn_inithead()
                  /* dom_rtattach() function for all protocol families */
{
    allocate and initialize one radix_node_head structure;
}
```

Figure 18.29 Steps involved in initialization of routing tables.

route.c

route.c

route.c

route.c

radix.c

domaininit is called once by the kernel's main function when the system is initialized. The linked list of domain structures is built by the ADDDOMAIN macro and the linked list is traversed, calling each domain's dom\_init function, if defined. As we saw in Figure 18.27, the only dom\_init function is route\_init, which is shown in Figure 18.30.

```
49 void
50 route_init()
51 {
52 rn_init(); /* initialize all zeros, all ones, mask table */
53 rtable_init((void **) rt_tables);
54 }
```

## Figure 18.30 route\_init function.

The function rn\_init, shown in Figure 18.32, is called only once.

The function rtable\_init, shown in Figure 18.31, is also called only once. It in turn calls all the dom\_rtattach functions, which initialize a routing table tree for that domain.

39	void
40	rtable_init(table)
41	void **table;
42	{
43	struct domain *dom;
44	for (dom = domains; dom; dom = dom->dom_next)
45	if (dom->dom_rtattach)
46	<pre>dom-&gt;dom_rtattach(&amp;table[dom-&gt;dom_family],</pre>
47	dom->dom_rtoffset);
48	}

Figure 18.31 rtable\_init function: call each domain's dom\_rtattach function.

We saw in Figure 18.27 that the only dom\_rtattach function is rn\_inithead, which we describe in the next section.

# 18.8 Initialization: rn\_init and rn\_inithead Functions

The function rn\_init, shown in Figure 18.32, is called once by route\_init to initialize some of the globals used by the radix functions.

750	void		
751	<pre>rn_init()</pre>		
752	{		
753	char	*cp, *cplim;	
754	struct	domain *dom;	

Section 18.8 Initialization: rn\_init and rn\_inithead Functions 585 18 755 for (dom = domains; dom; dom = dom->dom\_next) ni-756 if (dom->dom\_maxrtkey > max\_keylen) he 757 max\_keylen = dom->dom\_maxrtkey; we 758 if (max\_keylen == 0) { in 759 printf("rn\_init: radix functions require max\_keylen be set\n"); 760 return; 761 } te.c 762 R\_Malloc(rn\_zeros, char \*, 3 \* max\_keylen); 763 if (rn\_zeros == NULL) 764 panic("rn\_init"); 765 Bzero(rn\_zeros, 3 \* max\_keylen); 766 rn\_ones = cp = rn\_zeros + max\_keylen; 767 maskedKey = cplim = rn\_ones + max\_keylen; 768 while (cp < cplim) :te.c 769 \*cp++ = -1;770 if (rn\_inithead((void \*\*) &mask\_rnhead, 0) == 0) 771 panic("rn\_init 2"); 772 } radix.c t in Figure 18.32 rn\_init function. :hat Determine max\_keylen ıte.c All the domain structures are examined and the global max\_keylen is set to the 750-761 largest value of dom\_maxrtkey. In Figure 18.27 the largest value is 32 for AF\_ISO, but in a typical system that excludes the OSI and XNS protocols, max\_keylen is 16, the size of a sockaddr\_in structure. Allocate and initialize rn\_zeros, rn\_ones, and maskedKey A buffer three times the size of max\_keylen is allocated and the pointer stored in 762-769 the global rn\_zeros. R\_Malloc is a macro that calls the kernel's malloc function, specifying a type of M\_RTABLE and M\_DONTWAIT. We'll also encounter the macros ute.c Bcmp, Bcopy, Bzero, and Free, which call kernel functions of similar names, with the arguments appropriately type cast. This buffer is divided into three pieces, and each piece is initialized as shown in Figead. ure 18.33. max\_keylen bytes max\_keylen bytes max\_keylen bytes 0 0 0 1 1 1 1 0 0 0 0 0 0 10 0 0 1 1 . . . . . . . . . itialrn\_zeros rn\_ones maskedkey Figure 18.33 rn\_zeros, rn\_ones, and maskedKey arrays. adix.c rn\_zeros is an array of all zero bits, rn\_ones is an array of all one bits, and maskedKey is an array used to hold a temporary copy of a search key that has been masked.

# INTEL Ex.1013.611

i,

S

## Initialize tree of masks

770-772 The function rn\_inithead is called to initialize the head of the routing tree for the address masks; the radix\_node\_head structure pointed to by the global mask\_rnhead in Figure 18.8.

From Figure 18.27 we see that rn\_inithead is also the dom\_attach function for all the protocols that require a routing table. Instead of showing the source code for this function, Figure 18.34 shows the radix\_node\_head structure that it builds for the Internet protocols.

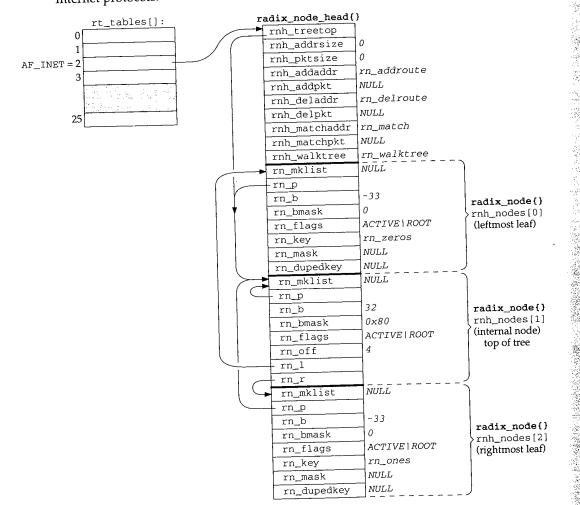


Figure 18.34 radix\_node\_head structure built by rn\_inithead for Internet protocols.

The three radix\_node structures form a tree: the middle of the three is the top (it is pointed to by rnh\_treetop), the first of the three is the leftmost leaf of the tree, and

the last of the three is the rightmost leaf of the tree. The parent pointer of all three nodes (rn\_p) points to the middle node.

The value 32 for rnh\_nodes [1].rn\_b is the bit position to test. It is from the dom\_rtoffset member of the Internet domain structure (Figure 18.27). Instead of performing shifts and masks during forwarding, the byte offset and corresponding byte mask are precomputed. The byte offset from the start of a socket address structure is in the rn\_off member of the radix\_node structure (4 in this case) and the byte mask is in the rn\_bmask member (0x80 in this case). These values are computed whenever a radix\_node structure is added to the tree, to speed up the comparisons during forwarding. As additional examples, the offset and byte mask for the two nodes that test bit 33 in Figure 18.4 would be 4 and 0x40, respectively. The offset and byte mask for the two nodes that test bit 63 would be 7 and 0x01.

The value of -33 for the rn\_b member of both leaves is negative one minus the index of the leaf.

The key of the leftmost node is all zero bits (rn\_zeros) and the key of the rightmost node is all one bits (rn\_ones).

All three nodes have the RNF\_ROOT flag set. (We have omitted the RNF\_ prefix.) This indicates that the node is one of the three original nodes used to build the tree. These are the only nodes with this flag.

One detail we have not mentioned is that the Network File System (NFS) also uses the routing table functions. For each mount point on the local host a radix\_node\_head structure is allocated, along with an array of pointers to these structures (indexed by the protocol family), similar to the rt\_tables array. Each time this mount point is exported, the protocol address of the host that can mount this filesystem is added to the appropriate tree for the mount point.

## 18.9 Duplicate Keys and Mask Lists

Before looking at the source code that looks up entries in a routing table we need to understand two fields in the radix\_node structure: rn\_dupedkey, which forms a linked list of additional radix\_node structures containing duplicate keys, and rn\_mklist, which starts a linked list of radix\_mask structures containing network masks.

We first return to Figure 18.4 and the two boxes on the far left of the tree labeled "end" and "default." These are duplicate keys. The leftmost node with the RNF\_ROOT flag set (rnh\_nodes [0] in Figure 18.34) has a key of all zero bits, but this is the same key as the default route. We would have the same problem with the rightmost end node in the tree, which has a key of all one bits, if an entry were created for 255.255.255, but this is the limited broadcast address, which doesn't appear in the routing table. In general, the radix node functions in Net/3 allow any key to be duplicated, if each occurrence has a unique mask.

Figure 18.35 shows the two nodes with a duplicate key of all zero bits. In this figure we have removed the RNF\_ prefix for the rn\_flags and omit nonnull parent, left, and right pointers, which add nothing to the discussion.

the bal for :his

the

·18

.] )

()

2]

Ð

588 Radix Tree Routing Tables Chapter 18

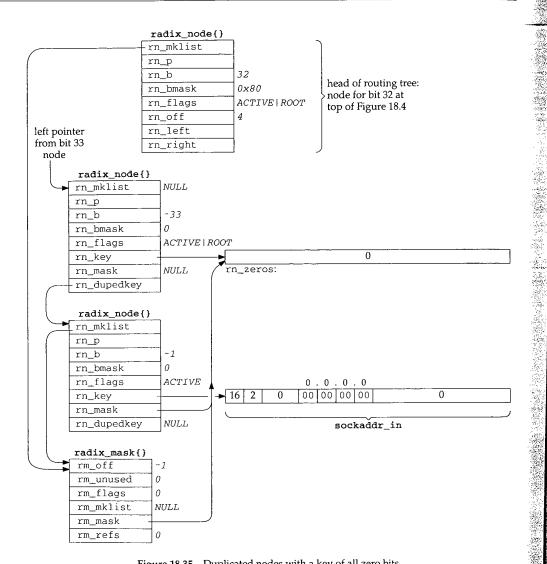


Figure 18.35 Duplicated nodes with a key of all zero bits.

The top node is the top of the routing tree-the node for bit 32 at the top of Figure 18.4. The next two nodes are leaves (their rn\_b values are negative) with the rn\_dupedkey member of the first pointing to the second. The first of these two leaves is the rnh\_nodes[0] structure from Figure 18.34, which is the left end marker of the tree—its RNF\_ROOT flag is set. Its key was explicitly set by rn\_inithead to rn\_zeros.

The second of these leaves is the entry for the default route. Its rn\_key points to a sockaddr\_in with the value 0.0.0.0, and it has a mask of all zero bits. Its rn\_mask points to rn\_zeros, since equivalent masks in the mask table are shared.

Duplicate Keys and Mask Lists 589

Section 18.9

Normally keys are not shared, let alone shared with masks. The rn\_key pointers of the two end markers (those with the RNF\_ROOT flag) are special since they are built by rn\_inithead (Figure 18.34). The key of the left end marker points to rn\_zeros and the key of the right end marker points to rn\_ones.

The final structure is a radix\_mask structure and is pointed to by both the top node of the tree and the leaf for the default route. The list from the top node of the tree is used with the backtracking algorithm when the search is looking for a network mask. The list of radix\_mask structures with an internal node specifies the masks that apply to subtrees starting at that node. In the case of duplicate keys, a mask list also appears with the leaves, as we'll see in the following example.

We now show a duplicate key that is added to the routing tree intentionally and the resulting mask list. In Figure 18.4 we have a host route for 127.0.0.1 and a network route for 127.0.0.0. The default mask for the class A network route is 0xff000000, as we show in the figure. If we divide the 24 bits following the class A network ID into a 16-bit subnet ID and an 8-bit host ID, we can add a route for the subnet 127.0.0 with a mask of 0xfffff00:

#### bsdi \$ route add 127.0.0.0 -netmask 0xffffff00 140.252.13.33

Although it makes little practical sense to use network 127 in this fashion, our interest is in the resulting routing table structure. Although duplicate keys are not common with the Internet protocols (other than the previous example with the default route), duplicate keys are required to provide routes to subnet 0 of any network.

There is an implied priority in these three entries with a network ID of 127. If the search key is 127.0.0.1 it matches all three entries, but the host route is selected because it is the *most specific*: its mask ( $0 \times ffffffff$ ) has the most one bits. If the search key is 127.0.0.2 it matches both network routes, but the route for subnet 0, with a mask of  $0 \times ffffff00$ , is more specific than the route with a mask of  $0 \times ff000000$ . The search key 127.1.2.3 matches only the entry with a mask of  $0 \times ff000000$ .

Figure 18.36 shows the resulting tree structure, starting at the internal node for bit 33 from Figure 18.4. We show two boxes for the entry with the key of 127.0.0.0 since there are two leaves with this duplicate key.

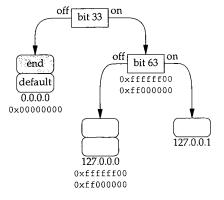


Figure 18.36 Routing tree showing duplicate keys for 127.0.0.0.

g-

he

es he

to

) a

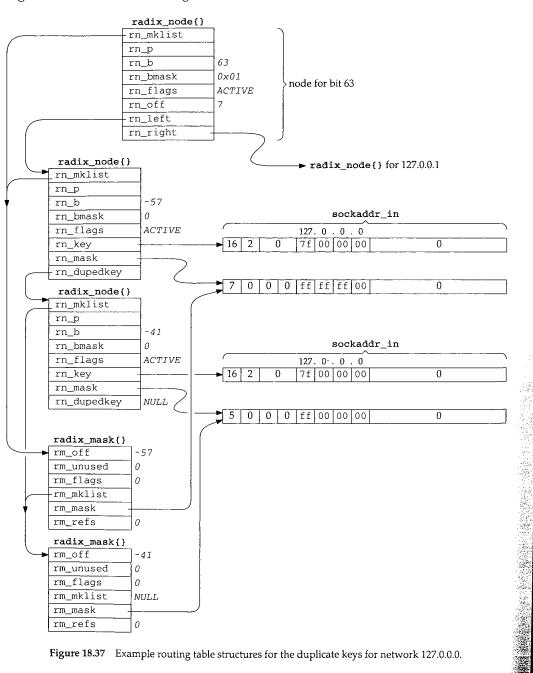
зk

in the second

P

a)

(2) A state of the state of



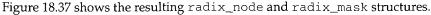


Figure 18.37 Example routing table structures for the duplicate keys for network 127.0.0.0.

Sec

Chapter 18

**INTEL Ex.1013.616** 

radix.c

INTEL Ex.1013.617

First look at the linked list of radix\_mask structures for each radix\_node. The mask list for the top node (bit 63) consists of the entry for  $0 \times ffffff00$  followed by  $0 \times ff000000$ . The more-specific mask comes first in the list so that it is tried first. The mask list for the second radix\_node (the one with the rn\_b of -57) is the same as that of the first. But the list for the third radix\_node consists of only the entry with a mask of  $0 \times ff000000$ .

Notice that masks with the same value are shared but keys with the same value are not. This is because the masks are maintained in their own routing tree, explicitly to be shared, because equal masks are so common (e.g., every class C network route has the same mask of 0xfffff00), while equal keys are infrequent.

## 18.10 rn\_match Function

We now show the rn\_match function, which is called as the rnh\_matchaddr function for the Internet protocols. We'll see that it is called by the rtalloc1 function, which is called by the rtalloc function. The algorithm is as follows:

- 1. Start at the top of the tree and go to the leaf corresponding to the bits in the search key. Check the leaf for an exact match (Figure 18.38).
- 2. Check the leaf for a network match (Figure 18.40).
- 3. Backtrack (Figure 18.43).

Figure 18.38 shows the first part of rn\_match.

```
135 struct radix node *
136 rn_match(v_arg, head)
137 void
         *v_arg;
138 struct radix_node_head *head;
139 (
       caddr_t v = v_arg;
140
141
       struct radix_node *t = head->rnh_treetop, *x;
142
       caddr_t cp = v, cp2, cp3;
143
       caddr_t cplim, mstart;
144
       struct radix_node *saved_t, *top = t;
145
              off = t->rn_off, vlen = *(u_char *) cp, matched_off;
       int
146
       /*
147
        * Open code rn_search(v, top) to avoid overhead of extra
        * subroutine call.
148
149
        */
150
       for (; t->rn_b >= 0;) {
151
          if (t->rn_bmask & cp[t->rn_off])
152
              153
           else
                                /* left if bit off */
154
              t = t->rn_l;
155
       }
```

l8

592 Radix Tree Routing Tables

Chapter 18

```
156
           See if we match exactly as a host destination
157
         */
158
159
        cp += off;
        cp2 = t->rn_key + off;
160
161
        cplim = v + vlen;
        for (; cp < cplim; cp++, cp2++)
162
            if (*cp != *cp2)
163
                goto on1;
164
165
        /*
         * This extra grot is in case we are explicitly asked
166
         * to look up the default. Ugh!
167
         * /
168
        if ((t->rn_flags & RNF_ROOT) && t->rn_dupedkey)
169
            t = t->rn_dupedkey;
170
171
        return t;
172
      on1:
```

- radix.c

Figure 18.38 rn\_match function: go down tree, check for exact host match.

The first argument v\_arg is a pointer to a socket address structure, and the second argument head is a pointer to the radix\_node\_head structure for the protocol. All protocols call this function (Figure 18.17) but each calls it with a different head argument.

In the assignment statements, off is the rn\_off member of the top node of the tree (4 for Internet addresses, from Figure 18.34), and vlen is the length field from the socket address structure of the search key (16 for Internet addresses).

### Go down the tree to the corresponding leaf

<sup>146-155</sup> This loop starts at the top of the tree and moves down the left and right branches until a leaf is encountered (rn\_b is less than 0). Each test of the appropriate bit is made using the precomputed byte mask in rn\_bmask and the corresponding precomputed offset in rn\_off. For Internet addresses, rn\_off will be 4, 5, 6, or 7.

#### Check for exact match

<sup>156–164</sup> When the leaf is encountered, a check is first made for an exact match. All bytes of the socket address structure, starting at the rn\_off value for the protocol family, are compared. This is shown in Figure 18.39 for an Internet socket address structure.

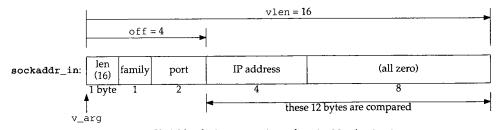


Figure 18.39 Variables during comparison of sockaddr\_in structures.

As soon as a mismatch is found, a jump is made to on1.

「「「「「「「「」」」」

のこうでは「読ん」語言語

Normally the final 8 bytes of the sockaddr\_in are 0 but proxy ARP (Section 21.12) sets one of these bytes nonzero. This allows two routing table entries for a given IP address: one for the normal IP address (with the final 8 bytes of 0) and a proxy ARP entry for the same IP address (with one of the final 8 bytes nonzero).

The length byte in Figure 18.39 was assigned to vlen at the beginning of the function, and we'll see that rtalloc1 uses the family member to select the routing table to search. The port is never used by the routing functions.

### Explicit check for default

165-172

Figure 18.35 showed that the default route is stored as a duplicate leaf with a key of 0. The first of the duplicate leaves has the RNF\_ROOT flag set. Hence if the RNF\_ROOT flag is set in the matching node and the leaf contains a duplicate key, the value of the pointer rn\_dupedkey is returned (i.e., the pointer to the node containing the default route in Figure 18.35). If a default route has not been entered and the search matches the left end marker (a key of all zero bits), or if the search encounters the right end marker (a key of all one bits), the returned pointer t points to a node with the RNF\_ROOT flag set. We'll see that rtalloc1 explicitly checks whether the matching node has this flag set, and considers such a match an error.

At this point in rn\_match a leaf has been reached but it is not an exact match with the search key. The next part of the function, shown in Figure 18.40, checks whether the leaf is a network match.

		radix.c
173	$matched_off = cp - v;$	
174	<pre>saved_t = t;</pre>	
175	do {	
176	if (t->rn_mask) {	
177	/*	
178	* Even if we don't match exactly as a host;	
179	* we may match if the leaf we wound up at is	
180	* a route to a net.	
181	*/	
182	cp3 = matched_off + t->rn_mask;	
183	cp2 = matched_off + t->rn_key;	
184	for (; cp < cplim; cp++)	
185	if ((*cp2++ ^ *cp) & *cp3++)	
186	break;	
187	if (cp == cplim)	
188	return t;	
189	$cp = matched_off + v;$	
190	}	
191	<pre>} while (t = t-&gt;rn_dupedkey);</pre>	
192	t = saved_t;	radix.c

Figure 18.40 rn\_match function: check for network match.

<sup>173–174</sup> cp points to the unequal byte in the search key. matched\_off is set to the offset of this byte from the start of the socket address structure.

<sup>175-183</sup> The do while loop iterates through all duplicate leaves and each one with a network mask is compared. Let's work through the code with an example. Assume we're

Chapter 18

-

Sec

191

looking up the IP address 140.252.13.60 in the routing table in Figure 18.4. The search will end up at the node labeled 140.252.13.32 (bits 62 and 63 are both off), which contains a network mask. Figure 18.41 shows the structures when the for loop in Figure 18.40 starts executing.

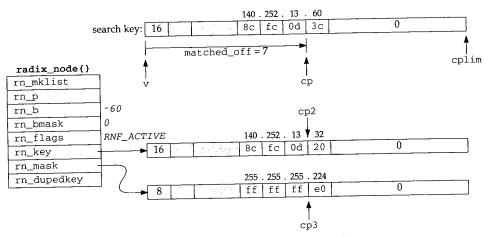


Figure 18.41 Example for network mask comparison.

The search key and the routing table key are both sockaddr\_in structures, but the length of the mask is different. The mask length is the minimum number of bytes containing nonzero values. All the bytes past this point, up through max\_keylen, are 0.

The search key is exclusive ORed with the routing table key, and the result logically ANDed with the network mask, one byte at a time. If the resulting byte is ever nonzero, the loop terminates because they don't match (Exercise 18.1). If the loop terminates normally, however, the search key ANDed with the network mask matches the routing table entry. The pointer to the routing table entry is returned.

Figure 18.42 shows how this example matches, and how the IP address 140.252.13.188 does not match, looking at just the fourth byte of the IP address. The search for both IP addresses ends up at this node since both addresses have bits 57, 62, and 63 off.

	search key = 140.252.13.60	search key = 140.252.13.188
search key byte (*cp): routing table key byte (*cp2):	$\begin{array}{r} 0011 \ 1100 \ = \ 3c \\ 0010 \ 0000 \ = \ 20 \end{array}$	$\begin{array}{rrrr} 1011 \ 1100 \ = \ bc \\ 0010 \ 0000 \ = \ 20 \end{array}$
exclusive OR: network mask byte (*cp3):	0001 1100 1110 0000 = e0	1001 1100 1110 0000 = e0
logical AND:	0000 0000	1000 0000

Figure 18.42 Example of search key match using network mask.

The first example (140.252.13.60) matches since the result of the logical AND is 0 (and all the remaining bytes in the address, the key, and the mask are all 0). The other example does not match since the result of the logical AND is nonzero.

184-190

19. 19.

19

Contraction of the

191

oter 18

earch

l COn-

n Fig-

plim

it the

3 con-

ically

ızero,

s nor-

uting

dress

7,62,

The

Į.

÷0.

If the routing table entry has duplicate keys, the loop is repeated for each key.

The final portion of rn\_match, shown in Figure 18.43, backtracks up the tree, looking for a network match or a match with the default.

```
– radix.c
193
         /* start searching up the tree */
194
        do {
195
             struct radix_mask *m;
196
             t = t - rn_p;
197
             if (m = t->rn_mklist) {
198
                 /*
199
                  * After doing measurements here, it may
200
                  * turn out to be faster to open code
                  * rn_search_m here instead of always
201
202
                  * copying and masking.
                  */
203
204
                 off = min(t->rn_off, matched_off);
                 mstart = maskedKey + off;
205
206
                 do {
207
                     cp2 = mstart;
208
                     cp3 = m->rm_mask + off;
209
                     for (cp = v + off; cp < cplim;)</pre>
210
                         *cp2++ = *cp++ & *cp3++;
211
                     x = rn_search(maskedKey, t);
212
                     while (x && x->rn_mask != m->rm_mask)
213
                         x = x->rn_dupedkey;
214
                     if (x &&
215
                         (Bcmp(mstart, x->rn_key + off,
216
                               vlen - off) == 0))
217
                         return x;
                 } while (m = m->rm_mklist);
218
219
             }
220
        } while (t != top);
221
        return 0;
222 };
                                                                               radix c
```

Figure 18.43 rn\_match function: backtrack up the tree.

<sup>193-195</sup> The do while loop continues up the tree, checking each level, until the top has been checked.

The pointer t is replaced with the pointer to the parent node, moving up one level. Having the parent pointer in each node simplifies backtracking.

197-210 Each level is checked only if the internal node has a nonnull list of masks. rn\_mklist is a pointer to a linked list of radix\_mask structures, each containing a mask that applies to the subtree starting at that node. The inner do while loop iterates through each radix\_mask structure on the list.

Using the previous example, 140.252.13.188, Figure 18.44 shows the various data structures when the innermost for loop starts. This loop logically ANDs each byte of the search key with each byte of the mask, storing the result in the global maskedKey. The mask value is 0xfffffe0 and the search would have backtracked from the leaf for 140.252.13.32 in Figure 18.4 two levels to the node that tests bit 62.

nd all mple

INTEL Ex.1013.621

a cira

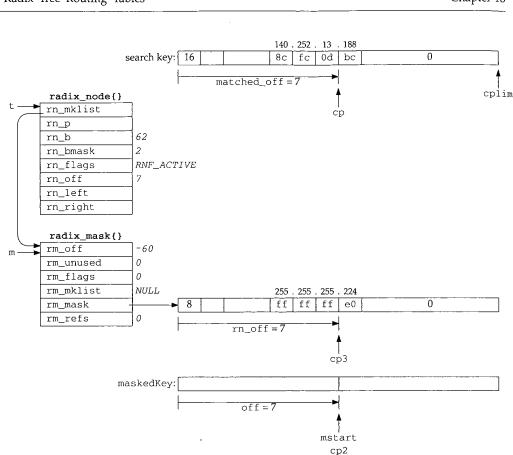
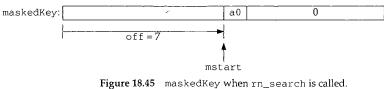


Figure 18.44 Preparation to search again using masked search key.

Once the for loop completes, the masking is complete, and rn\_search (shown in Figure 18.48) is called with maskedKey as the search key and the pointer t as the top of the subtree to search. Figure 18.45 shows the value of maskedKey for our example.



rigure 16.45 maskedkey when th\_search is caned.

The byte 0xa0 is the logical AND of 0xbc (188, the search key) and 0xe0 (the mask). rn\_search proceeds down the tree from its starting point, branching right or left depending on the key, until a leaf is reached. In this example the search key is the 9 bytes shown in Figure 18.45 and the leaf that's reached is the one labeled 140.252.13.32 in Figure 18.4, since bits 62 and 63 are off in the byte 0xa0. Figure 18.46 shows the data structures when Bcmp is called to check if a match has been found.

211

and the second second

S

212-221

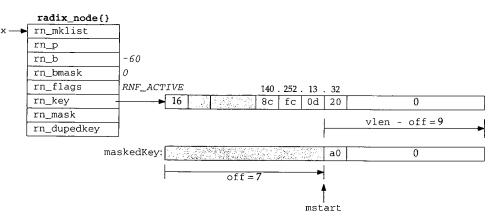


Figure 18.46 Comparison of maskedKey and new leaf.

Since the 9-byte strings are not the same, the comparison fails.

This while loop handles duplicate keys, each with a different mask. The only key of the duplicates that is compared is the one whose rn\_mask pointer equals m->rm\_mask. As an example, recall Figures 18.36 and 18.37. If the search starts at the node for bit 63, the first time through the inner do while loop m points to the radix\_mask structure for 0xffffff00. When rn\_search returns the pointer to the first of the duplicate leaves for 127.0.0.0, the rm\_mask of this leaf equals m->rm\_mask, so Bcmp is called. If the comparison fails, m is replaced with the pointer to the next radix\_mask structure on the list (the one with a mask of 0xff000000) and the do while loop iterates around again with the new mask. rn\_search again returns the pointer to the first of the duplicate leaves for 127.0.0.0, but its rn\_mask does not equal m->rm\_mask. The while steps to the next of the duplicate leaves and its rn\_mask is the right one.

Returning to our example with the search key of 140.252.13.188, since the search from the node that tests bit 62 failed, the backtracking continues up the tree until the top is reached, which is the next node up the tree with a nonnull rn\_mklist.

Figure 18.47 shows the data structures when the top node of the tree is reached. At this point maskedKey is computed (it is all zero bits) and rn\_search starts at this node (the top of the tree) and continues down the two left branches to the leaf labeled "default" in Figure 18.4.

When rn\_search returns, x points to the radix\_node with an rn\_b of -33, which is the first leaf encountered after the two left branches from the top of the tree. But x->rn\_mask (which is null) does not equal m->rm\_mask, so x is replaced with x->rn\_dupedkey. The test of the while loop occurs again, but now x->rn\_mask equals m->rm\_mask, so the while loop terminates. Bcmp compares the 12 bytes of 0 starting at mstart with the 12 bytes of 0 stating at x->rn\_key plus 4, and since they're equal, the function returns the pointer x, which points to the entry for the default route.

-

]

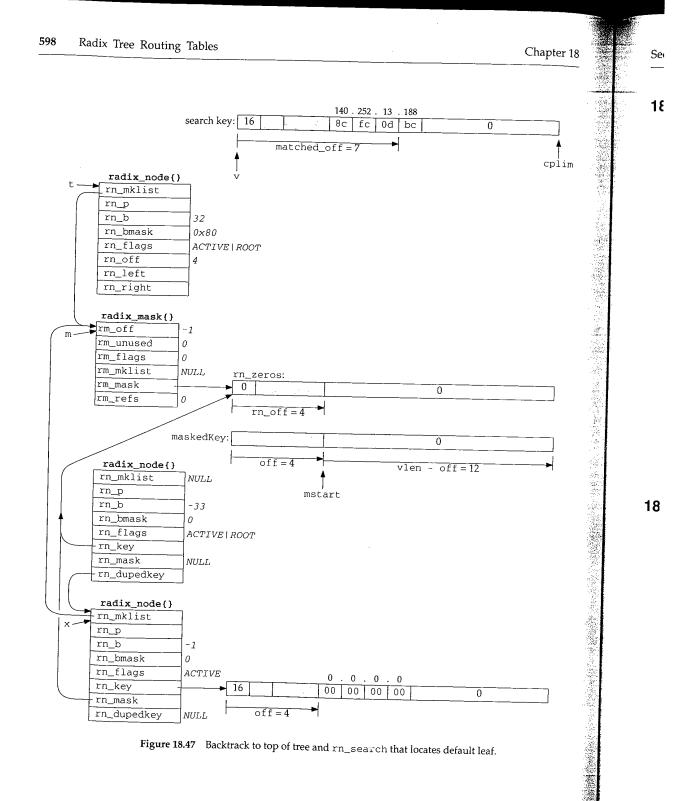
·ghe

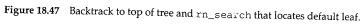
: 18

plim



INTEL Ex.1013.623





-----

Summary 599

radix.c

- radix.c

Section 18.12

### 18.11 rn\_search Function

rn\_search was called in the previous section from rn\_match to search a subtree of the routing table.

```
79 struct radix_node *
80 rn_search(v_arg, head)
81 void *v_arg;
82 struct radix_node *head;
83 (
84
       struct radix_node *x;
85
       caddr_t v;
86
       for (x = head, v = v_arg; x->rn_b >= 0;) {
87
          if (x->rn_bmask & v[x->rn_off])
                                   /* right if bit on */
88
               x = x - r;
89
           else
                                   /* left if bit off */
90
               x = x - rn_1;
91
       }
92
       return (x);
93 };
```

Figure 18.48 rn\_search function.

This loop is similar to the one in Figure 18.38. It compares one bit in the search key at each node, branching left if the bit is off or right if the bit is on, terminating when a leaf is encountered. The pointer to that leaf is returned.

### 18.12 Summary

Each routing table entry is identified by a key: the destination IP address in the case of the Internet protocols, which is either a host address or a network address with an associated network mask. Once the entry is located by searching for the key, additional information in the entry specifies the IP address of a router to which datagrams should be sent for the destination, a pointer to the interface to use, metrics, and so on.

The information maintained by the Internet protocols is the route structure, composed of just two elements: a pointer to a routing table entry and the destination address. We'll encounter one of these route structures in each of the Internet protocol control blocks used by UDP, TCP, and raw IP.

The Patricia tree data structure is well suited to routing tables. Routing table lookups occur much more frequently than adding or deleting routes, so from a performance standpoint using Patricia trees for the routing table makes sense. Patricia trees provide fast lookups at the expense of additional work in adding and deleting. Measurements in [Sklower 1991] comparing the radix tree approach to the Net/1 hash table show that the radix tree method is about two times faster in building a test tree and four times faster in searching.

# r 18

im

なるので、なるないないである。

学会部の言語を

1

## Exercises

- **18.1** We said with Figure 18.3 that the general condition for matching a routing table entry is that the search key logically ANDed with the routing table mask equal the routing table key. But in Figure 18.40 a different test is used. Build a logic truth table showing that the two tests are the same.
- **18.2** Assume a Net/3 system needs a routing table with 20,000 entries (IP addresses). Approximately how much memory is required for this, ignoring the space required for the masks?
- 18.3 What is the limit imposed on the length of a routing table key by the radix\_node structure?

# 19

at the ɔroxi-

try is table

:er 18

.sks? struc-

# Routing Requests and Routing Messages

### 19.1 Introduction

The various protocols within the kernel don't access the routing trees directly, using the functions from the previous chapter, but instead call a few functions that we describe in this chapter: rtalloc and rtalloc1 are two that perform routing table lookups, rtrequest adds and deletes routing table entries, and rtinit is called by most interfaces when the interface goes up or down.

Routing messages communicate information in two directions. A process such as the route command or one of the routing daemons (routed or gated) writes routing messages to a routing socket, causing the kernel to add a new route, delete an existing route, or modify an existing route. The kernel also generates routing messages that can be read by any routing socket when events occur in which the processes might be interested: an interface has gone down, a redirect has been received, and so on. In this chapter we cover the formats of these routing messages and the information contained therein, and we save our discussion of routing sockets until the next chapter.

Another interface provided by the kernel to the routing tables is through the sysctl system call, which we describe at the end of this chapter. This system call allows a process to read the entire routing table or a list of all the configured interfaces and interface addresses.

# 19.2 rtalloc and rtalloc1 Functions

rtalloc and rtalloc1 are the functions normally called to look up an entry in the routing table. Figure 19.1 shows rtalloc.

601

Routing Requests and Routing Messages 602

Chapter 19

route.c

60 struct route *ro; 61 { 62 if (ro->ro_rt && ro->ro_rt->rt_ifp && (ro->ro_rt->rt_flags & RTF_UP)) 63 return; /* XXX */	59	59	void rtalloc					route.c
64 ro->ro_rt = rtalloc1(&ro->ro_dst, 1);	61 62 63 64	61 62 63 64	{ if	(ro->ro_rt && return;	ro->ro_rt	->rt_ifp && /* XXX */	(ro->ro_rt->rt_flags	& RTF_UP))

Figure 19.1 rtalloc function.

The argument ro is often the pointer to a route structure contained in an Internet PCB (Chapter 22) which is used by UDP and TCP. If ro already points to an rtentry structure (ro\_rt is nonnull), and that structure points to an interface structure, and the route is up, the function returns. Otherwise rtalloc1 is called with a second argument of 1. We'll see the purpose of this argument shortly.

rtalloc1, shown in Figure 19.2, calls the rnh\_matchaddr function, which is always rn\_match (Figure 18.17) for Internet addresses. 66-76

The first argument is a pointer to a socket address structure containing the address to search for. The sa\_family member selects the routing table to search.

Call rn\_match 77-78

58--65

If the following three conditions are met, the search is successful.

1. A routing table exists for the protocol family,

2. rn\_match returns a nonnull pointer, and

3. the matching radix\_node does not have the RNF\_ROOT flag set.

Remember that the two leaves that mark the end of the tree both have the RNF\_ROOT flag set.

# Search fails

94-101

79

If the search fails because any one of the three conditions is not met, the statistic rts\_unreach is incremented and if the second argument to rtalloc1 (report) is nonzero, a routing message is generated that can be read by any interested processes on a routing socket. The routing message has the type RTM\_MISS, and the function returns

If all three of the conditions are met, the lookup succeeded and the pointer to the matching radix\_node is stored in rt and newrt. Notice that in the definition of the rtentry structure (Figure 18.24) the two radix\_node structures are at the beginning, and, as shown in Figure 18.8, the first of these two structures contains the leaf node. Therefore the pointer to a radix\_node structure returned by rn\_match is really a pointer to an rtentry structure, which is the matching leaf node.

er 19

ute.c

)

ute.c

rnet try . the rgu-

h is

ress

OOT

istic :) is s on ırns

the

the

ing,

ode.

ly a

66 str	cuct rtentry *	— route
	alloc1(dst, report)	
	ruct sockaddr *dst;	
69 int		
70 {		
71	struct radix_node_head *rnh = rt_tables[dst->sa_family];	
72	struct rtentry *rt;	
73	struct radix_node *rn;	
74	<pre>struct rtentry *newrt = 0;</pre>	
75	<pre>struct rt_addrinfo info;</pre>	
76	<pre>int s = splnet(), err = 0, msgtype = RTM_MISS;</pre>	
77	if (rnh && (rn = rnh->rnh_matchaddr((caddr_t) dst, rnh)) &&	
78	$((rn -> rn_flags \& RNF_ROOT) == 0))$ {	
79	newrt = rt = (struct rtentry *) rn;	
80	if (report && (rt->rt_flags & RTF_CLONING)) {	
81	err = rtrequest(RTM_RESOLVE, dst, SA(0),	
82	SA(0), 0, &newrt);	
83	if (err) {	
84	newrt = rt;	
85	rt->rt_refcnt++;	
86	goto miss;	
87	}	
88	if ((rt = newrt) && (rt->rt_flags & RTF_XRESOLVE)) {	
89	msgtype = RTM_RESOLVE;	
90	goto miss;	
91	}	
92	} else	
93	rt->rt_refcnt++;	
94	} else {	
95	<pre>rtstat.rts_unreach++;</pre>	
96	miss:if (report) {	
97	<pre>bzero((caddr_t) &amp; info, sizeof(info));</pre>	
98	info.rti_info[RTAX_DST] = dst;	
99	rt_missmsg(msgtype, &info, 0, err);	
100	}	
101	}	
102	<pre>splx(s);</pre>	
103	return (newrt);	
104 }		—— rout

### **Create clone entries**

If the caller specified a nonzero second argument, and if the RTF\_CLONING flag is 80-82 set, rtrequest is called with a command of RTM\_RESOLVE to create a new rtentry structure that is a clone of the one that was located. This feature is used by ARP and for multicast addresses.

### Clone creation fails

83-87 If rtrequest returns an error, newrt is set back to the entry returned by rn\_match and its reference count is incremented. A jump is made to miss where an RTM\_MISS message is generated.

### Check for external resolution

88-91 If rtrequest succeeds but the newly cloned entry has the RTF\_XRESOLVE flag set, a jump is made to miss, this time to generate an RTM\_RESOLVE message. The intent of this message is to notify a user process when the route is created, and it could be used with the conversion of IP addresses to X.121 addresses.

### Increment reference count for normal successful search

92–93 When the search succeeds but the RTF\_CLONING flag is not set, this statement increments the entry's reference count. This is the normal flow through the function, which then returns the nonnull pointer.

For a small function, rtalloc1 has many options in how it operates. There are seven different flows through the function, summarized in Figure 19.3.

	report argument	RTF CLONING flag	RTM RESOLVE return	RTF XRESOLVE flag	routing message generated	rt_refcnt	return value
	0						null
entry not found	1				RTM_MISS	++ ++ ++	null
		0				++	ptr
	0					++	ptr
entry found	1	1	OK	0		++	ptr
	1	1	OK	1	RTM_RESOLVE	++	ptr
	1	1	error		RTM_MISS	++	ptr

Figure 19.3 Summary of operation of rtalloc1.

We note that the first two rows (entry not found) are impossible if a default route exists. Also we show rt\_refent being incremented in the fifth and sixth rows when the call to rtrequest with a command of RTM\_RESOLVE is OK. The increment is done by rtrequest.

# 19.3 RTFREE Macro and rtfree Function

The RTFREE macro, shown in Figure 19.4, calls the rtfree function only if the reference count is less than or equal to 1, otherwise it just decrements the reference count.

209-213 The rtfree function, shown in Figure 19.5, releases an rtentry structure when there are no more references to it. We'll see in Figure 22.7, for example, that when a protocol control block is released, if it points to a routing entry, rtfree is called.

いたというなどないない

hapter 19 Section 19.3 RTFREE Macro and rtfree Function 605 209 #define RTFREE(rt) \ route.h rned by 210 if ((rt)->rt\_refcnt <= 1)  $\setminus$ vhere an 211 rtfree(rt); \ 212 else \ 213 (rt)->rt\_refcnt--; /\* no need for function call \*/ – route.h flag set, Figure 19.4 RTFREE macro. intent of be used 105 void - route.c 106 rtfree(rt) atement 107 struct rtentry \*rt; 108 { unction, 109 struct ifaddr \*ifa; 110 if (rt == 0) here are 111 panic("rtfree"); 112 rt->rt\_refcnt--; if (rt->rt\_refcnt <= 0 && (rt->rt\_flags & RTF\_UP) == 0) { 113 114if (rt->rt\_nodes->rn\_flags & (RNF\_ACTIVE | RNF\_ROOT)) return 115 panic("rtfree 2"); value 116 rttrash--: 117 if  $(rt - rt_refcnt < 0)$  { null printf("rtfree: %x not freed (neg refs)\n", rt); 118 null 119 return; 120 } ptr 121 ifa = rt->rt\_ifa; ptr 122 IFAFREE(ifa); ptr 123 Free(rt\_key(rt)); ptr 124 Free(rt); 125 } ptr 126 } - route.c Figure 19.5 rtfree function: release an rtentry structure. e exists. The entry's reference count is decremented and if it is less than or equal to 0 and the 105-115 the call route is not usable, the entry can be released. If either of the flags RNF\_ACTIVE or lone by RNF\_ROOT are set, this is an internal error. If RNF\_ACTIVE is set, this structure is still part of the routing table tree. If RNF\_ROOT is set, this structure is one of the end mark-ers built by rn\_inithead. rttrash is a debugging counter of the number of routing entries not in the routing 116 tree, but not released. It is incremented by rtrequest when it begins deleting a route, and then decremented here. Its value should normally be 0. le refer-Release interface reference unt. e when A check is made that the reference count is not negative, and then IFAFREE decre-117-122 ments the reference count for the ifaddr structure and releases it by calling ifafree when a when it reaches 0.

**INTEL Ex.1013.631** 

## **Release routing memory**

<sup>123-124</sup> The memory occupied by the routing entry key and its gateway is released. We'll see in rt\_setgate that the memory for both is allocated in one contiguous chunk, allowing both to be released with a single call to Free. Finally the rtentry structure itself is released.

# **Routing Table Reference Counts**

The handling of the routing table reference count, rt\_refcnt, differs from most other reference counts. We see in Figure 18.2 that most routes have a reference count of 0, yet the routing table entries without any references are not deleted. We just saw the reason in rtfree: an entry with a reference count of 0 is not deleted unless the entry's RTF\_UP flag is not set. The only time this flag is cleared is by rtrequest when a route is deleted from the routing tree.

Most routes are used in the following fashion.

• If the route is created automatically as a route to an interface when the interface is configured (which is typical for Ethernet interfaces, for example), then rtinit calls rtrequest with a command of RTM\_ADD, creating the new entry and setting the reference count to 1. rtinit then decrements the reference count to 0 before returning.

A point-to-point interface follows a similar procedure, so the route starts with a reference count of 0.

If the route is created manually by the route command or by a routing daemon, a similar procedure occurs, with route\_output calling rtrequest with a command of RTM\_ADD, setting the reference count to 1. This is then decremented by route\_output to 0 before it returns.

Therefore all newly created routes start with a reference count of 0.

• When an IP datagram is sent on a socket, be it TCP or UDP, we saw that ip\_output calls rtalloc, which calls rtalloc1. In Figure 19.3 we saw that the reference count is incremented by rtalloc1 if the route is found.

The located route is called a *held route*, since a pointer to the routing table entry is being held by the protocol, normally in a route structure contained within a protocol control block. An rtentry structure that is being held by someone else cannot be deleted, which is why rtfree doesn't release the structure until its reference count reaches 0.

• A protocol releases a held route by calling RTFREE or rtfree. We saw this in Figure 8.24 when ip\_output detects a change in the destination address. We'll encounter it in Chapter 22 when a protocol control block that holds a route is released.

Part of the confusion we'll encounter in the code that follows is that rtalloc1 is often called to look up a route in order to verify that a route to the destination exists, but

ter 19

We'll

hunk,

icture

other 0, yet eason

'F\_UP ute is

erface then entry erence

*w*ith a

emon, vith a

decre-

Section 19.4

when the caller doesn't want to hold the route. Since rtalloc1 increments the counter, the caller immediately decrements it.

Consider a route being deleted by rtrequest. The RTF\_UP flag is cleared, and if no one is holding the route (its reference count is 0), rtfree should be called. But rtfree considers it an error for the reference count to go below 0, so rtrequest checks whether its reference count is less than or equal to 0, and, if so, increments it and calls rtfree. Normally this sets the reference count to 1 and rtfree decrements it to 0 and deletes the route.

### 19.4 rtrequest Function

The rtrequest function is the focal point for adding and deleting routing table entries. Figure 19.6 shows some of the other functions that call it.

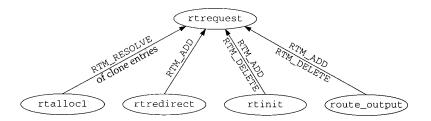


Figure 19.6 Summary of functions that call rtrequest.

rtrequest is a switch statement with one case per command: RTM\_ADD, RTM\_DELETE, and RTM\_RESOLVE. Figure 19.7 shows the start of the function and the RTM\_DELETE command.

		– route.c
v that	290 int	
w that	291 rtrequest(req, dst, gateway, netmask, flags, ret_nrt)	
W that	292 int req, flags;	
	293 struct sockaddr *dst, *gateway, *netmask;	
entry	294 struct rtentry **ret_nrt;	
ithin a	295 {	
	296 int $s = splnet();$	
neone	297 int error = 0;	
e until	298 struct rtentry *rt;	
	<pre>299 struct radix_node *rn;</pre>	
	<pre>300 struct radix_node_head *rnh;</pre>	
this in	<pre>301 struct ifaddr *ifa;</pre>	
We'll	<pre>302 struct sockaddr *ndst;</pre>	
oute is	<pre>303 #define senderr(x) { error = x ; goto bad; }</pre>	
	<pre>304 if ((rnh = rt_tables[dst-&gt;sa_family]) == 0)</pre>	
	305 senderr(ESRCH);	
oc1 is	306 if (flags & RTF_HOST)	
sts, but	307 netmask = 0;	

608 Routing Requests and Routing Messages

Chapter 19

route.c

308	switch (reg) {
309	case RTM_DELETE:
310	<pre>if ((rn = rnh-&gt;rnh_deladdr(dst, netmask, rnh)) == 0)</pre>
311	<pre>senderr(ESRCH);</pre>
312	if (rn->rn_flags & (RNF_ACTIVE   RNF_ROOT))
313	<pre>panic("rtrequest delete");</pre>
314	rt = (struct rtentry *) rn;
315	rt->rt_flags &= ~RTF_UP;
316	if (rt->rt_gwroute) {
317	<pre>rt = rt-&gt;rt_gwroute;</pre>
318	RTFREE(rt);
319	<pre>(rt = (struct rtentry *) rn)-&gt;rt_gwroute = 0;</pre>
320	}
321	if ((ifa = rt->rt_ifa) && ifa->ifa_rtrequest)
322	<pre>ifa-&gt;ifa_rtrequest(RTM_DELETE, rt, SA(0));</pre>
323	rttrash++;
324	if (ret_nrt)
325	<pre>*ret_nrt = rt;</pre>
326	else if (rt->rt_refcnt <= 0) {
327	rt->rt_refcnt++;
328	<pre>rtfree(rt);</pre>
329	}
330	break;

Figure 19.7 rtrequest function: RTM\_DELETE command.

290-307

The second argument, dst, is a socket address structure specifying the key to be added or deleted from the routing table. The sa\_family from this key selects the routing table. If the flags argument indicates a host route (instead of a route to a network), the netmask pointer is set to null, ignoring any value the caller may have passed.

# Delete from routing tree

309-315 The rnh\_deladdr function (rn\_delete from Figure 18.17) deletes the entry from the routing table tree and returns a pointer to the corresponding rtentry structure. The RTF\_UP flag is cleared.

# Remove reference to gateway routing table entry

316-320 If the entry is an indirect route through a gateway, RTFREE decrements the rt\_refcnt member of the gateway's entry and deletes it if the count reaches 0. The rt\_gwroute pointer is set to null and rt is set back to point to the entry that was deleted.

# Call interface request function

321--322 If an ifa\_rtrequest function is defined for this entry, that function is called. This function is used by ARP, for example, in Chapter 21 to delete the corresponding ARP entry.

# Return pointer or release reference

The rttrash global is incremented because the entry may not be released in the 323-330 code that follows. If the caller wants the pointer to the rtentry structure that was

e.c

be

ıt-

≥t-

ve

-m re.

he

he as

uis VP

he

as

deleted from the routing tree (if ret\_nrt is nonnull), then that pointer is returned, but the entry cannot be released: it is the caller's responsibility to call rtfree when it is finished with the entry. If ret\_nrt is null, the entry can be released: if the reference count is less than or equal to 0, it is incremented, and rtfree is called. The break causes the function to return.

Figure 19.8 shows the next part of the function, which handles the RTM\_RESOLVE command. This function is called with this command only from rtalloc1, when a new entry is to be created from an entry with the RTF\_CLONING flag set.

· · · · ·		route.c
331	case RTM_RESOLVE:	Tourca
332	if (ret_nrt == 0    (rt = *ret_nrt) == 0)	
333	senderr(EINVAL);	
334	ifa = rt->rt_ifa;	
335	flags = rt->rt_flags & ~RTF_CLONING;	
336	gateway = rt->rt_gateway;	
337	if ((netmask = rt->rt_genmask) == 0)	
338	<pre>flags  = RTF_HOST;</pre>	
339	goto makeroute;	
		——— route.c

Figure 19.8 rtrequest function: RTM\_RESOLVE command.

The final argument, ret\_nrt, is used differently for this command: it contains the pointer to the entry with the RTF\_CLONING flag set (Figure 19.2). The new entry will have the same rt\_ifa pointer, the same flags (with the RTF\_CLONING flag cleared), and the same rt\_gateway. If the entry being cloned has a null rt\_genmask pointer, the new entry has its RTF\_HOST flag set, because it is a host route; otherwise the new entry is a network route and the network mask of the new entry is copied from the rt\_genmask value. We give an example of cloned routes with a network mask at the end of this section. This case continues at the label makeroute, which is in the next figure.

Figure 19.9 shows the RTM\_ADD command.

### Locate corresponding interface

The function ifa\_ifwithroute finds the appropriate local interface for the destination (dst), returning a pointer to its ifaddr structure.

### Allocate memory for routing table entry

343-348 An rtentry structure is allocated. Recall that this structure contains both the two radix\_node structures for the routing tree and the other routing information. The structure is zeroed and the rt\_flags are set from the caller's flags, including the RTF\_UP flag.

### Allocate and copy gateway address

<sup>349-352</sup> The rt\_setgate function (Figure 19.11) allocates memory for both the routing table key (dst) and its gateway. It then copies gateway into the new memory and sets the pointers rt\_key, rt\_gateway, and rt\_gwroute.

# 610 Routing Requests and Routing Messages

Chapter 19

340	Case RTM_ADD:	— route.c
342	<pre>if ((ifa = ifa_ifwithroute(flags, dst, gateway)) == 0) senderr(ENETUNEFACU).</pre>	Toute.c
342	<pre>senderr(ENETUNREACH);</pre>	
343	marteroute.	
344	M_Hairoc(it, Struct rtentry * piperf(+ ))	
345	(20 0)	
346	Senderr(ENOBUES).	
347	Bzero(rt, sizeof(*rt)):	
348	rt->rt_flags = RTF UP   flags.	
349	if (rt_setgate(rt, dst, gateway)) {	
350	<pre>Free(rt);</pre>	
351	senderr(ENOBUFS);	
352	}	
353	<pre>ndst = rt_key(rt);</pre>	
354	if (netmask) (	
355	rt_maskedcopy(dst, ndst, netmask);	
356	} else	
357	<pre>Bcopy(dst, ndst, dst-&gt;sa_len);</pre>	
358	<pre>rn = rnh-&gt;rnh_addaddr((caddr_t) ndst, (caddr_t) netmask,</pre>	
359		
360	111 = 0	
361	if (rt->rt_gwroute)	
362	<pre>rtfree(rt-&gt;rt_gwroute);</pre>	
363	<pre>Free(rt_key(rt));</pre>	
364	Free(rt);	
365	senderr(EEXIST);	
366	}	
367	ifa->ifa_refcnt++;	
368	rt->rt_ifa = ifa;	
369	rt->rt_ifp = ifa->ifa_ifp;	
370	lf (req == RTM_RESOLVE)	
371	$rt - rt_rmx = (*ret nrt) - rt rms$	
372	(IId->IId_rtreguest)	
373	ifa->ifa_rtrequest(req. rtS)(ret_sub_s)	
374		
375	<pre>*ret_nrt = rt;</pre>	
376	rt->rt_refcnt++;	
377	}	
378	break;	-7
379	}	- - - - - - - - - - - - - - - - - - -
380	bad:	1. 
381	<pre>splx(s);</pre>	
382	return (error);	. 73 9
383 }		

Figure 19.9 rtrequest function: RTM\_ADD command.

# Copy destination address

353-357

The destination address (the routing table key dst) must now be copied into the memory pointed to by rn\_key. If a network mask is supplied, rt\_maskedcopy logically ANDs dst and netmask, forming the new key. Otherwise dst is copied into the

route.c

new key. The reason for logically ANDing dst and netmask is to guarantee that the key in the table has already been ANDed with its mask, so when a search key is compared against the key in the table only the search key needs to be ANDed. For example, the following command adds another IP address (an alias) to the Ethernet interface le0, with subnet 12 instead of 13:

### bsdi \$ ifconfig le0 inet 140.252.12.63 netmask 0xffffffe0 alias

The problem is that we've incorrectly specified all one bits for the host ID. Nevertheless, when the key is stored in the routing table we can verify with netstat that the address is first logically ANDed with the mask:

Destination	Gateway	Flags	Refs	Use	Interface
140.252.12.32	link#1	υC	0	0	le0

#### Add entry to routing tree

<sup>358-366</sup> The rnh\_addaddr function (rn\_addroute from Figure 18.17) adds this rtentry structure, with its destination and mask, to the routing table tree. If an error occurs, the structures are released and EEXIST returned (i.e., the entry is already in the routing table).

#### Store interface pointers

<sup>367-369</sup> The ifaddr structure's reference count is incremented and the pointers to its ifaddr and ifnet structures are stored.

## Copy metrics for newly cloned route

<sup>370–371</sup> If the command was RTM\_RESOLVE (not RTM\_ADD), the entire metrics structure is copied from the cloned entry into the new entry. If the command was RTM\_ADD, the caller can set the metrics after this function returns.

### Call interface request function

372-373 If an ifa\_rtrequest function is defined for this entry, that function is called. ARP uses this to perform additional processing for both the RTM\_ADD and RTM\_RESOLVE commands (Section 21.13).

### Return pointer and increment reference count

<sup>374–378</sup> If the caller wants a copy of the pointer to the new structure, it is returned through ret\_nrt and the rt\_reference count is incremented from 0 to 1.

### **Example: Cloned Routes with Network Masks**

The only use of the rt\_genmask value is with cloned routes created by the RTM\_RESOLVE command in rtrequest. If an rt\_genmask pointer is nonnull, then the socket address structure pointed to by this pointer becomes the network mask of the newly created route. In our routing table, Figure 18.2, the cloned routes are for the local Ethernet and for multicast addresses. The following example from [Sklower 1991] provides a different use of cloned routes. Another example is in Exercise 19.2.

Consider a class B network, say 128.1, that is behind a point-to-point link. The subnet mask is 0xffffff00, the typical value that uses 8 bits for the subnet ID and 8 bits

oute.c

o the logio the

r 19

for the host ID. We need a routing table entry for all possible 254 subnets, with a gateway value of a router that is directly connected to our host and that knows how to reach the link to which the 128.1 network is connected.

The easiest solution, assuming the gateway router isn't our default router, is a single entry with a destination of 128.1.0.0 and a mask of 0xffff0000. Assume, however, that the topology of the 128.1 network is such that each of the possible 254 subnets can have different operational characteristics: RTTs, MTUs, delays, and so on. If a separate routing table entry were used for each subnet, we would see that whenever a connection is closed, TCP would update the routing table entry with statistics about that route—its RTT, RTT variance, and so on (Figure 27.3). While we could create up to 254 entries by hand using the route command, one per subnet, a better solution is to use the cloning feature.

One entry is created by the system administrator with a destination of 128.1.0.0 and a network mask of 0xffff0000. Additionally, the RTF\_CLONING flag is set and the genmask is set to 0xffffff00, which differs from the network mask. If the routing table is searched for 128.1.2.3, and an entry does not exist for the 128.1.2 subnet, the entry for 128.1 with the mask of 0xffff0000 is the best match. A new entry is created (since the RTF\_CLONING flag is set) with a destination of 128.1.2 and a network mask of 0xfffff00 (the genmask value). The next time any host on this subnet is referenced, say 128.1.2.88, it will match this newly created entry.

### 19.5 rt\_setgate Function

Each leaf in the routing tree has a key (rt\_key, which is just the rn\_key member of the radix\_node structure contained at the beginning of the rtentry structure), and an associated gateway (rt\_gateway). Both are socket address structures specified when the routing table entry is created. Memory is allocated for both structures by rt\_setgate, as shown in Figure 19.10.

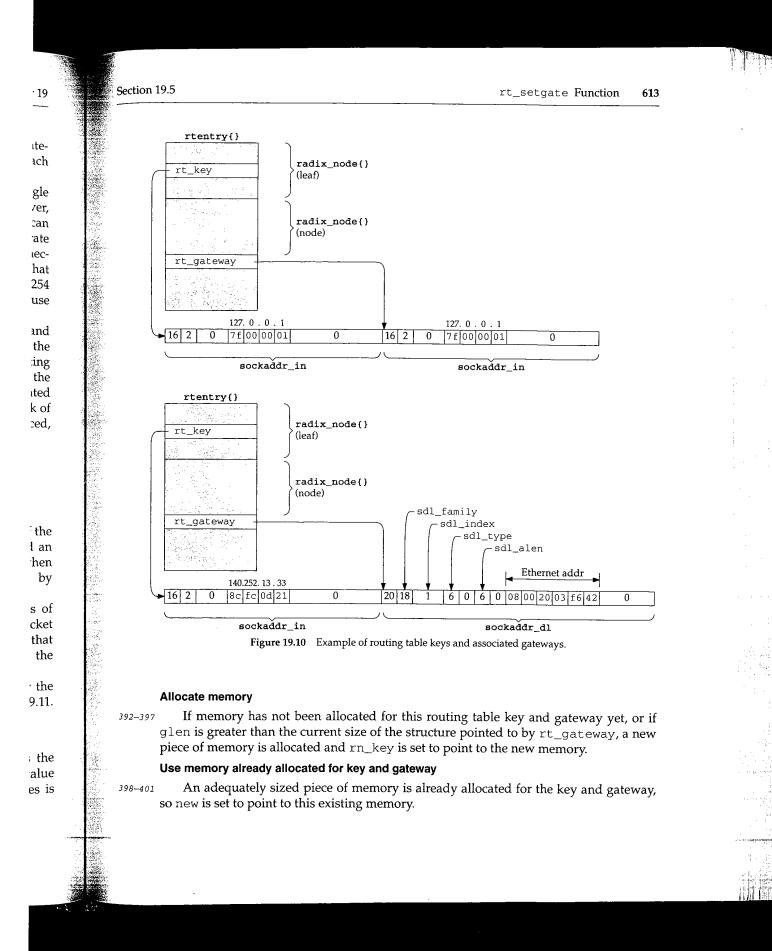
This example shows two of the entries from Figure 18.2, the ones with keys of 127.0.0.1 and 140.252.13.33. The former's gateway member points to an Internet socket address structure, while the latter's points to a data-link socket address structure that contains an Ethernet address. The former was entered into the routing table by the route system when the system was initialized, and the latter was created by ARP.

We purposely show the two structures pointed to by rt\_key one right after the other, since they are allocated together by rt\_setgate, which we show in Figure 19.11.

### Set lengths from socket address structures

384-391

dlen is the length of the destination socket address structure, and glen is the length of the gateway socket address structure. The ROUNDUP macro rounds the value up to the next multiple of 4 bytes, but the size of most socket address structures is already a multiple of 4.



INTEL Ex.1013.639

614 Routing Requests and Routing Messages

Chapter 19

route.c

```
384 int
                                                                               route.c
 385 rt_setgate(rt0, dst, gate)
 386 struct rtentry *rt0;
 387 struct sockaddr *dst, *gate;
 388 {
 389
         caddr_t new, old;
 390
                 dlen = ROUNDUP(dst->sa_len), glen = ROUNDUP(gate->sa_len);
         int
 391
         struct rtentry *rt = rt0;
         if (rt->rt_gateway == 0 ++ glen > ROUNDUP(rt->rt_gateway->sa_len)) {
 392
 393
             old = (caddr_t) rt_key(rt);
 394
             R_Malloc(new, caddr_t, dlen + glen);
 395
             if (new == 0)
396
                 return 1;
397
             rt->rt_nodes->rn_key = new;
398
         } else {
399
             new = rt->rt_nodes->rn_key;
400
             old \approx 0;
401
        }
402
        Bcopy(gate, (rt->rt_gateway = (struct sockaddr *) (new + dlen)), glen);
403
        if (old) {
404
             Bcopy(dst, new, dlen);
405
            Free(old);
406
        }
407
        if (rt->rt_gwroute) {
408
            rt = rt->rt_gwroute;
409
            RTFREE(rt);
410
            rt = rt0;
411
            rt->rt_gwroute = 0;
412
        }
413
        if (rt->rt_flags & RTF_GATEWAY) {
414
            rt->rt_gwroute = rtalloc1(gate, 1);
415
        }
416
        return 0;
417 }
```

Figure 19.11 rt\_setgate function.

### Copy new gateway

402

The new gateway structure is copied and rt\_gateway is set to point to the socket address structure.

# Copy key from old memory to new memory

<sup>403–406</sup> If a new piece of memory was allocated, the routing table key (dst) is copied right before the gateway field that was just copied. The old piece of memory is released.

# Release gateway routing pointer

407-412 If the routing table entry contains a nonnull rt\_gwroute pointer, that structure is released by RTFREE and the rt\_gwroute pointer is set to null.

452

453-459

45

Chapter 19

- route.c

:n);

n)) {

i, glen);

route.c

the socket

opied right

structure is

eased.

### Locate and store new gateway routing pointer

413-415 If the routing table entry is an indirect route, rtalloc1 locates the entry for the new gateway, which is stored in rt\_gwroute. If an invalid gateway is specified for an indirect route, an error is not returned by rt\_setgate, but the rt\_gwroute pointer will be null.

### 19.6 rtinit Function

There are four calls to rtinit from the Internet protocols to add or delete routes associated with interfaces.

- in\_control calls rtinit twice when the destination address of a point-topoint interface is set (Figure 6.21). The first call specifies RTM\_DELETE to delete any existing route to the destination; the second call specifies RTM\_ADD to add the new route.
- in\_ifinit calls rtinit to add a network route for a broadcast network or a host route for a point-to-point link (Figure 6.19). If the route is for an Ethernet interface, the RTF\_CLONING flag is automatically set by in\_ifinit.
- in\_ifscrub calls rtinit to delete an existing route for an interface.

Figure 19.12 shows the first part of the rtinit function. The cmd argument is always RTM\_ADD or RTM\_DELETE.

### Get destination address for route

If the route is to a host, the destination address is the other end of the point-to-point link. Otherwise we're dealing with a network route and the destination address is the unicast address of the interface (masked with ifa\_netmask).

### Mask network address with network mask

If a route is being deleted, the destination must be looked up in the routing table to locate its routing table entry. If the route being deleted is a network route and the interface has an associated network mask, an mbuf is allocated and the destination address is copied into the mbuf by rt\_maskedcopy, logically ANDing the caller's address with the mask. dst is set to point to the masked copy in the mbuf, and that is the destination looked up in the next step.

### Search for routing table entry

460-469 rtalloc1 searches the routing table for the destination address. If the entry is found, its reference count is decremented (since rtalloc1 incremented the reference count). If the pointer to the interface's ifaddr in the routing table does not equal the caller's argument, an error is returned.

#### Process request

470-473 rtrequest executes the command, either RTM\_ADD or RTM\_DELETE. When it returns, if an mbuf was allocated earlier, it is released.

### **INTEL Ex.1013.641**

route.c

4

- route.c

441 int 442 rtinit(ifa, cmd, flags) 443 struct ifaddr \*ifa; cmd, flags; 444 int 445 { struct rtentry \*rt; 446 struct sockaddr \*dst; 447 struct sockaddr \*deldst; 448 struct mbuf \*m = 0;449 struct rtentry \*nrt = 0; 450 error; 451 int dst = flags & RTF\_HOST ? ifa->ifa\_dstaddr : ifa->ifa\_addr; 452 if (cmd == RTM\_DELETE) { if ((flags & RTF\_HOST) == 0 && ifa->ifa\_netmask) { 453 454 m = m\_get(M\_WAIT, MT\_SONAME); deldst = mtod(m, struct sockaddr \*); 455 rt\_maskedcopy(dst, deldst, ifa->ifa\_netmask); 456 457 dst = deldst; 458 } 459 if (rt = rtalloc1(dst, 0)) { 460 rt->rt\_refcnt--; 461 if (rt->rt\_ifa != ifa) { 462 if (m) 463 (void) m\_free(m); return (flags & RTF\_HOST ? EHOSTUNREACH 464 465 : ENETUNREACH); 466 } 467 } 468 error = rtrequest(cmd, dst, ifa->ifa\_addr, ifa->ifa\_netmask, 469 flags | ifa->ifa\_flags, &nrt); 470 471 472 if (m) (void) m\_free(m); 473

Figure 19.12 rtinit function: call rtrequest to handle command.

Figure 19.13 shows the second half of rtinit.

# Generate routing message on successful delete

474-480

If a route was deleted, and rtrequest returned 0 along with a pointer to the rtentry structure that was deleted (in nrt), a routing socket message is generated by rt\_newaddrmsg. If the reference count is less than or equal to 0, it is incremented and the route is released by rtfree.

## Successful add

If a route was added, and rtrequest returned 0 along with a pointer to the rtentry structure that was added (in nrt), the reference count is decremented (since 481-482 rtrequest incremented it).

```
route.c
        if (cmd == RTM_DELETE && error == 0 && (rt = nrt)) {
474
475
            rt_newaddrmsg(cmd, ifa, error, nrt);
476
            if (rt->rt_refcnt <= 0) {
477
                rt->rt refcnt++;
478
                 rtfree(rt);
479
            }
480
        }
        if (cmd == RTM_ADD && error == 0 && (rt = nrt)) {
481
482
            rt->rt_refcnt--;
483
            if (rt->rt_ifa != ifa) (
                printf("rtinit: wrong ifa (%x) was (%x)\n", ifa,
484
485
                       rt->rt_ifa);
486
                if (rt->rt_ifa->ifa_rtrequest)
                    rt->rt_ifa->ifa_rtrequest(RTM_DELETE, rt, SA(0));
487
488
                IFAFREE(rt->rt_ifa);
489
                rt->rt_ifa = ifa;
490
                rt->rt_ifp = ifa->ifa_ifp;
491
                ifa->ifa_refcnt++;
492
                if (ifa->ifa_rtrequest)
493
                    ifa->ifa_rtrequest(RTM_ADD, rt, SA(0));
494
            }
495
            rt_newaddrmsg(cmd, ifa, error, nrt);
496
        }
497
        return (error);
498 }
                                                                              route.c
```

### Figure 19.13 rtinit function: second half.

### Incorrect interface

483-494 If the pointer to the interface's ifaddr in the new routing table entry does not equal the caller's argument, an error occurred. Recall that rtrequest determines the ifa pointer that is stored in the new entry by calling ifa\_ifwithroute (Figure 19.9). When this error occurs the following steps take place: an error message is output to the console, the ifa\_rtrequest function is called (if defined) with a command of RTM\_DELETE, the ifaddr structure is released, the rt\_ifa pointer is set to the value specified by the caller, the interface reference count is incremented, and the new interface's ifa\_rtrequest function (if defined) is called with a command of RTM\_ADD.

### Generate routing message

A routing socket message is generated by rt\_newaddrmsg for the RTM\_ADD command.

### 19.7 rtredirect Function

When an ICMP redirect is received, icmp\_input calls rtredirect and then calls pfctlinput (Figure 11.27). This latter function calls udp\_ctlinput and tcp\_ctlinput, which go through all the UDP and TCP protocol control blocks. If the

ıe 'y

d

495

.C

re

PCB is connected to the foreign address that has been redirected, and if the PCB holds a route to that foreign address, the route is released by rtfree. The next time any of these control blocks is used to send an IP datagram to that foreign address, rtalloc will be called and the destination will be looked up in the routing table, possibly finding a new (redirected) route.

The purpose of rtredirect, the first half of which is shown in Figure 19.14, is to validate the information in the redirect, update the routing table immediately, and then generate a routing socket message. – route.c

```
147 int
148 rtredirect(dst, gateway, netmask, flags, src, rtp)
149 struct sockaddr *dst, *gateway, *netmask, *src;
            flags;
150 int
151 struct rtentry **rtp;
152 {
        struct rtentry *rt;
153
                error = 0;
        int
154
        short *stat = 0;
155
        struct rt_addrinfo info;
156
        struct ifaddr *ifa;
157
         /* verify the gateway is directly reachable */
158
         if ((ifa = ifa_ifwithnet(gateway)) == 0) {
 159
             error = ENETUNREACH;
 160
             goto out;
 161
 162
         }
         rt = rtalloc1(dst, 0);
 163
          * If the redirect isn't from our current router for this dst,
         /*
 164
          * it's either old or wrong. If it redirects us to ourselves,
 165
          * we have a routing loop, perhaps as a result of an interface
 166
 167
           * going down recently.
 168
 170 #define equal(a1, a2) (bcmp((caddr_t)(a1), (caddr_t)(a2), (a1)->sa_len) == 0)
          */
         if (!(flags & RTF_DONE) && rt &&
 171
              (!equal(src, rt->rt_gateway) || rt->rt_ifa != ifa))
 172
              error = EINVAL;
 173
          else if (ifa_ifwithaddr(gateway))
 174
              error = EHOSTUNREACH;
 175
  176
          if (error)
              goto done;
  177
          /*
           * Create a new entry if we just got back a wildcard entry
  178
  179
           * or if the lookup failed. This is necessary for hosts
           * which use routing redirects generated by smart gateways
  180
  181
           * to dynamically build the routing tables.
  182
          if ((rt == 0) || (rt_mask(rt) && rt_mask(rt)->sa_len < 2))</pre>
           */
  183
  184
                                                                               - route.c
               goto create;
  185
                                                                                      1000
```

Figure 19.14 rtredirect function: validate received redirect.

3

holds a any of alloc finding

apter 19

nd then

– route.c

en) == 0)

– route.c

147-157 The arguments are dst, the destination IP address of the datagram that caused the redirect (HD in Figure 8.18); gateway, the IP address of the router to use as the new gateway field for the destination (R2 in Figure 8.18); netmask, which is a null pointer; flags, which is RTF\_GATEWAY and RTF\_HOST; src, the IP address of the router that sent the redirect (R1 in Figure 8.18); and rtp, which is a null pointer. We indicate that netmask and rtp are both null pointers when called by icmp\_input, but these arguments might be nonnull when called from other protocols.

### New gateway must be directly connected

<sup>158–162</sup> The new gateway must be directly connected or the redirect is invalid.

### Locate routing table entry for destination and validate redirect

163-177 rtalloc1 searches the routing table for a route to the destination. The following conditions must all be true, or the redirect is invalid and an error is returned. Notice that icmp\_input ignores any error return from rtredirect. ICMP does not generate an error in response to an invalid redirect—it just ignores it.

- the RTF\_DONE flag must not be set;
- rtalloc must have located a routing table entry for dst;
- the address of the router that sent the redirect (src) must equal the current rt\_gateway for the destination;
- the interface for the new gateway (the ifa returned by ifa\_ifwithnet) must equal the current interface for the destination (rt\_ifa), that is, the new gateway must be on the same network as the current gateway; and
- the new gateway cannot redirect this host to itself, that is, there cannot exist an attached interface with a unicast address or a broadcast address equal to gateway.

### Must create a new route

<sup>178–185</sup> If a route to the destination was not found, or if the routing table entry that was located is the default route, a new entry is created for the destination. As the comment indicates, a host with access to multiple routers can use this feature to learn of the correct router when the default is not correct. The test for finding the default route is whether the routing table entry has an associated mask and if the length field of the mask is less than 2, since the mask for the default route is rn\_zeros (Figure 18.35).

Figure 19.15 shows the second half of this function.

### Create new host route

186-195 If the current route to the destination is a network route and the redirect is a host redirect and not a network redirect, a new host route is created for the destination and the existing network route is left alone. We mentioned that the flags argument always specifies RTF\_HOST since the Net/3 ICMP considers all received redirects as host redirects. 620 Routing Requests and Routing Messages

Chapter 19

South States

route.c 186 /\* 187 \* Don't listen to the redirect if it's 188 \* for a route to an interface. 189 \*/ 190 if (rt->rt\_flags & RTF\_GATEWAY) { 191 if (((rt->rt\_flags & RTF\_HOST) == 0) && (flags & RTF\_HOST)) { 192 /\* 193 \* Changing from route to net => route to host. \* Create new route, rather than smashing route to net. 194 \*/ 195 196 create: 197 flags |= RTF\_GATEWAY | RTF\_DYNAMIC; 198 error = rtrequest((int) RTM\_ADD, dst, gateway, 199 netmask, flags, 200 (struct rtentry \*\*) 0); 201 stat = &rtstat.rts\_dynamic; 202 } else { 203 /\* \* Smash the current notion of the gateway to 204 205 \* this destination. Should check about netmask!!! \*/ 206 207 rt->rt\_flags i= RTF\_MODIFIED; 208 flags != RTF\_MODIFIED; 209 stat = &rtstat.rts\_newgateway; 210 rt\_setgate(rt, rt\_key(rt), gateway); 211 } 212 } else 213 error = EHOSTUNREACH; 214 done: 215 if (rt) { 216 if (rtp && !error) 217 \*rtp = rt;218 else 219 rtfree(rt); 220 } 221 out: 222 if (error) 223 rtstat.rts\_badredirect++; 224 else if (stat != NULL) 225 (\*stat)++; 226 bzero((caddr\_t) & info, sizeof(info)); 227 info.rti\_info[RTAX\_DST] = dst; 228 info.rti\_info[RTAX\_GATEWAY] = gateway; 229 info.rti\_info[RTAX\_NETMASK] = netmask; 230 info.rti\_info[RTAX\_AUTHOR] = src; 231 rt\_missmsg(RTM\_REDIRECT, &info, flags, error); 232 } - route.c

Figure 19.15 rtredirect function: second half.

INTEL Ex.1013.646

>ter 19

route c

route.c

Section 19.8

### Create route

196-201 rtrequest creates the new route, setting the RTF\_GATEWAY and RTF\_DYNAMIC flags. The netmask argument is a null pointer, since the new route is a host route with an implied mask of all one bits. stat points to a counter that is incremented later.

### Modify existing host route

202-211 This code is executed when the current route to the destination is already a host route. A new entry is not created, but the existing entry is modified. The RTF\_MODIFIED flag is set and rt\_setgate changes the rt\_gateway field of the routing table entry to the new gateway address.

### Ignore if destination is directly connected

212-213 If the current route to the destination is a direct route (the RTF\_GATEWAY flag is not set), it is a redirect for a destination that is already directly connected. EHOSTUNREACH is returned.

### **Return pointer and increment statistic**

If a routing table entry was located, it is either returned (if rtp is nonnull and there were no errors) or released by rtfree. The appropriate statistic is incremented.

### Generate routing message

226-232 An rt\_addrinfo structure is cleared and a routing socket message is generated by rt\_missmsg. This message is sent by raw\_input to any processes interested in the redirect.

### **19.8 Routing Message Structures**

Routing messages consist of a fixed-length header followed by up to eight socket address structures. The fixed-length header is one of the following three structures:

- rt\_msghdr
- if\_msghdr
- ifa\_msghdr

Figure 18.11 provided an overview of which functions generated the different messages and Figure 18.9 showed which structure is used by each message type. The first three members of the three structures have the same data type and meaning: the message length, version, and type. This allows the receiver of the message to decode the message. Also, each structure has a member that encodes which of the eight potential socket address structures follow the structure (a bitmask): the rtm\_addrs, ifm\_addrs, and ifam addrs members.

Figure 19.16 shows the most common of the structures, rt\_msghdr. The RTM\_IFINFO message uses an if\_msghdr structure, shown in Figure 19.17. The RTM\_NEWADDR and RTM\_DELADDR messages use an ifa\_msghdr structure, shown in Figure 19.18.

### **INTEL Ex.1013.647**

Chapter 19 Routing Requests and Routing Messages - route.h 139 struct rt\_msghdr { /\* to skip over non-understood messages \*/ u\_short rtm\_msglen; 140 /\* future binary compatibility \*/ u\_char rtm\_version; 141 /\* message type \*/ u\_char rtm\_type; 142 /\* index for associated ifp \*/ u\_short rtm\_index; 143 /\* flags, incl. kern & message, e.g. DONE \*/ int rtm\_flags; 144 /\* bitmask identifying sockaddrs in msg \*/ rtm\_addrs; int 145 /\* identify sender \*/ rtm\_pid; pid\_t 146 /\* for sender to identify action \*/ rtm\_seq; 147 int /\* why failed \*/ rtm\_errno; int 148 /\* from rtentry \*/ 149 int rtm\_use; /\* which metrics we are initializing \*/ u\_long rtm\_inits; 150 struct rt\_metrics rtm\_rmx; /\* metrics themselves \*/ 151 152 }; — route.h Figure 19.16 rt\_msghdr structure. if.h 235 struct if\_msghdr { /\* to skip over non-understood messages \*/ 236 u\_short ifm\_msglen; /\* future binary compatability \*/ u\_char ifm\_version; 237 /\* message type \*/ u\_char ifm\_type; 238 /\* like rtm\_addrs \*/ ifm\_addrs; 239 int /\* value of if\_flags \*/ 240 int ifm\_flags; /\* index for associated ifp \*/ u\_short ifm\_index; 241 /\* statistics and other data about if \*/ struct if\_data ifm\_data; 242 243 }; - if.h Figure 19.17 if\_msghdr structure. if.h 248 struct ifa\_msghdr { /\* to skip over non-understood messages \*/ 249 u\_short ifam\_msglen; /\* future binary compatability \*/ 250 u\_char ifam\_version; /\* message type \*/ 251 u\_char ifam\_type; /\* like rtm\_addrs \*/ 252 int ifam\_addrs; /\* value of ifa\_flags \*/ 253 int ifam\_flags; /\* index for associated ifp \*/ u\_short ifam\_index; 254 /\* value of ifa\_metric \*/ 255 int ifam\_metric; 256 }; - if.h

TP 1177

622

Figure 19.18 ifa\_msghdr structure.

Note that the first three members across the three different structures have the same data types and meanings.

The three variables rtm\_addrs, ifm\_addrs, and ifam\_addrs are bitmasks defining which socket address structures follow the header. Figure 19.19 shows the constants used with these bitmasks. Sec

R'I R'I

RI

RI

RΊ

Rl

R1

R'l

Bitmask Constant Value		Array index	<				
		Constant Valu		Name in rtsock.c	Description		
TA DST	0x01	RTAX DST	0	dst	destination socket address structure		
RTA GATEWAY	0x02	RTAX GATEWAY	1	gate	gateway socket address structure		
RTA_NETMASK	0x04	RTAX_NETMASK	2	netmask	netmask socket address structure		
RTA_GENMASK	0x08	RTAX_GENMASK	3	genmask	cloning mask socket address structure		
TA IFP	0x10	RTAX IFP	4	ifpaddr	interface name socket address structure		
RTA IFA	0x20	RTAX IFA	5	ifaaddr	interface address socket address structure		
RTA_AUTHOR	0x40	RTAX_AUTHOR	6		socket address structure for author of redirect		
RTA BRD	0x80	RTAX_BRD	7	brdaddr	broadcast or point-to-point destination address		
		RTAX MAX	8	1	#elements in an rti_info[] array		

Figure 19.19 Constants used to refer to members of rti\_info array.

The bitmask value is always the constant 1 left shifted by the number of bits specified by the array index. For example, 0x20 (RTA\_IFA) is 1 left shifted by five bits (RTAX\_IFA). We'll see this fact used in the code.

The socket address structures that are present always occur in order of increasing array index, one right after the other. For example, if the bitmask is 0x87, the first socket address structure contains the destination, followed by the gateway, followed by the network mask, followed by the broadcast address.

The array indexes in Figure 19.19 are used within the kernel to refer to its rt\_addrinfo structure, shown in Figure 19.20. This structure holds the same bitmask that we described, indicating which addresses are present, and pointers to those socket address structures.

									—— route.h
199 str	uct rt_a	addrinfo {							
200	int	rti_addr:	-,		same	as	rtm_addrs	*/	
201	struct	sockaddr	*rti_info[RTA	<_MAX];					
202 };									
									10410.11

Figure 19.20 rt\_addrinfo structure: encode which addresses are present and pointers to them.

For example, if the RTA\_GATEWAY bit is set in the rti\_addrs member, then the member rti\_info[RTAX\_GATEWAY] is a pointer to a socket address structure containing the gateway's address. In the case of the Internet protocols, the socket address structure is a sockaddr\_in containing the gateway's IP address.

The fifth column in Figure 19.19 shows the names used for the corresponding members of an rti\_info array throughout the file rtsock.c. These definitions look like

#define dst info.rti\_info[RTAX\_DST]

We'll encounter these names in many of the source files later in this chapter. The RTAX\_AUTHOR element is not assigned a name because it is never passed from a process to the kernel.

We've already encountered this rt\_addrinfo structure twice: in rtalloc1 (Figure 19.2) and rtredirect (Figure 19.14). Figure 19.21 shows the format of this

ter 19

\* /

NE \*/

route.h

```
— if.h
```

```
3 */
```

if.h

```
— if.h
```

—— if.h

same

definıstants

(1/3)

structure when built by rtalloc1, after a routing table lookup fails, when rt\_missmsg is called.

rt_addrinfo{}	_	
rti_addrs	0	sockaddr_in{}
rti_info[RTAX_DST] -	┣►	IP address that was not found
rti_info[RTAX_GATEWAY]		
rti_info[RTAX_NETMASK]		
rti_info[RTAX_GENMASK]	NULL	
rti_info[RTAX_IFP]	NULL	
rti_info[RTAX_IFA]	NULL	
rti_info[RTAX_AUTHOR]	NULL	
rti_info[RTAX_BRD]	NULL	

Figure 19.21 rt\_addrinfo structure passed by rtalloc1 to rt\_missmsg.

All the unused pointers are null because the structure is set to 0 before it is used. Also note that the rti\_addrs member is not initialized with the appropriate bitmask because when this structure is used within the kernel, a null pointer in the rti\_info array indicates a nonexistent socket address structure. The bitmask is needed only for messages between a process and the kernel.

Figure 19.22 shows the format of the structure built by rtredirect when it calls rt\_missmsg.

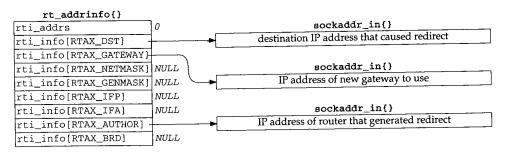


Figure 19.22 rt\_addrinfo structure passed by rtredirect to rt\_missmsg.

The following sections show how these structures are placed into the messages sent to a process.

Figure 19.23 shows the route\_cb structure, which we'll encounter in the following sections. It contains four counters; one each for the IP, XNS, and OSI protocols, and an "any" counter. Each counter is the number of routing sockets currently in existence for that domain.

203-208

By keeping track of the number of routing socket listeners, the kernel avoids building a routing message and calling raw\_input to send the message when there aren't any processes waiting for a message.

高麗市が法し

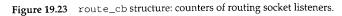
1

Å.

– rtsock.c

rtsock.c

203 st	ruct rou	ute_cb {	-						— route.l
204 205 206	int int int	ip_count; ns_count; iso_count;	/*	IP */ XNS * ISO *					
200 207 208 };	int	any_count;	/*	sum o	f above	three	counters	*/	



# 19.9 rt\_missmsg Function

The function rt\_missmsg, shown in Figure 19.24, takes the structures shown in Figures 19.21 and 19.22, calls rt\_msg1 to build a corresponding variable-length message for a process in an mbuf chain, and then calls raw\_input to pass the mbuf chain to all appropriate routing sockets.

516 void 517 rt\_missmsg(type, rtinfo, flags, error) type, flags, error; 518 int 519 struct rt\_addrinfo \*rtinfo; 520 { struct rt\_msghdr \*rtm; 521 struct mbuf \*m; 522 struct sockaddr \*sa = rtinfo->rti\_info[RTAX\_DST]; 523 if (route\_cb.any\_count == 0) 524 525 return; m = rt\_msg1(type, rtinfo); 526 if (m == 0) 527 return; 528 rtm = mtod(m, struct rt\_msghdr \*); 529 530 rtm->rtm\_flags = RTF\_DONE | flags; rtm->rtm\_errno = error; 531 rtm->rtm\_addrs = rtinfo->rti\_addrs; 532 route\_proto.sp\_protocol = sa ? sa->sa\_family : 0; 533 raw\_input(m, &route\_proto, &route\_src, &route\_dst); 534 535 }

Figure 19.24 rt\_missmsg function.

516–525 If there aren't any routing socket listeners, the function returns immediately.

# Build message in mbuf chain

526-528 rt\_msg1 (Section 19.12) builds the appropriate message in an mbuf chain, and returns the pointer to the chain. Figure 19.25 shows an example of the resulting mbuf chain, using the rt\_addrinfo structure from Figure 19.22. The information needs to be in an mbuf chain because raw\_input calls sbappendaddr to append the mbuf chain to a socket's receive buffer.

o k or

ιs



١g

m

626 Routing Requests and Routing Messages

Chapter 19

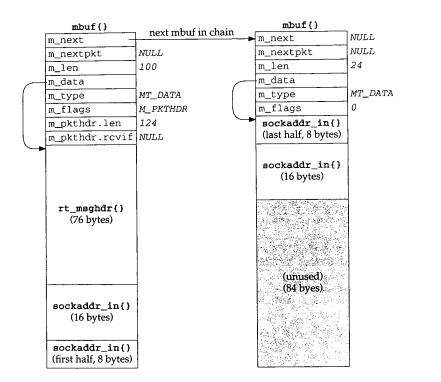


Figure 19.25 Mbuf chain built by rt\_msg1 corresponding to Figure 19.22.

### Finish building message

529-532 The two members rtm\_flags and rtm\_errno are set to the values passed by the caller. The rtm\_addrs member is copied from the rti\_addrs value. We showed this value as 0 in Figures 19.21 and 19.22, but rt\_msg1 calculates and stores the appropriate bitmask, based on which pointers in the rti\_info array are nonnull.

### Set protocol of message, call raw\_input

533-534 The final three arguments to raw\_input specify the protocol, source, and destination of the routing message. These three structures are initialized as

> struct sockaddr route\_dst = { 2, PF\_ROUTE, }; struct sockaddr route\_src = { 2, PF\_ROUTE, }; struct sockproto route\_proto = { PF\_ROUTE, };

The first two structures are never modified by the kernel. The sockproto structure, shown in Figure 19.26, is one we haven't seen before.

```
128 struct sockproto {

129 u_short sp_family; /* address family */

130 u_short sp_protocol; /* protocol */

131 };
```

Figure 19.26 sockproto structure.

INTEL Ex.1013.652

socket.h

socket.h

The family is never changed from its initial value of PF\_ROUTE, but the protocol is set each time raw\_input is called. When a process creates a routing socket by calling socket, the third argument (the protocol) specifies the protocol in which the process is interested. The caller of raw\_input sets the sp\_protocol member of the route\_proto structure to the protocol of the routing message. In the case of rt\_missmsg, it is set to the sa\_family of the destination socket address structure (if specified by the caller), which in Figures 19.21 and 19.22 would be AF\_INET.

### 19.10 rt\_ifmsg Function

In Figure 4.30 we saw that if\_up and if\_down both call rt\_ifmsg, shown in Figure 19.27, to generate a routing socket message when an interface goes up or down.

```
    rtsock.c

540 void
541 rt_ifmsg(ifp)
542 struct ifnet *ifp;
543 {
        struct if_msghdr *ifm;
544
545
        struct mbuf *m;
546
        struct rt_addrinfo info;
547
        if (route_cb.any_count == 0)
548
            return;
        bzero((caddr_t) & info, sizeof(info));
549
        m = rt_msgl(RTM_IFINFO, &info);
550
        if (m == 0)
551
552
            return;
        ifm = mtod(m, struct if_msghdr *);
553
        ifm->ifm_index = ifp->if_index;
554
555
        ifm->ifm_flags = ifp->if_flags;
        ifm->ifm_data = ifp->if_data;
                                          /* structure assignment */
556
557
        ifm->ifm_addrs = 0;
        route_proto.sp_protocol = 0;
558
        raw_input(m, &route_proto, &route_src, &route_dst);
559
560 }
                                                                               rtsock.c
```

cture,

y the

d this

oriate

stina-

ocket.h

ocket.h

Figure 19.27 rt\_ifmsg function.

547-548 If there aren't any routing socket listeners, the function returns immediately.

### Build message in mbuf chain

549-552 An rt\_addrinfo structure is set to 0 and rt\_msg1 builds an appropriate message in an mbuf chain. Notice that all socket address pointers in the rt\_addrinfo structure are null, so only the fixed-length if\_msghdr structure becomes the routing message; there are no addresses.

### Finish building message

553-557 The interface's index, flags, and if\_data structure are copied into the message in the mbuf and the ifm\_addrs bitmask is set to 0.

### Set protocol of message, call raw\_input

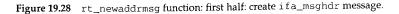
The protocol of the routing message is set to 0 because this message can apply to all protocol suites. It is a message about an interface, not about some specific destination. raw\_input delivers the message to the appropriate listeners.

# 19.11 rt\_newaddrmsg Function

In Figure 19.13 we saw that rtinit calls rt\_newaddrmsg with a command of RTM\_ADD or RTM\_DELETE when an interface has an address added or deleted. Figure 19.28 shows the first half of the function.

```
- rtsock.c
```

```
569 void
570 rt_newaddrmsg(cmd, ifa, error, rt)
571 int
           cmd, error;
572 struct ifaddr *ifa;
573 struct rtentry *rt;
574 {
        struct rt_addrinfo info;
575
        struct sockaddr *sa;
576
577
        int
                pass;
578
        struct mbuf *m;
        struct ifnet *ifp = ifa->ifa_ifp;
579
        if (route_cb.any_count == 0)
580
            return:
581
582
        for (pass = 1; pass < 3; pass++) {</pre>
            bzero((caddr_t) & info, sizeof(info));
583
            if ((cmd == RTM_ADD && pass == 1) ||
584
                 (cmd == RTM_DELETE && pass == 2)) {
585
                 struct ifa_msghdr *ifam;
586
                        ncmd = cmd == RTM_ADD ? RTM_NEWADDR : RTM_DELADDR;
587
                 int
                 ifaaddr = sa = ifa->ifa_addr;
588
                 ifpaddr = ifp->if_addrlist->ifa_addr;
589
                 netmask = ifa->ifa_netmask;
590
591
                 brdaddr = ifa->ifa_dstaddr;
                 if ((m = rt_msg1(ncmd, &info)) == NULL)
592
593
                     continue;
                 ifam = mtod(m, struct ifa_msghdr *);
594
                 ifam->ifam_index = ifp->if_index;
595
                 ifam->ifam_metric = ifa->ifa_metric;
596
                 ifam->ifam_flags = ifa->ifa_flags;
597
598
                 ifam->ifam_addrs = info.rti_addrs;
599
             }
                                                                              rtsock.c
```



hapter 19

Section 19.11

580-581

582

583

×.

essage in

ply to all stination.

mand of :ed. Fig-

— rtsock.c

ıR;

If there aren't any routing socket listeners, the function returns immediately.

### Generate two routing messages

The for loop iterates twice because two messages are generated. If the command is RTM\_ADD, the first message is of type RTM\_NEWADDR and the second message is of type RTM\_DELETE, the first message is of type RTM\_DELETE and the second message is of type RTM\_DELADDR. The RTM\_NEWADDR and RTM\_DELADDR messages are built from an ifa\_msghdr structure, while the RTM\_ADD and RTM\_DELETE messages are built from an rt\_msghdr structure. The function generates two messages because one message provides information about the interface and the other about the addresses.

An rt\_addrinfo structure is set to 0.

### Generate message with up to four addresses

<sup>588-591</sup> Pointers to four socket address structures containing information about the interface address that has been added or deleted are stored in the rti\_info array. Recall from Figure 19.19 that ifaaddr, ifpaddr, netmask, and brdaddr reference elements in the rti\_info array in info. rt\_msg1 builds the appropriate message in an mbuf chain. Notice that sa is set to point to the ifa\_addr structure, and we'll see at the end of the function that the family of this socket address structure becomes the protocol of the routing message.

Remaining members of the ifa\_msghdr structure are filled in with the interface's index, metric, and flags, along with the bitmask set by rt\_msg1.

Figure 19.29 shows the second half of rt\_newaddrmsg, which creates an rt\_msghdr message with information about the routing table entry that was added or deleted.

### **Build message**

Pointers to three socket address structures are stored in the rti\_info array: the rt\_mask, rt\_key, and rt\_gateway structures. sa is set to point to the destination address, and its family becomes the protocol of the routing message. rt\_msg1 builds the appropriate message in an mbuf chain.

Additional fields in the rt\_msghdr structure are filled in, including the bitmask set by rt\_msg1.

### Set protocol of message, call raw\_input

<sup>616–619</sup> The protocol of the routing message is set and raw\_input passes the message to the appropriate listeners. The function returns after two iterations through the loop.

— rtsock.c

INTEL Ex.1013.655

630 Routing Requests and Routing Messages

Chapter 19

rtsock.c

– rtsock.c

Sec

600	if ((cmd == RTM_ADD && pass == 2)
601	$(cmd == RTM_DELETE \&\& pass == 1))$ {
602	struct rt_msghdr *rtm;
603	if (rt == 0)
604	continue;
605	<pre>netmask = rt_mask(rt);</pre>
606	$dst = sa = rt_key(rt);$
607	<pre>gate = rt-&gt;rt_gateway;</pre>
608	if ((m = rt_msgl(cmd, &info)) == NULL)
609	continue;
610	rtm = mtod(m, struct rt_msghdr *);
611	rtm->rtm_index = ifp->if_index;
612	rtm->rtm_flags  = rt->rt_flags;
613	rtm->rtm_errno = error;
614	rtm->rtm_addrs = info.rti_addrs;
615	}
616	route proto sp protocol = sa ? sa->sa_family : 0;
617	raw_input(m, &route_proto, &route_src, &route_dst);
618	
618	1

Figure 19.29 rt\_newaddrmsg function: second half, create rt\_msghdr message.

# 19.12 rt\_msg1 Function

The functions described in the previous three sections each called rt\_msg1 to build the appropriate routing message. In Figure 19.25 we showed the mbuf chain that was built by rt\_msg1 from the rt\_msghdr and rt\_addrinfo structures in Figure 19.22. Figure 19.30 shows the function.

# Get mbuf and determine fixed size of message

An mbuf with a packet header is obtained and the length of the fixed-size message is stored in len. Two of the message types in Figure 18.9 use an ifa\_msghdr structure, one uses an if\_msghdr structure, and the remaining nine use an rt\_msghdr structure.

## Verify structure fits in mbuf

The size of the fixed-length structure must fit entirely within the data portion of the packet header mbuf, because the mbuf pointer is cast to a structure pointer using mtod and the structure is then referenced through the pointer. The largest of the three structures is if\_msghdr, which at 84 bytes is less than MHLEN (100).

# Initialize mbuf packet header and zero structure

<sup>425–428</sup> The two fields in the packet header are initialized and the structure in the mbuf is set to 0.

ock.c

;ock.c

i the built Fig-

sage ture, ture.

**√f the** ntod **∶truc-**

ouf is

r 19

a file			rtsock.c
		atic struct mbuf *	
	400 rt_msg1(type, rtinfo)		
語) (記: ///	401 in		
1. 12		ruct rt_addrinfo *rtinfo;	
	403 {		
	404	struct rt_msghdr *rtm;	
	405		
5	406	int i;	
z ( -	407	struct sockaddr *sa;	
	408	int len, dlen;	
	409	<pre>m = m_gethdr(M_DONTWAIT, MT_DATA);</pre>	
ŝ	410	if $(m == 0)$	
ξ. · .	411	return (m);	
	412	switch (type) {	
5/ -	412	Switcom (0) 50) (	
	413	case RTM_DELADDR:	
	414	case RTM_NEWADDR:	
4	415	<pre>len = sizeof(struct ifa_msghdr);</pre>	
	416	break;	
	417	case RTM_IFINFO:	
	418	len = sizeof(struct if_msghdr);	
	410	break;	
	419	Dreak,	
	420	default:	
	421	<pre>len = sizeof(struct rt_msghdr);</pre>	
	422	}	
	423	if (len > MHLEN)	
	424	<pre>panic("rt_msgl");</pre>	
	425	<pre>m-&gt;m_pkthdr.len = m-&gt;m_len = len;</pre>	
	426	m->m_pkthdr.rcvif = 0;	
	427	<pre>rtm = mtod(m, struct rt_msghdr *);</pre>	
	428	<pre>bzero((caddr_t) rtm, len);</pre>	
	429	for $(i = 0; i < RTAX_MAX; i++)$ {	
	430	if ((sa = rtinfo->rti_info[i]) == NULL)	
	431	continue;	
	432	<pre>rtinfo-&gt;rti_addrs  = (1 &lt;&lt; i);</pre>	
	433	<pre>dlen = ROUNDUP(sa-&gt;sa_len);</pre>	
	434	m_copyback(m, len, dlen, (caddr_t) sa);	
	435	len += dlen;	
	436	}	
	437	if (m->m_pkthdr.len != len) {	
	438	m_freem(m);	
	439	return (NULL);	
	440	}	
	441	rtm->rtm_msglen = len;	
	442	rtm->rtm_version = RTM_VERSION;	
	443	rtm->rtm_type = type;	
	444	return (m);	
	445 }		

Figure 19.30 rt\_msg1 function: obtain and initialize mbuf.

# Copy socket address structures into mbuf chain

- The caller passes a pointer to an rt\_addrinfo structure. The socket address structures corresponding to all the nonnull pointers in the rti\_info are copied into the mbuf by m\_copyback. The value 1 is left shifted by the RTAX\_xxx index to generate the corresponding RTA\_xxx bitmask (Figure 19.19), and each individual bitmask is logically ORed into the rti\_addrs member, which the caller can store on return into the corresponding member of the message structure. The ROUNDUP macro rounds the size of each socket address structure up to the next multiple of 4 bytes.
- 437-440 If, when the loop terminates, the length in the mbuf packet header does not equal len, the function m\_copyback wasn't able to obtain a required mbuf.

# Store length, version, and type

<sup>441-445</sup> The length, version, and message type are stored in the first three members of the message structure. Again, all three *xxx\_msghdr* structures start with the same three members, so this code works with all three structures even though the pointer rtm is a pointer to an rt\_msghdr structure.

### 19.13 rt\_msg2 Function

rt\_msgl constructs a routing message in an mbuf chain, and the three functions that called it then called raw\_input to append the mbuf chain to one or more socket's receive buffer. rt\_msg2 is different—it builds a routing message in a memory buffer, not an mbuf chain, and has as an argument a pointer to a walkarg structure that is used when rt\_msg2 is called by the two functions that handle the sysctl system call for the routing domain. rt\_msg2 is called in two different scenarios:

- 1. from route\_output to process the RTM\_GET command, and
- from sysctl\_dumpentry and sysctl\_iflist to process a sysctl system call.

Before looking at rt\_msg2, Figure 19.31 shows the walkarg structure that is used in scenario 2. We go through all these members as we encounter them.

42 43 44 45	<pre>truct walkarg {     int w_op;     int w_arg;     int w_given;     int w_needed;     int w_tmemsize;</pre>	<pre>/* NET_RT_xxx */ /* RTF_xxx for FLAGS, if_index for IFLIST */ /* size of process' buffer */ /* #bytes actually needed (at end) */ /* size of buffer pointed to by w_tmem */</pre>
46 47	int w_tmemsize; caddr_t w_where;	<pre>/* size of buffer pointed to by w_tmem */ /* ptr to process' buffer (maybe null) */ /* ptr to our malloc'ed buffer */</pre>
48	caddr_t w_tmem;	
49 }	;	rtsock.c

Figure 19.31 walkarg structure: used with the sysctl system call in the routing domain.

Figure 19.32 shows the first half of the rt\_msg2 function. This portion is similar to the first half of rt\_msg1.

rtsock.c

Section 19.13

rt msg2	Function	633
---------	----------	-----

ructhe rate logi-) the size qual f the hree i is a that ket's lfer, iat is ۱ call stem used sock.c ST \*/ \* / \* /

tsock.c

lar to

rtsock.c 446 static int 447 rt\_msg2(type, rtinfo, cp, w) 448 int type; 449 struct rt\_addrinfo \*rtinfo; 450 caddr\_t cp; 451 struct walkarg \*w; 452 { 453 int i; 454 int len, dlen, second\_time = 0; 455 caddr\_t cp0; 456 rtinfo->rti\_addrs = 0; 457 again: switch (type) { 458 459 case RTM\_DELADDR: 460 case RTM\_NEWADDR: 461 len = sizeof(struct ifa\_msghdr); 462 break; case RTM\_IFINFO: 463 464 len = sizeof(struct if\_msghdr); 465 break; default: 466 467 len = sizeof(struct rt\_msghdr); 468 3 469 if (cp0 = cp)cp += len; 470 for (i = 0; i < RTAX\_MAX; i++) { 471 struct sockaddr \*sa; 472 473 if ((sa = rtinfo->rti\_info[i]) == 0) continue: 474 rtinfo->rti\_addrs |= (1 << i);</pre> 475 476 dlen = ROUNDUP(sa->sa\_len); 477 if (cp) { 478 bcopy((caddr\_t) sa, cp, (unsigned) dlen); 479 cp += dlen; 480 3 len += dlen; 481 482 } rtsock.c

Figure 19.32 rt\_msg2 function: copy socket address structures.

Since this function stores the resulting message in a memory buffer, the caller specifies the start of that buffer in the cp argument. It is the caller's responsibility to ensure that the buffer is large enough for the message that is generated. To help the caller determine this size, if the cp argument is null, rt\_msg2 doesn't store anything but processes the input and returns the total number of bytes required to hold the result. We'll see that route\_output uses this feature and calls this function twice: first to determine the size and then to store the result, after allocating a buffer of the correct size. When rt\_msg2 is called by route\_output, the final argument is null. This final argument is nonnull when called as part of the sysct1 system call processing.

er 19

Chapter 19

rtsock.c

Se

### Determine size of structure

The size of the fixed-length message structure is set based on the message type. If the cp pointer is nonnull, it is incremented by this size.

## Copy socket address structures

471-482

The for loop goes through the rti\_info array, and for each element that is a nonnull pointer it sets the appropriate bit in the rti\_addrs bitmask, copies the socket address structure (if cp is nonnull), and updates the length.

Figure 19.33 shows the second half of rt\_msg2, most of which handles the optional walkarg structure.

483	if (cp == 0 && w != NULL && !second_time) {	
484	<pre>struct walkarg *rw = w;</pre>	
485	rw->w_needed += len;	
486	if (rw->w_needed <= 0 && rw->w_where) {	
487	if (rw->w_tmemsize < len) {	
488	if (rw->w_tmem)	
489	<pre>free(rw-&gt;w_tmem, M_RTABLE);</pre>	
490	if (rw->w_tmem = (caddr_t)	
491	malloc(len, M_RTABLE, M_NOWAIT))	
492	rw->w_tmemsize = len;	
493	}	
494	if (rw->w_tmem) {	
495	cp = rw->w_tmem;	
496	second_time = 1;	
497	goto again;	
498	} else	
499	rw->w_where = 0;	
500	}	
501	}	
502	if (cp) {	
503	struct rt_msghdr *rtm = (struct rt_msghdr *) cp0;	
504	rtm->rtm_version = RTM_VERSION;	
505	rtm->rtm_type = type;	
506	rtm->rtm_msglen = len;	
507	}	
508	return (len);	
509 }		rtsock.c

Figure 19.33 rt\_msg2 function: handle optional walkarg argument.

This if statement is true only when a pointer to a walkarg structure was passed and this is the first loop through the function. The variable second\_time was initialized to 0 but can be set to 1 within this if statement, and a jump made back to the label again in Figure 19.32. The test for cp being a null pointer is superfluous since whenever the w pointer is nonnull, the cp pointer is null, and vice versa.

## Check if data to be stored

485-486 w\_needed is incremented by the size of the message. This variable is initialized to
 0 minus the size of the user's buffer to the sysctl function. For example, if the buffer

:**r** 19

2. If

10ncket

onal

sock.c

の時代であるのない

size is 500 bytes, w\_needed is initialized to -500. As long as it remains negative, there is room in the buffer. w\_where is a pointer to the buffer in the calling process. It is null if the process doesn't want the result—the process just wants sysctl to return the size of the result, so the process can allocate a buffer and call sysctl again. rt\_msg2 doesn't copy the data back to the process—that is up to the caller—but if the w\_where pointer is null, there's no need for rt\_msg2 to malloc a buffer to hold the result and loop back through the function again, storing the result in this buffer. There are really five different scenarios that this function handles, summarized in Figure 19.34.

called from	ср	W	w.w_where	second_time	Description
	null	null			wants return length
route_output	nonnull	null	· · · · · · · · · · · · · · · · · · ·		wants result
sysctl_rtable	null	nonnull	null	0	process wants return length
	null	nonnull	nonnull	0	first time around to calculate length
	nonnull	nonnull	nonnull	1	second time around to store result

Figure 19.34 Summary of different scenarios for rt\_msg2.

### Allocate buffer first time or if message length increases

487-493 w\_tmemsize is the size of the buffer pointed to by w\_tmem. It is initialized to 0 by sysctl\_rtable, so the first time rt\_msg2 is called for a given sysctl request, the buffer must be allocated. Also, if the size of the result increases, the existing buffer must be released and a new (larger) buffer allocated.

### Go around again and store result

If w\_tmem is nonnull, a buffer already exists or one was just allocated. cp is set to point to this buffer, second\_time is set to 1, and a jump is made to again. The if statement at the beginning of this figure won't be true during this second pass, since second\_time is now 1. If w\_tmem is null, the call to malloc failed, so the pointer to the buffer in the process is set to null, preventing anything from being returned.

### Store length, version, and type

502-509 If cp is nonnull, the first three elements of the message header are stored. The function returns the length of the message.

19.14 sysctl\_rtable Function

This function handles the sysctl system call on a routing socket. It is called by net\_sysctl as shown in Figure 18.11.

Before going through the source code, Figure 19.35 shows the typical use of this system call with respect to the routing table. This example is from the arp program.

The first three elements in the mib array cause the kernel to call sysctl\_rtable to process the remaining elements.

rtsock.c

passed nitiale label when-

zed to buffer

ant

636 Routing Requests and Routing Messages

Chapter 19

```
mib[6];
int
size_t needed;
        *buf, *lim, *next;
char
struct rt_msghdr *rtm;
mib[0] = CTL_NET;
mib[1] = PF_ROUTE;
mib[2] = 0;
                        /* address family; can be 0 */
mib[3] = AF_INET;
mib[4] = NET_RT_FLAGS; /* operation */
                        /* flags; can be 0 */
mib[5] = RTF_LLINFO;
if (sysctl(mib, 6, NULL, &needed, NULL, 0) < 0)
    quit("sysctl error, estimate");
if ( (buf = malloc(needed)) == NULL)
    quit("malloc");
if (sysctl(mib, 6, buf, &needed, NULL, 0) < 0)
    quit("sysctl error, retrieval");
lim = buf + needed;
for (next = buf; next < lim; next += rtm->rtm_msglen) {
    rtm = (struct rt_msghdr *)next;
    ... /* do whatever */
}
```

Figure 19.35 Example of sysctl with routing table.

mib[4] specifies the operation. Three operations are supported.

1. NET\_RT\_DUMP: return the routing table corresponding to the address family specified by mib[3]. If the address family is 0, all routing tables are returned.

An RTM\_GET routing message is returned for each routing table entry containing two, three, or four socket address structures per message: those addresses pointed to by rt\_key, rt\_gateway, rt\_netmask, and rt\_genmask. The final two pointers might be null.

- 2. NET\_RT\_FLAGS: the same as the previous command except mib[5] specifies an RTF\_*xxx* flag (Figure 18.25), and only entries with this flag set are returned.
- 3. NET\_RT\_IFLIST: return information on all the configured interfaces. If the mib[5] value is nonzero it specifies an interface index and only the interface with the corresponding if\_index is returned. Otherwise all interfaces on the ifnet linked list are returned.

For each interface one RTM\_IFINFO message is returned, with information about the interface itself, followed by one RTM\_NEWADDR message for each ifaddr structure on the interface's if\_addrlist linked list. If the mib[3] value is nonzero, RTM\_NEWADDR messages are returned for only the addresses Section 19.14

19

nily

ain-

sses

The

s an

the

face

the

tion

each

[3]

sses

1.

with an address family that matches the mib[3] value. Otherwise mib[3] is 0 and information on all addresses is returned.

This operation is intended to replace the SIOCGIFCONF ioctl (Figure 4.26).

One problem with this system call is that the amount of information returned can vary, depending on the number of routing table entries or the number of interfaces. Therefore the first call to sysctl typically specifies a null pointer as the third argument, which means: don't return any data, just return the number of bytes of return information. As we see in Figure 19.35, the process then calls malloc, followed by sysctl to fetch the information. This second call to sysctl again returns the number of bytes through the fourth argument (which might have changed since the previous call), and this value provides the pointer lim that points just beyond the final byte of data that was returned. The process then steps through the routing messages in the buffer, using the rtm\_msglen member to step to the next message.

Figure 19.36 shows the values for these six mib variables that various Net/3 programs specify to access the routing table and interface list.

mib[]	arp	route	netstat	routed	gated	rwhod
0	CTL_NET	CTL_NET	CTL_NET .	CTL_NET	CTL_NET	CTL_NET
1	PF_ROUTE	PF_ROUTE	PF_ROUTE	PF_ROUTE	PF_ROUTE	PF_ROUTE
2	0	0	0	0	0	0
	AF_INET	0		AF_INET	-	AF_INET
4	NET_RT_FLAGS	NET_RT_DUMP	NET_RT_DUMP	NET_RT_IFLIST	NET_RT_IFLIST	NET_RT_IFLIST
5	RTF_LLINFO	0	0	0	0	0

Figure 19.36 Examples of programs that call sysctl to obtain routing table and interface list.

The first three programs fetch entries from the routing table and the last three fetch the interface list. The routed program supports only the Internet routing protocols, so it specifies a mib[3] value of AF\_INET, while gated supports other protocols, so its value for mib[3] is 0.

Figure 19.37 shows the organization of the three  $sysctl_xxx$  functions that we cover in the following sections.

Figure 19.38 shows the sysctl\_rtable function.

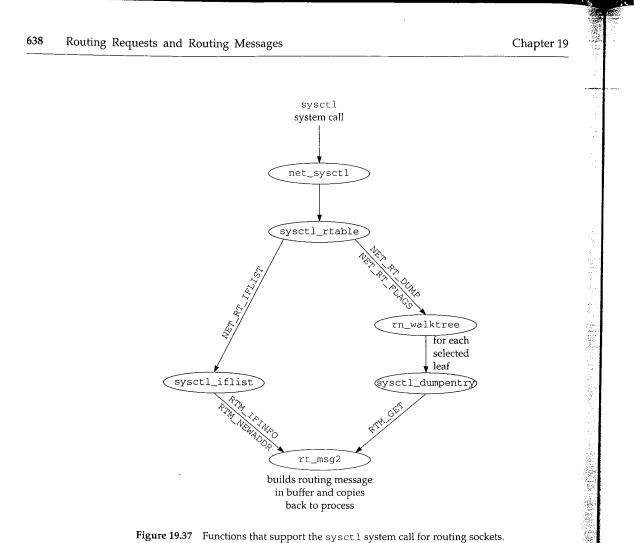
### Validate arguments

<sup>705–719</sup> The new argument is used when the process is calling sysctl to set the value of a variable, which isn't supported with the routing tables. Therefore this argument must be a null pointer.

namelen must be 3 because at this point in the processing of the system call, three elements in the name array remain: name[0], the address family (what the process specifies as mib[3]); name[1], the operation (mib[4]); and name[2], the flags (mib[5]).

INTEL Ex.1013.663

bi 1



. . .

Figure 19.37 Functions that support the sysctl system call for routing sockets.

```
705 int
706 sysctl_rtable(name, namelen, where, given, new, newlen)
707 int
           *name;
708 int
           namelen;
709 caddr_t where;
710 size_t *given;
711 caddr_t *new;
712 size_t newlen;
713 {
714
       struct radix_node_head *rnh;
715
        int
               i, s, error = EINVAL;
716
       u_char af;
717
       struct walkarg w;
718
       if (new)
719
            return (EPERM);
```

and the state of the

**INTEL Ex.1013.664** 

通常が

Section Section

「二十三三三」

rtsock.c

\_\_\_\_\_

720	if (namelen != 3)
721	return (EINVAL);
722	af = name[0];
723	<pre>Bzero(&amp;w, sizeof(w));</pre>
724	w.w_where = where;
725	w.w_given = *given;
726	w.w_needed = 0 - w.w_given;
727	w.w_op = name[1];
728	w.w_arg = name[2];
729	s = splnet();
730	switch (w.w_op) {
731	case NET_RT_DUMP:
732	case NET_RT_FLAGS:
733	for (i = 1; i <= AF_MAX; i++)
734	if ((rnh = rt_tables[i]) && (af == 0    af == i) &&
735	(error = rnh->rnh_walktree(rnh,
736	<pre>sysctl_dumpentry, &amp;w)))</pre>
737	break;
738	break;
739	case NET_RT_IFLIST:
740	error = sysctl_iflist(af, &w);
741	}
742	splx(s);
743	if (w.w_tmem)
744	free(w.w_tmem, M_RTABLE);
745	w.w_needed += w.w_given;
746	if (where) {
747	*given = w.w_where - where;
748	if (*given < w.w_needed)
749	return (ENOMEM);
750	<pre>} else {</pre>
751	*given = (11 * w.w_needed) / 10;
752	}
753	return (error);
754 )	rtsc

— rtsock.c

Figure 19.38 sysctl\_rtable function: process sysctl system call requests.

### Initialize walkarg structure

A walkarg structure (Figure 19.31) is set to 0 and the following members are initialized: w\_where is the address in the calling process of the buffer for the results (this can be a null pointer, as we mentioned); w\_given is the size of the buffer in bytes (this is meaningless on input if w\_where is a null pointer, but it must be set on return to the amount of data that would have been returned); w\_needed is set to the negative of the buffer size; w\_op is the operation (the NET\_RT\_xxx value); and w\_arg is the flags value.

#### Dump routing table

731-738 The NET\_RT\_DUMP and NET\_RT\_FLAGS operations are handled the same way: a loop is made through all the routing tables (the rt\_tables array), and if the routing

table is in use and either the address family argument was 0 or the address family argument matches the family of this routing table, the rnh\_walktree function is called to process the entire routing table. In Figure 18.17 we show that this function is normally rn\_walktree. The second argument to this function is the address of another function that is called for each leaf of the routing tree (sysctl\_dumpentry). The third argument is just a pointer to anything that rn\_walktree passes to the sysctl\_dumpentry function. This argument is a pointer to the walkarg structure that contains all the information about this sysctl call.

### **Return interface list**

739-740 The NET\_RT\_IFLIST operation calls the function sysctl\_iflist, which goes through all the ifnet structures.

#### **Release buffer**

743-744 If a buffer was allocated by rt\_msg2 to contain a routing message, it is now released.

#### Update w\_needed

745

The size of each message was added to w\_needed by rt\_msg2. Since this variable was initialized to the negative of w\_given, its value can now be expressed as

w\_needed = 0 - w\_given + totalbytes

where totalbytes is the sum of all the message lengths added by rt\_msg2. By adding the value of w\_given back into w\_needed, we get

w\_needed = 0 - w\_given + totalbytes + w\_given = totalbytes

the total number of bytes. Since the two values of  $w\_given$  in this equation end up canceling each other, when the process specifies  $w\_where$  as a null pointer it need not initialize the value of  $w\_given$ . Indeed, we see in Figure 19.35 that the variable needed was not initialized.

#### Return actual size of message

If where is nonnull, the number of bytes stored in the buffer is returned through the given pointer. If this value is less than the size of the buffer specified by the process, an error is returned because the return information has been truncated.

#### Return estimated size of message

<sup>750–752</sup> When the where pointer is null, the process just wants the total number of bytes returned. A 10% fudge factor is added to the size, in case the size of the desired tables increases between this call to sysct1 and the next.

## 19.15 sysctl\_dumpentry Function

In the previous section we described how this function is called by rn\_walktree, which in turn is called by sysctl\_rtable. Figure 19.39 shows the function.

623

Section 19.15

er 19

- rtsock.c rgu-623 int d to 624 sysctl\_dumpentry(rn, w) vally 625 struct radix\_node \*rn; :tion 626 struct walkarg \*w; 627 { rgu-34 struct rtentry \*rt = (struct rtentry \*) rn; 628 the error = 0, size; 629 int :ture struct rt\_addrinfo info; 630 if (w->w\_op == NET\_RT\_FLAGS && !(rt->rt\_flags & w->w\_arg)) 631 return 0; 632 goes bzero((caddr\_t) & info, sizeof(info)); 633 dst = rt\_key(rt); 634 qate = rt->rt\_gateway; 635 netmask = rt\_mask(rt); 636 genmask = rt->rt\_genmask; 637 now size = rt\_msg2(RTM\_GET, &info, 0, w); 638 if (w->w\_where && w->w\_tmem) { 639 struct rt\_msghdr \*rtm = (struct rt\_msghdr \*) w->w\_tmem; 640 iable 641 rtm->rtm\_flags = rt->rt\_flags; rtm->rtm\_use = rt->rt\_use; 642 rtm->rtm\_rmx = rt->rt\_rmx; 643 rtm->rtm\_index = rt->rt\_ifp->if\_index; 644 rtm->rtm\_errno = rtm->rtm\_pid = rtm->rtm\_seq = 0; 645 . By rtm->rtm\_addrs = info.rti\_addrs; 646 if (error = copyout((caddr\_t) rtm, w->w\_where, size)) 647 w->w\_where = NULL; 648 649 else 650 w->w\_where += size; 651 id up } 652 return (error); d not 653 } 電影器 riable rtsock.c Figure 19.39 sysctl\_dumpentry function: process one routing table entry. ;h the

Each time this function is called, its first argument points to a radix\_node structure, which is also a pointer to a rtentry structure. The second argument points to the walkarg structure that was initialized by sysctl\_rtable.

## Check flags of routing table entry

<sup>631–632</sup> If the process specified a flag value (mib[5]), this entry is skipped if the rt\_flags member doesn't have the desired flag set. We see in Figure 19.36 that the arp program uses this to select only those entries with the RTF\_LLINFO flag set, since these are the entries of interest to ARP.

#### Form routing message

Ę.

ocess,

bytes

tables

tree,

<sup>633-638</sup> The following four pointers in the rti\_info array are copied from the routing table entry: dst, gate, netmask, and genmask. The first two are always nonnull, but the other two can be null. rt\_msg2 forms an RTM\_GET message.

INTEL Ex.1013.667

調査する

ģ,

×.

rtsock.c

#### Copy message back to process

639-651 If the process wants the message returned and a buffer was allocated by rt\_msg2, the remainder of the routing message is formed in the buffer pointed to by w\_tmem and copyout copies the message back to the process. If the copy was successful, w\_where is incremented by the number of bytes copied.

## 19.16 sysctl\_iflist Function

This function, shown in Figure 19.40, is called directly by sysctl\_rtable to return the interface list to the process.

```
654 int
655 sysctl_iflist(af, w)
656 int
           af;
657 struct walkarg *w;
658 {
        struct ifnet *ifp;
659
        struct ifaddr *ifa;
660
661
        struct rt_addrinfo info;
662
        int
                len, error = 0;
        bzero((caddr_t) & info, sizeof(info));
663
664
        for (ifp = ifnet; ifp; ifp = ifp->if_next) {
665
            if (w->w_arg && w->w_arg != ifp->if_index)
666
                continue;
            ifa = ifp->if_addrlist;
667
668
            ifpaddr = ifa->ifa_addr;
669
            len = rt_msg2(RTM_IFINFO, &info, (caddr_t) 0, w);
670
            ifpaddr = 0;
671
            if (w->w_where && w->w_tmem) {
672
                struct if_msghdr *ifm;
673
                ifm = (struct if_msghdr *) w->w_tmem;
674
                ifm->ifm_index = ifp->if_index;
675
                ifm->ifm_flags = ifp->if_flags;
676
                ifm->ifm_data = ifp->if_data;
677
                ifm_addrs = info.rti_addrs;
678
                if (error = copyout((caddr_t) ifm, w->w_where, len))
679
                    return (error);
680
                w->w_where += len;
681
            }
682
            while (ifa = ifa->ifa_next) {
683
                if (af && af != ifa->ifa_addr->sa_family)
684
                    continue;
685
                ifaaddr = ifa->ifa_addr;
686
                netmask = ifa->ifa_netmask;
687
                brdaddr = ifa->ifa_dstaddr;
688
                len = rt_msg2(RTM_NEWADDR, &info, 0, w);
689
                if (w->w_where && w->w_tmem) {
690
                    struct ifa_msghdr *ifam;
```

Chapter 19

rt\_msg2, \_tmem and ,w\_where

return the

rtsock.c

Section 19.16 643 sysctl iflist Function ifam = (struct ifa\_msghdr \*) w->w\_tmem; 691 ifam->ifam\_index = ifa->ifa\_ifp->if\_index; 692 ifam->ifam\_flags = ifa->ifa\_flags; 693 ifam->ifam\_metric = ifa->ifa\_metric; 694 ifam->ifam\_addrs = info.rti\_addrs; 695 if (error = copyout(w->w\_tmem, w->w\_where, len)) 696 return (error); 697 w->w where += len; 698 699 } 700 } ifaaddr = netmask = brdaddr = 0; 701 702 } 703 return (0); 704 } - rtsock.c

Figure 19.40 sysctl\_iflist function: return list of interfaces and their addresses.

This function is a for loop that iterates through each interface starting with the one pointed to by ifnet. Then a while loop proceeds through the linked list of ifaddr structures for each interface. An RTM\_IFINFO routing message is generated for each interface and an RTM\_NEWADDR message for each address.

#### Check interface index

The process can specify a nonzero flags argument (mib[5] in Figure 19.36) to select 654-666 only the interface with a matching if\_index value.

#### **Build routing message**

The only socket address structure returned with the RTM\_IFINFO message is 667-670 ifpaddr. The message is built by rt\_msg2. The pointer ifpaddr in the info structure is then set to 0, since the same info structure is used for generating the subsequent RTM\_NEWADDR messages.

#### Copy message back to process

If the process wants the message returned, the remainder of the if\_msghdr struc-671-681 ture is filled in, copyout copies the buffer to the process, and w\_where is incremented.

#### Iterate through address structures, check address family

Each ifaddr structure for the interface is processed and the process can specify a 682-684 nonzero address family (mib[3] in Figure 19.36) to select only the interface addresses of the given family.

#### **Build routing message**

701

Up to three socket address structures are returned in each RTM\_NEWADDR message: 685-688 ifaaddr, netmask, and brdaddr. The message is built by rt\_msg2.

#### Copy message back to process

If the process wants the message returned, the remainder of the ifa\_msghdr struc-689-699 ture is filled in, copyout copies the buffer to the process, and w\_where is incremented. These three pointers in the info array are set to 0, since the same array is used for the next interface message.

INTEL Ex.1013.669

 $\dot{\alpha}_1$ 

## 19.17 Summary

Routing messages all have the same format—a fixed-length structure followed by a variable number of socket address structures. There are three different types of messages, each corresponding to a different fixed-length structure, and the first three elements of each structure identify the length, version, and type of message. A bitmask in each structure identifies which socket address structures follow the fixed-length structure.

These messages are passed between a process and the kernel in two different ways. Messages can be passed in either direction, one message per read or write, across a routing socket. This allows a superuser process complete read and write access to the kernel's routing tables. This is how routing daemons such as routed and gated implement their desired routing policy.

Alternatively any process can read the contents of the kernel's routing tables using the sysctl system call. This does not involve a routing socket and does not require special privileges. The entire result, normally consisting of many routing messages, is returned as part of the system call. Since the process does not know the size of the result, a method is provided for the system call to return this size without returning the actual result.

### Exercises

- **19.1** What is the difference in the RTF\_DYNAMIC and RTF\_MODIFIED flags? Can both be set for a given routing table entry?
- **19.2** What happens when the default route is entered with a command of the form

bsdi \$ route add default -cloning -genmask 255.255.255.255 sun

**19.3** Estimate the space required by sysctl to dump a routing table that contains 15 ARP entries and 20 routes.

19

a :slein IC-

ys. ut-

ered

ng ire , is :he :he

for

١RP

# **Routing Sockets**

## 20.1 Introduction

A process sends and receives the routing messages described in the previous chapter by using a socket in the *routing domain*. The socket system call is issued specifying a family of PF\_ROUTE and a socket type of SOCK\_RAW.

The process can then send five routing messages to the kernel:

- 1. RTM\_ADD: add a new route.
- 2. RTM\_DELETE: delete an existing route.
- 3. RTM\_GET: fetch all the information about a route.
- 4. RTM\_CHANGE: change the gateway, interface, or metrics of an existing route.
- 5. RTM\_LOCK: specify which metrics the kernel should not modify.

Additionally, the process can receive any of the other seven types of routing messages that are generated by the kernel when some event, such as interface down, redirect received, etc., occurs.

This chapter looks at the routing domain, the routing control blocks that are created for each routing socket, the function that handles messages from a process (route\_output), the function that sends routing messages to one or more processes (raw\_input), and the various functions that support all the socket operations on a routing socket.

INTEL Ex.1013.671

Sales in

語を言語を見ていた言語

# 20.2 routedomain and protosw Structures

Before describing the routing socket functions, we need to discuss additional details about the routing domain; the SOCK\_RAW protocol supported in the routing domain; and routing control blocks, one of which is associated with each routing socket.

Figure 20.1 lists the domain structure for the PF\_ROUTE domain, named routedomain.

Member	Value	Description
dom_family dom_name dom_init dom_externalize dom_dispose dom_protosw dom_protoswNPROTOSW dom_next dom_rtattach dom_rtaffset dom_maxrtkey	PF_ROUTE route route_init 0 0 routesw 0 0 0	protocol family for domain name domain initialization, Figure 18.30 not used in routing domain not used in routing domain protocol switch structure, Figure 20.2 pointer past end of protocol switch structure filled in by domaininit, Figure 7.15 not used in routing domain not used in routing domain not used in routing domain

Figure 20.1 routedomain structure.

Unlike the Internet domain, which supports multiple protocols (TCP, UDP, ICMP, etc.), only one protocol (of type SOCK\_RAW) is supported in the routing domain. Figure 20.2 lists the protocol switch entry for the PF\_ROUTE domain.

Member	routesw[0]	Description
pr_type pr_domain pr_protocol pr_flags pr_input pr_output pr_ctlinput pr_ctloutput pr_usrreq pr_init pr_fasttimo pr_slowtimo pr_drain pr_sysctl	SOCK_RAW &routedomain 0 PR_ATOMIC PR_ADDR raw_input route_output raw_ctlinput 0 route_usrreq raw_init 0 0 0 0 sysctl_rtable	raw socket part of the routing domain socket layer flags, not used by protocol processing this entry not used; raw_input called directly called for PRU_SEND requests control input function not used respond to communication requests from a process initialization not used not used not used for sysct1(8) system call

Figure 20.2 The routing protocol protosw structure.

raw\_init Function 647

### Section 20.4

## 20.3 Routing Control Blocks

Each time a routing socket is created with a call of the form

socket(PF\_ROUTE, SOCK\_RAW, protocol);

the corresponding PRU\_ATTACH request to the protocol's user-request function (route\_usrreq) allocates a routing control block and links it to the socket structure. The *protocol* can restrict the messages sent to the process on this socket to one particular family. If a *protocol* of AF\_INET is specified, for example, only routing messages containing Internet addresses will be sent to the process. A *protocol* of 0 causes all routing messages from the kernel to be sent on the socket.

Recall that we call these structures *routing control blocks*, not *raw control blocks*, to avoid confusion with the raw IP control blocks in Chapter 32.

Figure 20.3 shows the definition of the rawcb structure.

```
– raw_cb.h
39 struct rawcb {
      struct rawcb *rcb_next;
                                   /* doubly linked list */
40
       struct rawcb *rcb_prev;
41
      struct socket *rcb_socket; /* back pointer to socket */
42
       struct sockaddr *rcb_faddr; /* destination address */
43
       struct sockaddr *rcb_laddr; /* socket's address */
44
       struct sockproto rcb_proto; /* protocol family, protocol */
45
46 };
                                ((struct rawcb *)(so)->so_pcb)
47 #define sotorawcb(so)
                                                                         – raw_cb.h
```

Figure 20.3 rawcb structure.

Additionally, a global of the same name, rawcb, is allocated as the head of the doubly linked list. Figure 20.4 shows the arrangement.

<sup>39-47</sup> We showed the sockproto structure in Figure 19.26. Its sp\_family member is set to PF\_ROUTE and its sp\_protocol member is set to the third argument to the socket system call. The rcb\_faddr member is permanently set to point to route\_src, which we described with Figure 19.26. rcb\_laddr is always a null pointer.

## 20.4 raw\_init Function

The raw\_init function, shown in Figure 20.5, is the protocol initialization function in the protosw structure in Figure 20.2. We described the entire initialization of the routing domain with Figure 18.29.

<sup>38–42</sup> The function initializes the doubly linked list of routing control blocks by setting the next and previous pointers of the head structure to point to itself.

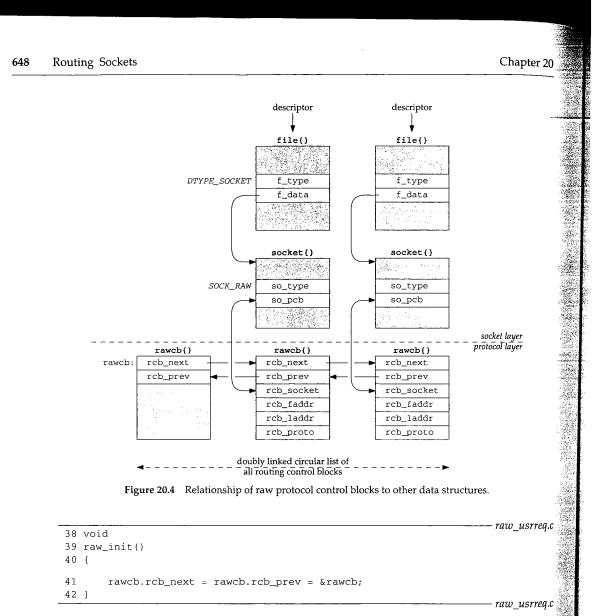


Figure 20.5 raw\_init function: initialize doubly linked list of routing control blocks.

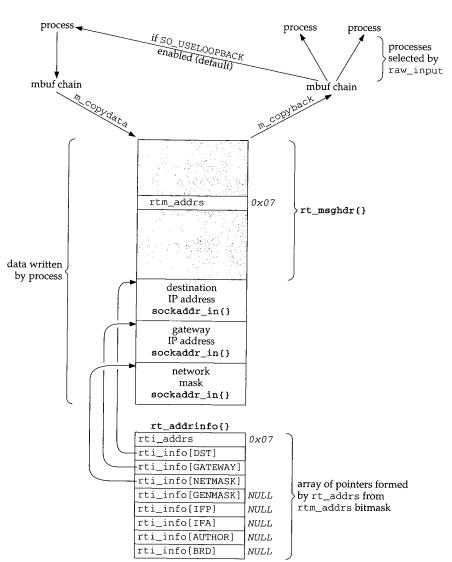
## 20.5 route\_output Function

As we showed in Figure 18.11, route\_output is called when the PRU\_SEND request is issued to the protocol's user-request function, which is the result of a write operation by a process to a routing socket. In Figure 18.9 we indicated that five different types of routing messages are accepted by the kernel from a process.

Since this function is invoked as a result of a write by a process, the data from the process (the routing message to process) is in an mbuf chain from sosend. Figure 20.6

Section 20.5

÷



shows an overview of the processing steps, assuming the process sends an RTM\_ADD command, specifying three addresses: the destination, its gateway, and a network mask (hence this is a network route, not a host route).

Figure 20.6 Example processing of an RTM\_ADD command from a process.

There are numerous points to note in this figure, most of which we'll cover as we proceed through the source code for route\_output. Also note that, to save space, we omit the RTAX\_ prefix for each array index in the rt\_addrinfo structure.

**INTEL Ex.1013.675** 

Chapter 20

- The process specifies which socket address structures follow the fixed-length rt\_msghdr structure by setting the bitmask rtm\_addrs. We show a bitmask of 0x07, which corresponds to a destination address, a gateway address, and a network mask (Figure 19.19). The RTM\_ADD command requires the first two; the third is optional. Another optional address, the genmask specifies the mask to be used for generating cloned routes.
- The write system call (the sosend function) copies the buffer from the process into an mbuf chain in the kernel.
- m\_copydata copies the mbuf chain into a buffer that route\_output obtains using malloc. It is easier to access all the information in the structure and the socket address structures that follow when stored in a single contiguous buffer than it is when stored in an mbuf chain.
- The function rt\_xaddrs is called by route\_output to take the bitmask and build the rt\_addrinfo structure that points into the buffer. The code in route\_output references these structures using the names shown in the fifth column in Figure 19.19. The bitmask is also copied into the rti\_addrs member.
- route\_output normally modifies the rt\_msghdr structure. If an error occurs, the corresponding errno value is returned in rtm\_errno (for example, EEXIST if the route already exists); otherwise the flag RTF\_DONE is logically ORed into the rtm\_flags supplied by the process.
- The rt\_msghdr structure and the addresses that follow become input to 0 or more processes that are reading from a routing socket. The buffer is first converted back into an mbuf chain by m\_copyback. raw\_input goes through all the routing PCBs and passes a copy to the appropriate processes. We also show that a process with a routing socket receives a copy of each message it writes to that socket unless it disables the SO\_USELOOPBACK socket option.

To avoid receiving a copy of their own routing messages, some programs, such as route, call shutdown with a second argument of 0 to prevent any data from being received on the routing socket.

We examine the source code for route\_output in seven parts. Figure 20.7 shows an overview of the function.

```
int
route_output()
{
    R_Malloc() to allocate buffer;
    m_copydata() to copy from mbuf chain into buffer;
    rt_xaddrs() to build rt_addrinfo();
    switch (message type) {
    case RTM_ADD:
        rtrequest(RTM_ADD);
        rt_setmetrics();
        break;
```

Chapter 20

Section 20.5

route\_output Function 651

ed-length bitmask ss, and a two; the mask to

e process

: obtains e and the us buffer

nask and code in the fifth rs mem-

or occurs, example, logically

it to 0 or first conrough all lso show writes to

h as route, received on

shows an

case RTM\_DELETE: rtrequest(RTM\_DELETE); break; case RTM\_GET: case RTM\_CHANGE: case RTM\_LOCK: rtalloc1(); switch (message type) { case RTM\_GET: rt\_msg2(RTM\_GET); break; case RTM\_CHANGE: change appropriate fields; /\* fall through \*/ case RTM\_LOCK: set rmx\_locks; break; 3 break; } set rtm\_error if error, else set RTF\_DONE flag; m\_copyback() to copy from buffer into mbuf chain; /\* mbuf chain to appropriate processes \*/ raw\_input();

Figure 20.7 Summary of route\_output processing steps.

The first part of route\_output is shown in Figure 20.8.

#### Check mbuf for validity

<sup>113-136</sup> The mbuf chain is checked for validity: its length must be at least the size of an rt\_msghdr structure. The first longword is fetched from the data portion of the mbuf, which contains the rtm\_msglen value.

#### Allocate buffer

A buffer is allocated to hold the entire message and m\_copydata copies the message from the mbuf chain into the buffer.

#### Check version number

<sup>143-146</sup> The version of the message is checked. In the future, should a new version of the routing messages be introduced, this member could be used to provide support for older versions.

147-149 The process ID is copied into rtm\_pid and the bitmask supplied by the process is copied into info.rti\_addrs, a structure local to this function. The function rt\_xaddrs (shown in the next section) fills in the eight socket address pointers in the info structure to point into the buffer now containing the message.

652 Routing Sockets

Chapter 20

rtsock.c

```
113 int
114 route_output(m, so)
115 struct mbuf *m;
116 struct socket *so;
117 {
        struct rt_msghdr *rtm = 0;
118
        struct rtentry *rt = 0;
119
        struct rtentry *saved_nrt = 0;
120
        struct rt_addrinfo info;
121
                len, error = 0;
122
        int
        struct ifnet *ifp = 0;
123
        struct ifaddr *ifa = 0;
124
125 #define senderr(e) { error = e; goto flush;}
        if (m == 0 |! ((m->m_len < sizeof(long)) &&
126
                                (m = m_pullup(m, sizeof(long))) == 0))
127
                     return (ENOBUFS);
128
        if ((m->m_flags \& M_PKTHDR) == 0)
129
            panic("route_output");
130
        len = m->m_pkthdr.len;
131
        if (len < sizeof(*rtm) ||
132
            len != mtod(m, struct rt_msghdr *)->rtm_msglen) {
133
134
             dst = 0;
             senderr(EINVAL);
135
136
         3
        R_Malloc(rtm, struct rt_msghdr *, len);
137
         if (rtm == 0) {
138
139
             dst = 0;
             senderr(ENOBUFS);
140
141
         }
         m_copydata(m, 0, len, (caddr_t) rtm);
142
         if (rtm->rtm_version != RTM_VERSION) {
143
144
             dst = 0;
             senderr(EPROTONOSUPPORT);
145
 146
         }
         rtm->rtm_pid = curproc->p_pid;
 147
         info.rti_addrs = rtm->rtm_addrs;
 148
         rt_xaddrs((caddr_t) (rtm + 1), len + (caddr_t) rtm, &info);
 149
 150
         if (dst == 0)
 151
             senderr(EINVAL);
 152
         if (genmask) {
             struct radix_node *t;
 153
             t = rn_addmask((caddr_t) genmask, 1, 2);
 154
             if (t && Bcmp(genmask, t->rn_key, *(u_char *) genmask) == 0)
 155
                 genmask = (struct sockaddr *) (t->rn_key);
 156
 157
             else
                  senderr(ENOBUFS);
 158
                                                                              rtsock.c
 159
         }
```

Figure 20.8 route\_output function: initial processing, copy message from mbuf chain.

INTEL Ex.1013.678

apter 20

rtsock.c

Section 20.5

rtsock c

#### Destination address required

150-151 A destination address is a required address for all commands. If the info.rti\_info[RTAX\_DST] element is a null pointer, EINVAL is returned. Remember that dst refers to this array element (Figure 19.19).

#### Handle optional genmask

A genmask is optional and is used as the network mask for routes created when the RTF\_CLONING flag is set (Figure 19.8). rn\_addmask adds the mask to the tree of masks, first searching for an existing entry for the mask and then referencing that entry if found. If the mask is found or added to the mask tree, an additional check is made that the entry in the mask tree really equals the genmask value, and, if so, the genmask pointer is replaced with a pointer to the mask in the mask tree.

Figure 20.9 shows the next part of route\_output, which handles the RTM\_ADD and RTM\_DELETE commands.

160	<pre>switch (rtm-&gt;rtm_type) {</pre>
161	case RTM_ADD:
162	if $(gate == 0)$
163	senderr(EINVAL);
164	error = rtrequest(RTM_ADD, dst, gate, netmask,
165	rtm->rtm_flags, &saved_nrt);
166	if (error == 0 && saved_nrt) {
167	rt_setmetrics(rtm->rtm_inits,
168	&rtm->rtm_rmx, &saved_nrt->rt_rmx);
169	<pre>saved_nrt-&gt;rt_refcnt;</pre>
170	saved_nrt->rt_genmask = genmask;
171	. )
172	break;
173	case RTM_DELETE:
174	error = rtrequest(RTM_DELETE, dst, gate, netmask,
175	<pre>rtm-&gt;rtm_flags, (struct rtentry **) 0);</pre>
176	break;rtsock.c

Figure 20.9 route\_output function: process RTM\_ADD and RTM\_DELETE commands.

An RTM\_ADD command requires the process to specify a gateway.

164-165 rtrequest processes the request. The netmask pointer can be null if the route being entered is a host route. If all is OK, the pointer to the new routing table entry is returned through saved\_nrt.

<sup>166-172</sup> The rt\_metrics structure is copied from the caller's buffer into the routing table entry. The reference count is decremented and the genmask pointer is stored (possibly a null pointer).

<sup>173–176</sup> Processing the RTM\_DELETE command is simple because all the work is done by rtrequest. Since the final argument is a null pointer, rtrequest calls rtfree if the reference count is 0, deleting the entry from the routing table (Figure 19.7).

INTEL Ex.1013.679

654 Routing Sockets Chapter 20

5

1

The next part of the processing is shown in Figure 20.10, which handles the common code for the RTM\_GET, RTM\_CHANGE, and RTM\_LOCK commands.

	rtsock.
177	case RTM_GET:
178	case RTM_CHANGE:
179	case RTM_LOCK:
180	rt = rtalloc1(dst, 0);
181	if $(rt == 0)$
182	senderr(ESRCH);
183	if (rtm->rtm_type != RTM_GET) {
184	struct radix_node *rn;
185	extern struct radix_node_head *mask_rnhead;
186	if (Bcmp(dst, rt_key(rt), dst->sa_len) != 0)
187	senderr(ESRCH);
188	if (netmask && (rn = rn_search(netmask,
189	<pre>mask_rnhead-&gt;rnh_treetop)))</pre>
190	netmask = (struct sockaddr *) rn->rn_key;
191	for (rn = rt->rt_nodes; rn; rn = rn->rn_dupedkey)
192	if (netmask == (struct sockaddr *) rn->rn_mask)
193	break;
194	if (rn == 0)
195	senderr(ETOOMANYREFS);
196	<pre>rt = (struct rtentry *) rn;</pre>
197	} rtsock

Figure 20.10 route\_output function: common processing for RTM\_GET, RTM\_CHANGE, and RTM\_LOCK.

#### Locate existing entry

Since all three commands reference an existing entry, rtalloc1 locates the entry. If 177-182 the entry isn't found, ESRCH is returned.

#### Do not allow network match

For the RTM\_CHANGE and RTM\_LOCK commands, a network match is inadequate: an 183-187 exact match with the routing table key is required. Therefore, if the dst argument doesn't equal the routing table key, the match was a network match and ESRCH is returned.

#### Use network mask to find correct entry

Even with an exact match, if there are duplicate keys, each with a different network 188 - 193mask, the correct entry must still be located. If a netmask argument was supplied, it is looked up in the mask table (mask\_rnhead). If found, the netmask pointer is replaced with the pointer to the mask in the mask tree. Each leaf node in the duplicate key list is examined, looking for an entry with an rn\_mask pointer that equals netmask. This test compares the pointers, not the structures that they point to. This works because all masks appear in the mask tree, and only one copy of each unique mask is stored in this tree. In the common case, keys are not duplicated, so the for loop iterates once. If a host entry is being modified, a mask must not be specified and then both netmask and rn\_mask are null pointers (which are equal). But if an entry that has an associated mask is being modified, that mask must be specified as the netmask argument.

Section 20.5

he com-

napter 20

- rtsock.c

194-195 If the for loop terminates without finding a matching network mask, ETOOMANYREFS is returned.

The comment XXX is because this function must go to all this work to find the desired entry. All these details should be hidden in another function similar to rtalloc1 that detects a network match and handles a mask argument.

The next part of this function, shown in Figure 20.11, continues processing the RTM\_GET command. This command is unique among the commands supported by route\_output in that it can return more data than it was passed. For example, only a single socket address structure is required as input, the destination, but at least two are returned: the destination and its gateway. With regard to Figure 20.6, this means the buffer allocated for m\_copydata to copy into might need to be increased in size.

w.			rtsock.c
	198	<pre>switch (rtm-&gt;rtm_type) {</pre>	
	199	case RTM_GET:	
	200	dst = rt_key(rt);	
	201	<pre>gate = rt-&gt;rt_gateway;</pre>	
	202	<pre>netmask = rt_mask(rt);</pre>	
	203	<pre>genmask = rt-&gt;rt_genmask;</pre>	-
	204	if (rtm->rtm_addrs & (RTA_IFP   RTA_IFA)) {	
	205	if (ifp = rt->rt_ifp) {	
	206	ifpaddr = ifp->if_addrlist->ifa_addr;	
2	207	ifaaddr = rt->rt_ifa->ifa_addr;	
V	208	rtm->rtm_index = ifp->if_index;	
-	209	} else {	
d.	210	ifpaddr = 0;	
2	211	ifaaddr = 0;	
	212	}	
	213	}	
	214	len = rt_msg2(RTM_GET, &info, (caddr_t) 0,	
	215	<pre>(struct walkarg *) 0);</pre>	
- 	216	if (len > rtm->rtm_msglen) {	
	217	<pre>struct rt_msghdr *new_rtm;</pre>	
	218	R_Malloc(new_rtm, struct rt_msghdr *, len);	
7	219	if (new_rtm == 0)	
	220	senderr(ENOBUFS);	
	221	<pre>Bcopy(rtm, new_rtm, rtm-&gt;rtm_msglen);</pre>	
	222	<pre>Free(rtm);</pre>	
	223	<pre>rtm = new_rtm;</pre>	
	224	}	
	225	<pre>(void) rt_msg2(RTM_GET, &amp;info, (caddr_t) rtm,</pre>	
	226	<pre>(struct walkarg *) 0);</pre>	
7	227	rtm->rtm_flags = rt->rt_flags;	
	228	<pre>rtm-&gt;rtm_rmx = rt-&gt;rt_rmx;</pre>	
-	229	rtm->rtm_addrs = info.rti_addrs;	
· ·	230	break;	

Figure 20.11 route\_output function: RTM\_GET processing.

))

— *rtsock.c* 4\_lock.

entry. If

luate: an rgument SRCH is

network lied, it is replaced ey list is sk. This cause all d in this ace. If a ask and

sociated

INTEL Ex.1013.681

#### Return destination, gateway, and masks

Four pointers are stored in the rti\_info array: dst, gate, netmask, and 198-203 genmask. The latter two might be null pointers. These pointers in the info structure point to the socket address structures that will be returned to the process.

#### **Return interface information**

The process can set the masks RTA\_IFP and RTA\_IFA in the rtm\_flags bitmask. 204-213 If either or both are set, the process wants to receive the contents of both the ifaddr structures pointed to by this routing table entry: the link-level address of the interface (pointed to by rt\_ifp->if\_addrlist) and the protocol address for this entry (pointed to by rt\_ifa->ifa\_addr). The interface index is also returned.

#### Construct reply

- rt\_msg2 is called with a null third pointer to calculate the length of the routing 214-224 message corresponding to RTM\_GET and the addresses pointed to by the info structure. If the length of the result message exceeds the length of the input message, then a new buffer is allocated, the input message is copied into the new buffer, the old buffer is released, and rtm is set to point to the new buffer.
- rt\_msg2 is called again, this time with a nonnull third pointer, which builds the 225-230 result message in the buffer. The final three members in the rt\_msghdr structure are then filled in.

Figure 20.12 shows the processing of the RTM\_CHANGE and RTM\_LOCK commands.

#### Change gateway

If a gate address was passed by the process, rt\_setgate is called to change the 231-233 gateway for the entry.

#### Locate new interface

The new gateway (if changed) can also require new rt\_ifp and rt\_ifa pointers. 234-244 The process can specify these new values by passing either an ifpaddr socket address structure or an ifaaddr socket address structure. The former is tried first, and then the latter. If neither is passed by the process, the rt\_ifp and rt\_ifa pointers are left alone.

#### Check if interface changed

If an interface was located (ifa is nonnull), then the existing rt\_ifa pointer for 245-256 the route is compared to the new value. If it has changed, new values for rt\_ifp and rt\_ifa are stored in the routing table entry. Before doing this the interface request function (if defined) is called with a command of RTM\_DELETE. The delete is required because the link-layer information from one type of network to another can be quite different, say changing a route from an X.25 network to an Ethernet, and the output routines must be notified.

#### Update metrics

The metrics in the routing table entry are updated by rt\_setmetrics. 257-258

「「「「「「「「「「」」」」

Section 20.5

いたまで

## route\_output Function 657

231	case RTM_CHANGE:
231	
	<pre>if (gate &amp;&amp; rt_setgate(rt, rt_key(rt), gate)) </pre>
233	senderr(EDQUOT);
234	/* new gateway could require new ifaddr, ifp; flags may also be
235	different; ifp may be specified by ll sockaddr when protocol
236	address is ambiguous */
237	if (ifpaddr && (ifa = ifa_ifwithnet(ifpaddr)) &&
238	(ifp = ifa->ifa_ifp))
239	ifa = ifaof_ifpforaddr(ifaaddr ? ifaaddr : gate,
240	ifp);
241	else if ((ifaaddr && (ifa = ifa_ifwithaddr(ifaaddr)))
242	<pre>(ifa = ifa_ifwithroute(rt-&gt;rt_flags,</pre>
243	rt_key(rt), gate)))
244	<pre>ifp = ifa-&gt;ifa_ifp;</pre>
245	if (ifa) {
246	<pre>struct ifaddr *oifa = rt-&gt;rt_ifa;</pre>
247	if (oifa != ifa) {
248	if (oifa && oifa->ifa_rtrequest)
249	oifa->ifa_rtrequest(RTM_DELETE,
250	rt, gate);
251	IFAFREE(rt->rt_ifa);
252	rt ->rt ifa = ifa;
252	— · · · · · · · · · · · · · · · · · · ·
	<pre>ifa-&gt;ifa_refcnt++;</pre>
254	<pre>rt-&gt;rt_ifp = ifp;</pre>
255	}
256	)
257	rt_setmetrics(rtm->rtm_inits, &rtm->rtm_rmx,
258	&rt->rt_rmx);
259	if (rt->rt_ifa && rt->rt_ifa->ifa_rtrequest)
260	rt->rt_ifa->ifa_rtrequest(RTM_ADD, rt, gate);
261	if (genmask)
262	rt->rt_genmask = genmask;
263	/*
264	* Fall into
265	*/
266	case RTM_LOCK:
267	rt->rt_rmx.rmx_locks &= ~(rtm->rtm_inits);
268	rt->rt_rmx.rmx_locks /=
269	(rtm->rtm_inits & rtm->rtm_rmx.rmx_locks);
270	break;
271	}
272	, break;
273	default:
274	senderr(EOPNOTSUPP);

 $Figure \ 20.12 \quad \texttt{route\_output function: RTM\_CHANGE and RTM\_LOCK processing.}$ 

## Call interface request function

259–260 If an interface request function is defined, it is called with a command of RTM\_ADD.

INTEL Ex.1013.683

Se

### Store clone generation mask

261-262 If the process specifies the genmask argument, the pointer to the mask that was obtained in Figure 20.8 is saved in rt\_genmask.

## Update bitmask of locked metrics

<sup>266–270</sup> The RTM\_LOCK command updates the bitmask stored in rt\_rmx.rmx\_locks. Figure 20.13 shows the values of the different bits in this bitmask, one value per metric.

Constant	Value	Description
RTV MTU	0x01	initialize or lock rmx_mtu
RTV_HOPCOUNT	0x02	initialize or lock rmx_hopcount
RTV EXPIRE	0x04	initialize or lock rmx_expire
RTV_RPIPE	0x08	initialize or lock rmx_recvpipe
RTV_SPIPE	0x10	initialize or lock rmx_sendpipe
RTV_SSTHRESH	0x20	initialize or lock rmx_ssthresh
RTV_RTT	0x40	initialize or lock rmx_rtt
RTV_RTTVAR	0x80	initialize or lock rmx_rttvar

Figure 20.13 Constants to initialize or lock metrics.

The rmx\_locks member of the rt\_metrics structure in the routing table entry is the bitmask telling the kernel which metrics to leave alone. That is, those metrics specified by rmx\_locks won't be updated by the kernel. The only use of these metrics by the kernel is with TCP, as noted with Figure 27.3. The rmx\_pksent metric cannot be locked or initialized, but it turns out this member is never even referenced or updated by the kernel.

The rtm\_inits value in the message from the process specifies the bitmask of which metrics were just initialized by rt\_setmetrics. The rtm\_rmx.rmx\_locks value in the message specifies the bitmask of which metrics should now be locked. The value of rt\_rmx.rmx\_locks is the bitmask in the routing table of which metrics are currently locked. First, any bits to be initialized (rtm\_inits) are unlocked. Any bits that are both initialized (rtm\_inits) and locked (rtm\_rmx.rmx\_locks) are locked.

273–275

This default is for the switch at the beginning of Figure 20.9 and catches any of the routing commands other than the five that are supported in messages from a process.

The final part of route\_output, shown in Figure 20.14, sends the reply to raw\_input.

2

apter 20	
"PICI 20	

Section 20.5

route\_output Function 659

nat was

:s. Figric.

 $\frac{1}{20}$ 

pecified ; by the nnot be updated nask of \_locks ed. The

'y is the

\_locks ed. The trics are any bits cked. s any of

eply to

1 a pro-

#### rtsock.c flush: 276 if (rtm) { 277 if (error) 278 rtm->rtm\_errno = error; 279 280 else rtm->rtm\_flags |= RTF\_DONE; 281 282 } if (rt) 283 284 rtfree(rt); 285 { struct rawcb \*rp = 0;286 287 /\* \* Check to see if we don't want our own messages. 288 289 \*/ if ((so->so\_options & SO\_USELOOPBACK) == 0) { 290 if (route\_cb.any\_count <= 1) { 291 if (rtm) 292 Free(rtm); 293 m\_freem(m); 294 return (error); 295 } 296 /\* There is another listener, so construct message \*/ 297 rp = sotorawcb(so); 298 299 } if (rtm) { 300 m\_copyback(m, 0, rtm->rtm\_msglen, (caddr\_t) rtm); 301 Free(rtm); 302 303 } if (rp) 304 /\* Avoid us \*/ rp->rcb\_proto.sp\_family = 0; 305 if (dst) 306 route\_proto.sp\_protocol = dst->sa\_family; 307 raw\_input(m, &route\_proto, &route\_src, &route\_dst); 308 if (rp) 309 rp->rcb\_proto.sp\_family = PF\_ROUTE; 310 311 }

Figure 20.14 route\_output function: pass results to raw\_input.

Return error or OK

312

313 }

return (error);

276-282 flush is the label jumped to by the senderr macro defined at the beginning of the function. If an error occurred it is returned in the rtm\_errno member; otherwise the RTF\_DONE flag is set.

## **Release held route**

<sup>283–284</sup> If a route is being held, it is released. The call to rtalloc1 at the beginning of Figure 20.10 holds the route, if found.

1611

rtsock.c

#### No process to receive message

285--296

The SO\_USELOOPBACK socket option is true by default and specifies that the sending process is to receive a copy of each routing message that it writes to a routing socket. (If the sender doesn't receive a copy, it can't receive any of the information returned by RTM\_GET.) If that option is not set, and the total count of routing sockets is less than or equal to 1, there are no other processes to receive the message and the sender doesn't want a copy. The buffer and mbuf chain are both released and the function returns.

#### Other listeners but no loopback copy

There is at least one other listener but the sending process does not want a copy. 297-299 The pointer rp, which defaults to null, is set to point to the routing control block for the sender and is also used as a flag that the sender doesn't want a copy.

#### Convert buffer into mbuf chain

The buffer is converted back into an mbuf chain (Figure 20.6) and the buffer 300-303 released.

#### Avoid loopback copy

If rp is set, some other process might want the message but the sender does not 304-305 want a copy. The sp\_family member of the sender's routing control block is temporarily set to 0, but the sp\_family of the message (the route\_proto structure, shown with Figure 19.26) has a family of PF\_ROUTE. This trick prevents raw\_input from passing a copy of the result to the sending process because raw\_input does not pass a copy to any socket with an sp\_family of 0.

#### Set address family of routing message

If dst is a nonnull pointer, the address family of that socket address structure 306-308 becomes the protocol of the routing message. With the Internet protocols this value would be PF\_INET. A copy is passed to the appropriate listeners by raw\_input.

309-313 If the sp\_family member in the calling process was temporarily set to 0, it is reset to PF\_ROUTE, its normal value.

#### 20.6 rt\_xaddrs Function

The rt\_xaddrs function is called only once from route\_output (Figure 20.8) after the routing message from the process has been copied from the mbuf chain into a buffer and after the bitmask from the process (rtm\_addrs) has been copied into the rti\_info member of an rt\_addrinfo structure. The purpose of rt\_xaddrs is to take this bitmask and set the pointers in the rti\_info array to point to the corresponding address in the buffer. Figure 20.15 shows the function.

rtsock.c

<sup>330 #</sup>define ROUNDUP(a) \

<sup>((</sup>a) > 0 ? (1 + (((a) - 1) ) (sizeof(long) - 1))) : sizeof(long)) 331 332 #define ADVANCE(x, n) (x += ROUNDUP((n)->sa\_len))

Chapter 20

that the sendto a routing e information ting sockets is sage and the and the func-

want a copy. l block for the

nd the buffer

nder does not block is temoto structure, s raw\_input uput does not

lress structure ols this value \_input. to 0, it is reset

ure 20.8) after in into a buffer pied into the \_xaddrs is to he correspond-

— rtsock.c

(long))

```
333 static void
334 rt_xaddrs(cp, cplim, rtinfo)
335 caddr_t cp, cplim;
336 struct rt_addrinfo *rtinfo;
337 {
        struct sockaddr *sa;
338
339
        int
                i:
        bzero(rtinfo->rti_info, sizeof(rtinfo->rti_info));
340
        for (i = 0; (i < RTAX_MAX) && (cp < cplim); i++) {
341
            if ((rtinfo->rti_addrs & (1 << i)) == 0)
342
                continue;
343
            rtinfo->rti_info[i] = sa = (struct sockaddr *) cp;
344
            ADVANCE(cp, sa);
345
        }
346
347 }
```

- rtsock.c

661

rt\_setmetrics Function

Figure 20.15 rt\_xaddrs function: fill rti\_into array with pointers.

330-340

Section 20.7

The array of pointers is set to 0 so all the pointers to address structures not appearing in the bitmask will be null.

Each of the 8 (RTAX\_MAX) possible bits in the bitmask is tested and, if set, a pointer 341-347 is stored in the rti\_info array to the corresponding socket address structure. The ADVANCE macro takes the sa\_len field of the socket address structure, rounds it up to the next multiple of 4 bytes, and increments the pointer cp accordingly.

#### rt\_setmetrics Function 20.7

This function was called twice from route\_output: when a new route was added and when an existing route was changed. The rtm\_inits member in the routing message from the process specifies which of the metrics the process wants to initialize from the rtm\_rmx array. The bit values in the bitmask are shown in Figure 20.13.

Notice that both rtm\_addrs and rtm\_inits are bitmasks in the message from the process, the former specifying the socket address structures that follow, and the latter specifying which metrics are to be initialized. Socket address structures whose bits don't appear in rtm\_addrs don't even appear in the routing message, to save space. But the entire rt\_metrics array always appears in the fixed-length rt\_msghdr structure-elements in the array whose bits are not set in rtm\_inits are ignored.

Figure 20.16 shows the rt\_setmetrics function.

The which argument is always the rtm\_inits member of the routing message 314-318 from the process. in points to the rt\_metrics structure from the process, and out points to the rt\_metrics structure in the routing table entry that is being created or modified.

Each of the 8 bits in the bitmask is tested and if set, the corresponding metric is 319-329 copied. Notice that when a new routing table entry is being created with the RTM\_ADD command, route\_output calls rtrequest, which sets the entire routing table entry to 0 (Figure 19.9). Hence, any metrics not specified by the process in the routing message default to 0.

662 Routing Sockets

Chapter 20

rtsock.c

rtsock.c

raw\_usrreq.c

\$

金の等意語

314 void
315 rt_setmetrics(which, in, out)
316 u_long which;
317 struct rt_metrics *in, *out;
318 (
319 #define metric(f, e) if (which & (f)) out->e = in->e;
<pre>320 metric(RTV_RPIPE, rmx_recvpipe);</pre>
<pre>321 metric(RTV_SPIPE, rmx_sendpipe);</pre>
<pre>322 metric(RTV_SSTHRESH, rmx_ssthresh);</pre>
<pre>323 metric(RTV_RTT, rmx_rtt);</pre>
<pre>324 metric(RTV_RTTVAR, rmx_rttvar);</pre>
<pre>325 metric(RTV_HOPCOUNT, rmx_hopcount);</pre>
<pre>326 metric(RTV_MTU, rmx_mtu);</pre>
<pre>327 metric(RTV_EXPIRE, rmx_expire);</pre>
328 #undef metric
329 }

Figure 20.16 rt\_setmetrics function: set elements of the rt\_metrics structure.

## 20.8 raw\_input Function

All routing messages destined for a process—those that originate from within the kernel and those that originate from a process—are given to raw\_input, which selects the processes to receive the message. Figure 18.11 summarizes the four functions that call raw\_input.

When a routing socket is created, the family is always PF\_ROUTE and the protocol, the third argument to socket, can be 0, which means the process wants to receive all routing messages, or a value such as AF\_INET, which restricts the socket to messages containing addresses of that specific protocol family. A routing control block is created for each routing socket (Section 20.3) and these two values are stored in the sp\_family and sp\_protocol members of the rcb\_proto structure.

Figure 20.17 shows the raw\_input function.

```
51 void
52 raw_input(m0, proto, src, dst)
53 struct mbuf *m0;
54 struct sockproto *proto;
55 struct sockaddr *src, *dst;
56 {
57 struct rawcb *rp;
58 struct mbuf *m = m0;
59 int sockets = 0;
60 struct socket *last;
```

Section 20	0
Section ZU	.o

er 20

ock.c

:ock.c

ker-3 the call

col, e all ages ated iily

rreq.c

ģ,

663

	61 last = 0;
	62 for (rp = rawcb.rcb_next; rp != &rawcb rp = rp->rcb_next) {
	63 if (rp->rcb_proto.sp_family != proto->sp_family)
	$64 \qquad \text{continue};$
	65 if (rp->rcb_proto.sp_protocol &&
	66 rp->rcb_proto.sp_protocol != proto->sp_protocol)
	67 continue;
	68 /*
	69 * We assume the lower level routines have
	70 * placed the address in a canonical format
	71 * suitable for a structure comparison.
	72 *
	73 * Note that if the lengths are not the same
	74 * the comparison will fail at the first byte.
	75 */
	75 #define equal(a1, a2) \
	77 $(\text{bcmp}((\text{caddr_t})(a1), (\text{caddr_t})(a2), a1->sa_len) == 0)$
	77 (bemp((caddr_t)(dr), (caddr_c)(dr), is real addr, dst)) 78 if (rp->rcb_laddr && !equal(rp->rcb_laddr, dst))
	79 continue;
	80 if (rp->rcb_faddr && !equal(rp->rcb_faddr, src))
	80 II (IP )ICD_Idddf dd Ioquar(IP) 81 continue;
250	82 if (last) {
	83 struct mbuf *n;
	84 if (n = m_copy(m, 0, (int) M_COPYALL)) {
	85 if (sbappendaddr(&last->so_rcv, src,
	86 n, (struct mbuf *) 0) == 0)
	87 /* should notify about lost packet */
	88 m_freem(n);
	89 else {
an ann an Anna an Anna Anna Anna Anna A	90 sorwakeup(last);
	91 sockets++;
	92 }
	93 }
1996) 1996) 1986)	94 )
1972)) 1985)	95 last = rp->rcb_socket;
. Martin	96 }
1993) 1988	90 / 97 if (last) {
	98 if (sbappendaddr(&last->so_rcv, src,
2	m, (struct mbuf *) 0) == 0
Salary Salary	
	100 m_freem(m); 101 else {
19.434 19.4	
lé l	
	104 }
	105 } else
	106 m_freem(m);
I.	107 } raw_usrr.

Figure 20.17 raw\_input function: pass routing messages to 0 or more processes.

10.0

In all four calls to raw\_input that we've seen, the proto, src, and dst arguments are pointers to the three globals route\_proto, route\_src, and route\_dst, which are declared and initialized as shown with Figure 19.26.

## Compare address family and protocol

<sup>62–67</sup> The for loop goes through every routing control block checking for a match. The family in the control block (normally PF\_ROUTE) must match the family in the sockproto structure or the control block is skipped. Next, if the protocol in the control block (the third argument to socket) is nonzero, it must match the family in the sockproto structure, or the message is skipped. Hence a process that creates a routing socket with a protocol of 0 receives all routing messages.

## Compare local and foreign addresses

<sup>68–81</sup> These two tests compare the local address in the control block and the foreign address in the control block, if specified. Currently the process is unable to set the rcb\_laddr or rcb\_faddr members of the control block. Normally a process would set the former with bind and the latter with connect, but that is not possible with routing sockets in Net/3. Instead, we'll see that route\_usrreq permanently connects the socket to the route\_src socket address structure, which is OK since that is always the src argument to this function.

## Append message to socket receive buffer

<sup>82-107</sup> If last is nonnull, it points to the most recently seen socket structure that should receive this message. If this variable is nonnull, a copy of the message is appended to that socket's receive buffer by m\_copy and sbappendaddr, and any processes waiting on this receive buffer are awakened. Then last is set to point to this socket that just matched the previous tests. The use of last is to avoid calling m\_copy (an expensive operation) if only one process is to receive the message.

If *N* processes are to receive the message, the first N - 1 receive a copy and the final one receives the message itself.

The variable sockets that is incremented within this function is not used. Since it is incremented only when a message is passed to a process, if it is 0 at the end of the function it indicates that no process received the message (but the value isn't stored anywhere).

## 20.9 route\_usrreg Function

route\_usrreq is the routing protocol's user-request function. It is called for a variety of operations. Figure 20.18 shows the function.

```
64 int
65 route_usrreq(so, req, m, nam, control)
66 struct socket *so;
67 int req;
68 struct mbuf *m, *nam, *control;
69 {
```

rtsock.c

Section 20.9

20

nts

ich

The

the

rol

the

ing

70 int error = 0;struct rawcb \*rp = sotorawcb(so); 71 72 int s; 73 if (req == PRU\_ATTACH) { MALLOC(rp, struct rawcb \*, sizeof(\*rp), M\_PCB, M\_WAITOK); 74 if (so->so\_pcb = (caddr\_t) rp) 75 bzero(so->so\_pcb, sizeof(\*rp)); 76 77 } if (req == PRU\_DETACH && rp) { 78 af = rp->rcb\_proto.sp\_protocol; 79 int if (af == AF\_INET) 80 route\_cb.ip\_count--; 81 else if (af == AF\_NS) 82 route\_cb.ns\_count--; 83 else if (af == AF\_ISO) 84 ign route\_cb.iso\_count--; 85 the route\_cb.any\_count--; 86 uld 87 } *ith* s = splnet(); 88 error = raw\_usrreq(so, req, m, nam, control); ects 89 rp = sotorawcb(so); 90 ays if (req == PRU\_ATTACH && rp) { 91 af = rp->rcb\_proto.sp\_protocol; 92 int if (error) { 93 free((caddr\_t) rp, M\_PCB); 94 uld 95 splx(s); 1 to return (error); 96 ting 97 } if (af == AF\_INET) just 98 route\_cb.ip\_count++; 99 sive else if (af == AF\_NS) 100 route\_cb.ns\_count++; 101 inal else if (af == AF\_ISO) 102 route\_cb.iso\_count++; 103 route\_cb.any\_count++; ce it 104 the rp->rcb\_faddr = &route\_src; 105 ored soisconnected(so); 106 so->so\_options |= SO\_USELOOPBACK; 107 } 108 109 splx(s); return (error); 110 111 } rtsock.c Figure 20.18 route\_usrreq function: process PRU\_xxx requests.

ciety

#### :ock.c

## PRU\_ATTACH: allocate control block

The PRU\_ATTACH request is issued when the process calls socket. Memory is allo-64-77 cated for a routing control block. The pointer returned by MALLOC is stored in the so\_pcb member of the socket structure, and if the memory was allocated, the rawcb structure is set to 0.

INTEL Ex.1013.691

おおに 間

#### PRU\_DETACH: decrement counters

78-87 The close system call issues the PRU\_DETACH request. If the socket structure points to a protocol control block, two of the counters in the route\_cb structure are decremented: one is the any\_count and one is based on the protocol.

#### Process request

88–90 The function raw\_usrreq is called to process the PRU\_xxx request further.

#### Increment counters

<sup>91-104</sup> If the request is PRU\_ATTACH and the socket points to a routing control block, a check is made for an error from raw\_usrreq. Two of the counters in the route\_cb structure are then incremented: one is the any\_count and one is based on the protocol.

#### Connect socket

<sup>105–106</sup> The foreign address in the routing control block is set to route\_src. This permanently connects the new socket to receive routing messages from the PF\_ROUTE family.

#### Enable SO\_USELOOPBACK by default

<sup>107-111</sup> The SO\_USELOOPBACK socket option is enabled. This is a socket option that defaults to being enabled—all others default to being disabled.

## 20.10 raw\_usrreq Function

raw\_usrreq performs most of the processing for the user request in the routing domain. It was called by route\_usrreq in the previous section. The reason the user-request processing is divided between these two functions is that other protocols (e.g., the OSI CLNP) call raw\_usrreq but not route\_usrreq. raw\_usrreq is not intended to be the pr\_usrreq function for a protocol. Instead it is a common subroutine called by the various pr\_usrreq functions.

Figure 20.19 shows the beginning and end of the raw\_usrreq function. The body of the switch is discussed in separate figures following this figure.

#### PRU\_CONTROL requests invalid

119-129 The PRU\_CONTROL request is from the ioctl system call and is not supported in the routing domain.

#### Control information invalid

<sup>130–133</sup> If control information was passed by the process (using the sendmsg system call) an error is returned, since the routing domain doesn't use this optional information.

#### Socket must have a control block

- 134-137 If the socket structure doesn't point to a routing control block, an error is returned. If a new socket is being created, it is the caller's responsibility (i.e., route\_usrreq) to allocate this control block and store the pointer in the so\_pcb member before calling this function.
- 262-269 The default for this switch catches two requests that are not handled by case statements: PRU\_BIND and PRU\_CONNECT. The code for these two requests is present but commented out in Net/3. Therefore issuing the bind or connect system calls on a

1

1:

Section 20.10

apter 20 — raw\_usrreq.c 119 int 120 raw\_usrreq(so, req, m, nam, control) ructure ure are 121 struct socket \*so; 122 int req; 123 struct mbuf \*m, \*nam, \*control; 124 { struct rawcb \*rp = sotorawcb(so); 125 int error = 0; 126 len; 127 int block, a if (req == PRU\_CONTROL) 128 return (EOPNOTSUPP); ite\_cb 129 if (control && control->m\_len) { rotocol. 130 error = EOPNOTSUPP; 131 132 goto release; 133 perma-} if (rp == 0) { 134 family. error = EINVAL; 135 goto release; 136 137 } on that switch (req) { 138 /\* switch cases \*/ 262 default: routing panic("raw\_usrreg"); 263 :he user-} 264 ols (e.g., release: 265 1 is not 266 if (m != NULL) m\_freem(m); subrou-267 return (error); 268 — raw\_usrreq.c 269 } 'he body Figure 20.19 Body of raw\_usrreq function. ported in routing socket causes a kernel panic. This is a bug. Fortunately it requires a superuser process to create this type of socket.

for the PRU\_ATTACH and PRU\_DETACH requests.

must be created by a superuser process.

We now discuss the individual case statements. Figure 20.20 shows the processing

The PRU\_ATTACH request is a result of the socket system call. A routing socket

The function raw\_attach (Figure 20.24) links the control block into the doubly

The PRU\_DETACH is issued by the close system call. The test of a null rp pointer

raw\_detach (Figure 20.25) removes the control block from the doubly linked list.

linked list. The nam argument is the third argument to socket and gets stored in the

is superfluous, since the test was already done before the switch statement.

tem call) tion.

error is ility (i.e., so\_pcb 139-148

149-150

151-159

160-161

control block.

by case is present calls on a

Routing Sockets         139       /*         140       * Allocate a raw control block and         141       * necessary info to allow packets         142       * the appropriate raw interface ro         143       */         144       case PRU_ATTACH:         145       if ((so->so_state & SS_PRIV) == 0)         146       error = EACCES;         147       break;         148       )         149       error = raw_attach(so, (int) nam);	to be foulded to nutine.		186–188 189–196 197–202
<pre>140 * Allocate a raw control block and 140 * necessary info to allow packets 142 * the appropriate raw interface ro 143 */ 144 case PRU_ATTACH: 145 if ((so-&gt;so_state &amp; SS_PRIV) == 0) 146 error = EACCES; 147 break; 148 }</pre>	fill in the to be routed to putine.		189–196 197–202
<pre>140 * Allocate a raw control block and 140 * necessary info to allow packets 141 * the appropriate raw interface ro 143 */ 144 case PRU_ATTACH: 145 if ((so-&gt;so_state &amp; SS_PRIV) == 0) 146 error = EACCES; 147 break; 148 }</pre>	fill in the to be routed to putine.		189–196 197–202
<pre>140 * Allocate a raw control block and 140 * necessary info to allow packets 142 * the appropriate raw interface ro 143 */ 144 case PRU_ATTACH: 145 if ((so-&gt;so_state &amp; SS_PRIV) == 0) 146 error = EACCES; 147 break; 148 }</pre>	utine.		189–196 197–202
<pre>140 * Allocate a raw control block and 141 * necessary info to allow packets 142 * the appropriate raw interface ro 143 */ 144 case PRU_ATTACH: 145 if ((so-&gt;so_state &amp; SS_PRIV) == 0) 146 error = EACCES; 147 break; 148 }</pre>	utine.		197–202
<pre>141 * necessary info to allow packets 142 * the appropriate raw interface ro 143 */ 144 case PRU_ATTACH: 145 if ((so-&gt;so_state &amp; SS_PRIV) == 0) 146 error = EACCES; 147 break; 148 }</pre>	utine.		197–202
<pre>142 * the appropriate raw interface ro 143 */ 144 case PRU_ATTACH: 145 if ((so-&gt;so_state &amp; SS_PRIV) == 0) 146 error = EACCES; 147 break; 148 }</pre>	utine.		197–202
<pre>143 */ 144 case PRU_ATTACH: 145 if ((so-&gt;so_state &amp; SS_PRIV) == 0) 146 error = EACCES; 147 break; 148 }</pre>			197–202
<pre>144 case PRU_ATTACH: 145 if ((so-&gt;so_state &amp; SS_PRIV) == 0) 146 error = EACCES; 147 break; 148 }</pre>	{		197–202
145 if ((so->so_state & SS_PRIV) == 0) 146 error = EACCES; 147 break;	{		197-202
146 error = EACCES; 147 break;			
147 break;			
148 }			
149 $\operatorname{error} = \operatorname{raw}_{\operatorname{attach}(SO, (IIIC), \operatorname{ham})}$		· 1	
150 break;			
151 /*	t deallocation.		
151 /* 152 * Destroy state just before socke	the options.		
152 * Flush data or not depending on			
154 */			
155 case PRU_DETACH:		÷	
156 if (rp == 0) {		- 12 C	
157 error = ENOTCONN;		2	
158 break;			
159 }		1.1	
160 raw_detach(rp);		1	
161 break;	raw_usrreq.c	41	
Figure 20.20 raw_usrreq function: PRU_ATT			
Figure 20.20 raw_usrreq function. FRO_AT			

	raw_usrreq.c		2
186	case PRU_CONNECT2:		2
187	error = EOPNOTSUPP;		_
188	goto release;		2
189	case PRU_DISCONNECT:		2
190	if (rp->rcb_faddr == 0) {		2
191	error = ENOTCONN;		2
192	break;		
193	}		2
194	raw_disconnect(rp);	1123 364	2
195	soisdisconnected(so);	-	2
196	break;		2
107	/*		2
197	/^ * Mark the connection as being incapable of further input.		
198	*/		
199	,		
200	case PRU_SHUTDOWN:		203-217
201 202	socantsendmore(so); break;		I

Figure 20.21 raw\_usrreq function: PRU\_CONNECT2, PRU\_DISCONNECT, and PRU\_SHUTDOWN requests.

r

186–188

usrreq.c

upter 20

The PRU\_CONNECT2 request is from the socketpair system call and is not supported in the routing domain.

Since a routing socket is always connected (Figure 20.18), the PRU\_DISCONNECT 189-196 request is issued by close before the PRU\_DETACH request. The socket must already be connected to a foreign address, which is always true for a routing socket. raw\_disconnect and soisdisconnected complete the processing.

The PRU SHUTDOWN request is from the shutdown system call when the argument 197-202 specifies that no more writes will be performed on the socket. socantsendmore disables further writes.

The most common request for a routing socket, PRU\_SEND, and the PRU\_ABORT and PRU\_SENSE requests are shown in Figure 20.22.

If a nam argument is specified, that is, the process specified a destination address using

either sendto or sendmsg, an error is returned because route\_usrreq always sets

	/* raw_usrreq.c
203 204	* Ship a packet out. The appropriate raw output
204	<ul> <li>* Ship a packet out. The appropriate faw output</li> <li>* routine handles any massaging necessary.</li> </ul>
T. * .	* Fourthe handles any massaging necessary.
206	
207	case PRU_SEND:
208	if (nam) {
209	if (rp->rcb_faddr) {
210	error = EISCONN;
211	break;
212	}
213	<pre>rp-&gt;rcb_faddr = mtod(nam, struct sockaddr *);</pre>
214	<pre>} else if (rp-&gt;rcb_faddr == 0) {</pre>
215	error = ENOTCONN;
216	break;
217	}
218	error = (*so->so_proto->pr_output) (m, so);
219	m = NULL;
220	if (nam)
221	rp->rcb_faddr = 0;
222	break;
223	case PRU_ABORT:
224	raw_disconnect(rp);
225	sofree(so);
226	<pre>soisdisconnected(so);</pre>
227	break;
228	case PRU_SENSE:
229	/*
230	* stat: don't bother with a blocksize.
	*/
231	

\_usrreq.c

rcb\_faddr for a routing socket.

quests.

670 Routing Sockets

The message in the mbuf chain pointed to by m is passed to the protocol's 218-222 pr\_output function, which is route\_output.

If a PRU\_ABORT request is issued, the control block is disconnected, the socket is 223-227 released, and the socket is disconnected.

The PRU\_SENSE request is issued by the fstat system call. The function returns 228-232 OK.

Figure 20.23 shows the remaining PRU\_xxx requests.

1 18		raw_usrreq.c
33	/*	
34	* Not supported.	
35	* /	
36	case PRU_RCVOOB:	
37	case PRU_RCVD:	
238	return (EOPNOTSUPP);	
239	case PRU_LISTEN:	
240	case PRU_ACCEPT:	
241	case PRU_SENDOOB:	
242	error = EOPNOTSUPP;	
243	break;	
244	case PRU_SOCKADDR:	
245	if (rp->rcb_laddr == 0) {	
246	error = EINVAL;	
247	break;	
248	}	
249	<pre>len = rp-&gt;rcb_laddr-&gt;sa_len;</pre>	(unsigned) len);
250	len = rp->rcb_laddr->sa_ten; bcopy((caddr_t) rp->rcb_laddr, mtod(nam, caddr_t)	, (010-5
251	nam->m_len = len;	2 1
252	break;	5 1
253	case PRU_PEERADDR:	1 
254	if (rp->rcb_faddr == 0) {	
255	error = ENOTCONN;	24
256	break;	5
257	}	
258	<pre>len = rp-&gt;rcb_faddr-&gt;sa_len; len = rp-&gt;rcb_faddr-&gt;sa_len;</pre>	) (unsigned) len);
259	len = rp->rcb_iaddr->sa_ien, bcopy((caddr_t) rp->rcb_faddr, mtod(nam, caddr_t)	,, (
260	nam->m_len = len;	- 
261	break;	raw_usrreq.c

Figure 20.23 raw\_usrreq function: final part.

These five requests are not supported. 233-243

The PRU\_SOCKADDR and PRU\_PEERADDR requests are from the getsockname and 244-261 getpeername system calls respectively. The former always returns an error, since the bind system call, which sets the local address, is not supported in the routing domain. The latter always returns the contents of the socket address structure route\_src, which was set by route\_usrreq as the foreign address.

Sectic

20.1

49-

65-

68-

75-

88

纝

Chapter 20

# 20.11 raw\_attach, raw\_detach, and raw\_disconnect Functions

cket is

tocol's

pter 20

eturns

usrreq.c

len);

len);

,\_usrreq.c

ame and ince the domain. te\_src, The raw\_attach function, shown in Figure 20.24, was called by raw\_input to finish processing the PRU\_ATTACH request.

raw\_cb.c 49 int 50 raw\_attach(so, proto) 51 struct socket \*so; proto; 52 int 53 { struct rawcb \*rp = sotorawcb(so); 54 55 int error; 56 \* It is assumed that raw\_attach is called 57 \* after space has been allocated for the 58 59 \* rawcb. 60 \* / if (rp == 0)61 return (ENOBUFS); 62 if (error = soreserve(so, raw\_sendspace, raw\_recvspace)) 63 64 return (error); rp->rcb\_socket = so; 65 rp->rcb\_proto.sp\_family = so->so\_proto->pr\_domain->dom\_family; 66 rp->rcb\_proto.sp\_protocol = proto; 67 insque(rp, &rawcb); 68 69 return (0); 70 } - raw\_cb.c

Figure 20.24 raw\_attach function.

<sup>49-64</sup> The caller must have already allocated the raw protocol control block. soreserve sets the high-water marks for the send and receive buffers to 8192. This should be more than adequate for the routing messages.

A pointer to the socket structure is stored in the protocol control block along with the dom\_family (which is PF\_ROUTE from Figure 20.1 for the routing domain) and the proto argument (which is the third argument to socket).

<sup>68-70</sup> insque adds the control block to the front of the doubly linked list headed by the global rawcb.

The raw\_detach function, shown in Figure 20.25, was called by raw\_input to finish processing the PRU\_DETACH request.

The so\_pcb pointer in the socket structure is set to null and the socket is released. The control block is removed from the doubly linked list by remque and the memory used for the control block is released by free.

The raw\_disconnect function, shown in Figure 20.26, was called by raw\_input to process the PRU\_DISCONNECT and PRU\_ABORT requests.

<sup>88-94</sup> If the socket does not reference a descriptor, raw\_detach releases the socket and control block.

672 Routing Sockets

Chapter 20

raw\_cb.c

raw\_cb.c

raw\_cb.c

raw cb.c

75 void 76 raw\_detach(rp) 77 struct rawcb \*rp; 78 { 79 struct socket \*so = rp->rcb\_socket; 80 so->so\_pcb = 0; 81 sofree(so); 82 remque(rp); 83 free((caddr\_t) (rp), M\_PCB); 84 }

Figure 20.25 raw\_detach function.

88 void
89 raw\_disconnect(rp)
90 struct rawcb \*rp;
91 {
92 if (rp->rcb\_socket->so\_state & SS\_NOFDREF)
93 raw\_detach(rp);
94 }

Figure 20.26 raw\_disconnect function.

#### 20.12 Summary

A routing socket is a raw socket in the PF\_ROUTE domain. Routing sockets can be created only by a superuser process. If a nonprivileged process wants to read the routing information contained in the kernel, the sysct1 system call supported by the routing domain can be used (we described this in the previous chapter).

This chapter was our first encounter with the protocol control blocks (PCBs) that are normally associated with each socket. In the routing domain a special rawcb contains information about the routing socket: the local and foreign addresses, the address family, and the protocol. We'll see in Chapter 22 that the larger Internet protocol control block (inpcb) is used with UDP, TCP, and raw IP sockets. The concepts are the same, however: the socket structure is used by the socket layer, and the PCB, a rawcb or an inpcb, is used by the protocol layer. The socket structure points to the PCB and vice versa.

The route\_output function handles the five routing requests that can be issued by a process. raw\_input delivers a routing message to one or more routing sockets, depending on the protocol and address family. The various PRU\_xxx requests for a routing socket are handled by raw\_usrreq and route\_usrreq. In later chapters we'll encounter additional xxx\_usrreq functions, one per protocol (UDP, TCP, and raw IP), each consisting of a switch statement to handle each request.

# **Exercises**

Chapter 20

- **20.1** List two ways a process can receive the return value from route\_output when the process writes a message to a routing socket. Which method is more reliable?
- 20.2 What happens when a process specifies a nonzero *protocol* argument to the socket system call, since the pr\_protocol member of the routesw structure is 0?
- **20.3** Routes in the routing table (other than ARP entries) never time out. Implement a timeout on routes.

20

:b.c

cb.c

cb.c

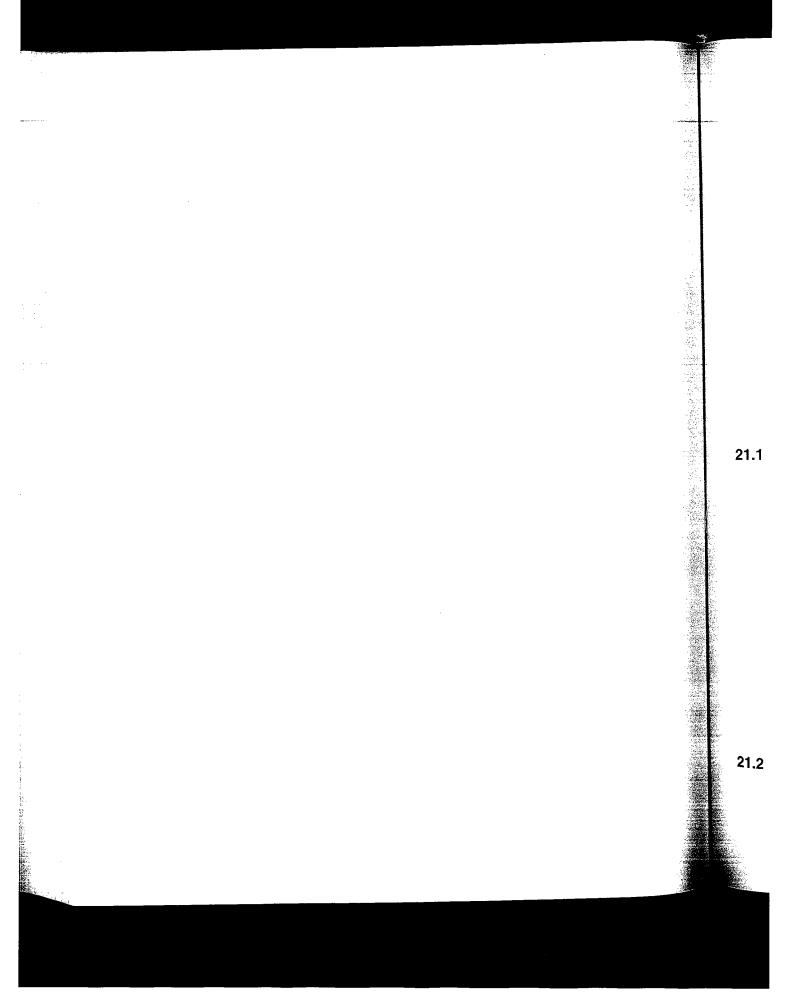
reing ing are ins umtrol

me, • an *r*ice

ıed ets, r a ters :aw

INTEL Ex.1013.699

žv E



# 21

# ARP: Address Resolution Protocol

# 21.1 Introduction

ARP, the Address Resolution Protocol, handles the translation of 32-bit IP addresses into the corresponding hardware address. For an Ethernet, the hardware addresses are 48-bit Ethernet addresses. In this chapter we only consider mapping IP addresses into 48-bit Ethernet addresses, although ARP is more general and can work with other types of data links. ARP is specified in RFC 826 [Plummer 1982].

When a host has an IP datagram to send to another host on a locally attached Ethernet, the local host first looks up the destination host in the *ARP cache*, a table that maps a 32-bit IP address into its corresponding 48-bit Ethernet address. If the entry is found for the destination, the corresponding Ethernet address is copied into the Ethernet header and the datagram is added to the appropriate interface's output queue. If the entry is not found, the ARP functions hold onto the IP datagram, broadcast an ARP request asking the destination host for its Ethernet address, and, when a reply is received, send the datagram to its destination.

This simple overview handles the common case, but there are many details that we describe in this chapter as we examine the Net/3 implementation of ARP. Chapter 4 of Volume 1 contains additional ARP examples.

# 21.2 ARP and the Routing Table

The Net/3 implementation of ARP is tied to the routing table, which is why we postponed discussing ARP until we had described the structure of the Net/3 routing tables. Figure 21.1 shows an example that we use in this chapter when describing ARP.



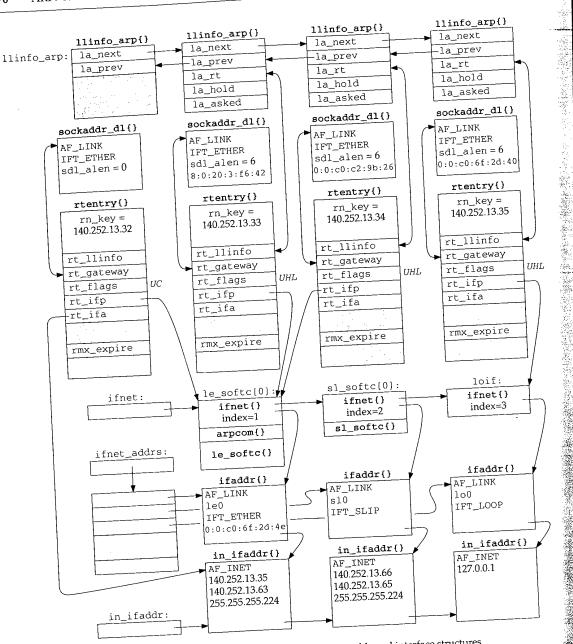


Figure 21.1 Relationship of ARP to routing table and interface structures.

The entire figure corresponds to the example network used throughout the text (Figure 1.17). It shows the ARP entries on the system bsdi. The ifnet, ifaddr, and in\_ifaddr structures are simplified from Figures 3.32 and 6.5. We have removed some of the details from these three structures, which were covered in Chapters 3 and 6.

For example, we don't show the two sockaddr\_dl structures that appear after each ifaddr structure—instead we summarize the information contained in these two structures. Similarly, we summarize the information contained in the three in\_ifaddr structures.

We briefly summarize some relevant points from this figure, the details of which we cover as we proceed through the chapter.

- 1. A doubly linked list of llinfo\_arp structures contains a minimal amount of information for each hardware address known by ARP. The global llinfo\_arp is the head of this list. Not shown in this figure is that the la\_prev pointer of the first entry points to the last entry, and the la\_next pointer of the last entry points to the first entry. This linked list is processed by the ARP timer function every 5 minutes.
- 2. For each IP address with a known hardware address, a routing table entry exists (an rtentry structure). The llinfo\_arp structure points to the corresponding rtentry structure, and vice versa, using the la\_rt and rt\_llinfo pointers. The three routing table entries in this figure with an associated llinfo\_arp structure are for the hosts sun (140.252.13.33), svr4 (140.252.13.34), and bsdi itself (140.252.13.35). These three are also shown in Figure 18.2.
- 3. We show a fourth routing table entry on the left, without an llinfo\_arp structure, which is the entry for the network route to the local Ethernet (140.252.13.32). We show its rt\_flags with the C bit on, since this entry is cloned to form the other three routing table entries. This entry is created by the call to rtinit when the IP address is assigned to the interface by in\_ifinit (Figure 6.19). The other three entries are host entries (the H flag) and are generated by ARP (the L flag) when a datagram is sent to that IP address.
- 4. The rt\_gateway member of the rtentry structure points to a sockaddr\_dl structure. This data-link socket address structure contains the hardware address if the sdl\_alen member equals 6.
- 5. The rt\_ifp member of the routing table entry points to the ifnet structure of the outgoing interface. Notice that the two routing table entries in the middle, for other hosts on the local Ethernet, both point to le\_softc[0], but the routing table entry on the right, for the host bsdi itself, points to the loopback structure. Since rt\_ifp.if\_output (Figure 8.25) points to the output routine, packets sent to the local IP address are routed to the loopback interface.
- 6. Each routing table entry also points to the corresponding in\_ifaddr structure. (Actually the rt\_ifa member points to an ifaddr structure, but recall from Figure 6.8 that the first member of an in\_ifaddr structure is an ifaddr structure.) We show only one of these pointers in the figure, although all four point to the same structure. Remember that a single interface, say 1e0, can have multiple IP addresses, each with its own in\_ifaddr structure, which is why the rt\_ifa pointer is required in addition to the rt\_ifp pointer.

Chapter 21

rp{}

\_d1 { }

R

= 6

y{}

13.35

Εo

з

√ау

ire

£:

t{} <=3

1r{}

iddr{}

)P

UHL

:2d:40



7. The la\_hold member is a pointer to an mbuf chain. An ARP request is broadcast because a datagram is sent to that IP address. While the kernel awaits the ARP reply it holds onto the mbuf chain for the datagram by storing its address in la\_hold. When the ARP reply is received, the mbuf chain pointed to by la\_hold is sent.

Chapter 21

金属のないないない

8. Finally, we show the variable rmx\_expire, which is in the rt\_metrics structure within the routing table entry. This value is the timer associated with each ARP entry. Some time after an ARP entry has been created (normally 20 minutes) the ARP entry is deleted.

Even though major routing table changes took place with 4.3BSD Reno, the ARP cache was left alone with 4.3BSD Reno and Net/2. 4.4BSD, however, removed the stand-alone ARP cache and moved the ARP information into the routing table.

The ARP table in Net/2 was an array of structures composed of the following members: an IP address, an Ethernet address, a timer, flags, and a pointer to an mbuf (similar to the la\_hold member in Figure 21.1). We see with Net/3 that the same information is now spread throughout multiple structures, all of which are linked.

# 21.3 Code Introduction

There are nine ARP functions in a single C file and definitions in two headers, as shown in Figure 21.2.

File	Description
net/if_arp.h netinet/if_ether.h	arphdr structure definition various structure and constant definitions
netinet/if_ether.c	ARP functions

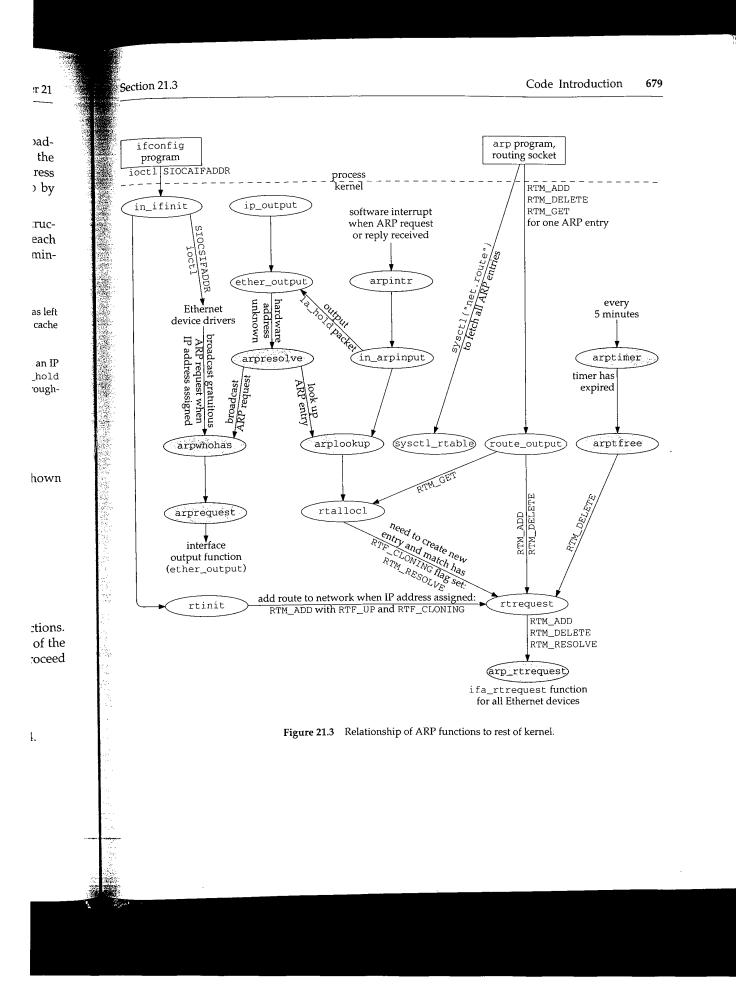
Figure 21.2 Files discussed in this chapter.

Figure 21.3 shows the relationship of the ARP functions to other kernel functions. In this figure we also show the relationship between the ARP functions and some of the routing functions from Chapter 19. We describe all these relationships as we proceed through the chapter.

# **Global Variables**

Ten global variables are introduced in this chapter, which are shown in Figure 21.4.

の方法である。



Chapter 21

Variable	Datatype	Description
llinfo_arp	struct llinfo_arp	head of llinfo_arp doubly linked list (Figure 21.1)
arpintrg	struct ifqueue	ARP input queue from Ethernet device drivers (Figure 4.9)
arpt_prune arpt_keep arpt_down	int int int	#seconds between checking ARP list (5 × 60) #seconds ARP entry valid once resolved (20 × 60) #seconds between ARP flooding algorithm (20)
arp_inuse arp_allocated arp_maxtries arpinit_done useloopback	int int int int int	#ARP entries currently in use #ARP entries ever allocated max #tries for an IP address before pausing (5) initialization-performed flag use loopback for local host (default true)

Figure 21.4 Global variables introduced in this chapter.

#### Statistics

The only statistics maintained by ARP are the two globals arp\_inuse and arp\_allocated, from Figure 21.4. The former counts the number of ARP entries currently in use and the latter counts the total number of ARP entries allocated since the system was initialized. Neither counter is output by the netstat program, but they can be examined with a debugger.

The entire ARP cache can be listed using the arp -a command, which uses the sysctl system call with the arguments shown in Figure 19.36. Figure 21.5 shows the output from this command, for the entries shown in Figure 18.2.

```
bsdi $ arp -a
sun.tuc.noao.edu (140.252.13.33) at 8:0:20:3:f6:42
svr4.tuc.noao.edu (140.252.13.34) at 0:0:c0:c2:9b:26
bsdi.tuc.noao.edu (140.252.13.35) at 0:0:c0:6f:2d:40 permanent
ALL-SYSTEMS.MCAST.NET (224.0.0.1) at (incomplete)
```

Figure 21.5 arp -a output corresponding to Figure 18.2.

Since the multicast group 224.0.0.1 has the L flag set in Figure 18.2, and since the arp program looks for entries with the RTF\_LLINFO flag set, the multicast groups are output by the program. Later in this chapter we'll see why this entry is marked as "incomplete" and why the entry above it is "permanent."

#### **SNMP Variables**

As described in Section 25.8 of Volume 1, the original SNMP MIB defined an address translation group that was the system's ARP cache. MIB-II deprecated this group and instead each network protocol group (i.e., IP) contains its own address translation tables. Notice that the change in Net/2 to Net/3 from a stand-alone ARP table to an integration of the ARP information within the IP routing table parallels this SNMP change.



hapter 21

use and ntries cursince the , but they

ι uses the shows the

ce the arp ps are outas "incom-

an address group and translation table to an this SNMP Figure 21.6 shows the IP address translation table from MIB-II, named ipNetToMediaTable. The values returned by SNMP for this table are taken from the routing table entry and its corresponding ifnet structure.

Name	Member	FoMediaIfIndex > << ipNetToMediaNetAddress > Description
ipNetToMediaIfIndex ipNetToMediaPhysAddress ipNetToMediaNetAddress ipNetToMediaType	if_index rt_gateway rt_key rt_flags	corresponding interface: ifIndex physical address IP address type of mapping: 1 = other, 2 = invalidated, 3 = dynamic, 4 = static (see text)

Figure 21.6	IP address translation table: ipNetToMediaTable.	
-------------	--	--

If the routing table entry has an expiration time of 0 it is considered permanent and hence "static." Otherwise the entry is considered "dynamic."

# 21.4 ARP Structures

Figure 21.7 shows the format of an ARP packet when transmitted on an Ethernet.

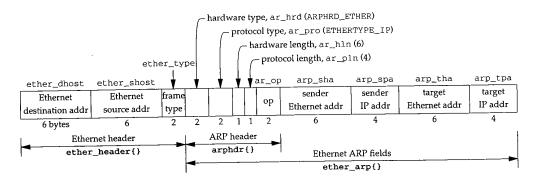


Figure 21.7 Format of an ARP request or reply when used on an Ethernet.

The ether\_header structure (Figure 4.10) defines the 14-byte Ethernet header; the arphdr structure defines the next five fields, which are common to ARP requests and ARP replies on any type of media; and the ether\_arp structure combines the arphdr structure with the sender and target addresses when ARP is used on an Ethernet.

Figure 21.8 shows the definition of the arphdr structure. Figure 21.7 shows the values of the first four fields in this structure when ARP is mapping IP addresses to Ethernet addresses.

Figure 21.9 shows the combination of the arphdr structure with the fields used with IP addresses and Ethernet addresses, forming the ether\_arp structure. Notice that ARP uses the terms *hardware* to describe the 48-bit Ethernet address, and *protocol* to describe the 32-bit IP address.

Chapter 21 ARP: Address Resolution Protocol if arp.h 45 struct arphdr { /\* format of hardware address \*/ u\_short ar\_hrd; 46 /\* format of protocol address \*/ u\_short ar\_pro; 47 /\* length of hardware address \*/ u\_char ar\_hln; 48 /\* length of protocol address \*/ u\_char ar\_pln; 49 /\* ARP/RARP operation, Figure 21.15 \*/ u\_short ar\_op; 50 - if\_arp.h 51 }; Figure 21.8 arphdr structure: common ARP request/reply header. - if\_ether.h 79 struct ether\_arp { /\* fixed-size header \*/ struct arphdr ea\_hdr; 80 /\* sender hardware address \*/ u\_char arp\_sha[6]; 81

682

22 u\_char arp\_spa[4]; /\* sender protocol address \*/
23 u\_char arp\_tha[6]; /\* target hardware address \*/
24 u\_char arp\_tpa[4]; /\* target protocol address \*/
25 };
26 #define arp\_hrd ea\_hdr.ar\_hrd
27 #define arp\_pro ea\_hdr.ar\_pro
28 #define arp\_hl ea\_hdr.ar\_pln
29 #define arp\_op ea\_hdr.ar\_op

Figure 21.9 ether\_arp structure.

One llinfo\_arp structure, shown in Figure 21.10, exists for each ARP entry. Additionally, one of these structures is allocated as a global of the same name and used as the head of the linked list of all these structures. We often refer to this list as the *ARP cache*, since it is the only data structure in Figure 21.1 that has a one-to-one correspondence with the ARP entries.

103 struct llinfo\_arp { struct llinfo\_arp \*la\_next; 104 struct llinfo\_arp \*la\_prev; 105 struct rtentry \*la\_rt; 106 /\* last packet until resolved/timeout \*/ struct mbuf \*la\_hold; 107 /\* #times we've queried for this addr \*/ 108 long la\_asked;  $109 \};$ /\* deletion time in seconds \*/ 110 #define la\_timer la\_rt->rt\_rmx.rmx\_expire if\_ether.h

Figure 21.10 llinfo\_arp structure.

With Net/2 and earlier systems it was easy to identify the structure called the *ARP cache*, since a single structure contained everything for each ARP entry. Since Net/3 stores the ARP information among multiple structures, no single structure can be called the *ARP cache*. Nevertheless, having the concept of an ARP cache, which is the collection of information describing a single ARP entry, simplifies the discussion.

的影响不知

1

if\_ether.h

if\_ether.h

107

110

÷.

ipter 21

if\_arp.h

if\_arp.h

f\_ether.h

if ether.h

? entry.

1d used

he ARP

respon-

if\_ether.h

ıds \*/

if\_ether.h

ache. since

\RP infor-

Neverthescribing a <sup>104-106</sup> The first two entries form the doubly linked list, which is updated by the insque and remque functions. la\_rt points to the associated routing table entry, and the rt\_llinfo member of the routing table entry points to this structure.

When ARP receives an IP datagram to send to another host but the destination's hardware address is not in the ARP cache, an ARP request must be sent and the ARP reply received before the datagram can be sent. While waiting for the reply the mbuf pointer to the datagram is saved in la\_hold. When the ARP reply is received, the packet pointed to by la\_hold (if any) is sent.

<sup>108–109</sup> la\_asked counts how many consecutive times an ARP request has been sent to this IP address without receiving a reply. We'll see in Figure 21.24 that when this counter reaches a limit, that host is considered down and another ARP request won't be sent for a while.

This definition uses the rmx\_expire member of the rt\_metrics structure in the routing table entry as the ARP timer. When the value is 0, the ARP entry is considered permanent. When nonzero, the value is the number of seconds since the Unix Epoch when the entry expires.

# 21.5 arpwhohas Function

The arpwhohas function is normally called by arpresolve to broadcast an ARP request. It is also called by each Ethernet device driver to issue a *gratuitous ARP* request when the IP address is assigned to the interface (the SIOCSIFADDR ioctl in Figure 6.28). Section 4.7 of Volume 1 describes gratuitous ARP—it detects if another host on the Ethernet is using the same IP address and also allows other hosts with ARP entries for this host to update their ARP entry if this host has changed its Ethernet address. arpwhohas simply calls arprequest, shown in the next section, with the correct arguments.

```
if_ether.c
if_eth
```

Figure 21.11 arpwhohas function: broadcast an ARP request.

<sup>196-202</sup> The arpcom structure (Figure 3.26) is common to all Ethernet devices and is part of the le\_softc structure, for example (Figure 3.20). The ac\_ipaddr member is a copy of the interface's IP address, which is set by the driver when the SIOCSIFADDR ioctl is executed (Figure 6.28). ac\_enaddr is the Ethernet address of the device.

The second argument to this function, addr, is the IP address for which the ARP request is being issued: the target IP address. In the case of a gratuitous ARP request, addr equals ac\_ipaddr, so the second and third arguments to arprequest are the same, which means the sender IP address will equal the target IP address in the gratuitous ARP request.

# 21.6 arprequest Function

The arprequest function is called by arpwhohas to broadcast an ARP request. It builds an ARP request packet and passes it to the interface's output function.

Before looking at the source code, let's examine the data structures built by the function. To send the ARP request the interface output function for the Ethernet device (ether\_output) is called. One argument to ether\_output is an mbuf containing the data to send: everything that follows the Ethernet type field in Figure 21.7. Another argument is a socket address structure containing the destination address. Normally this destination address is an IP address (e.g., when ip\_output calls ether\_output in Figure 21.3). For the special case of an ARP request, the sa\_family member of the socket address structure is set to AF\_UNSPEC, which tells ether\_output that it contains a filled-in Ethernet header, including the destination Ethernet address. This prevents ether\_output from calling arpresolve, which would cause an infinite loop. We don't show this loop in Figure 21.3, but the "interface output function" below arprequest is ether\_output. If ether\_output were to call arpresolve again, the infinite loop would occur.

Figure 21.12 shows the mbuf and the socket address structure built by this function. We also show the two pointers eh and ea, which are used in the function.

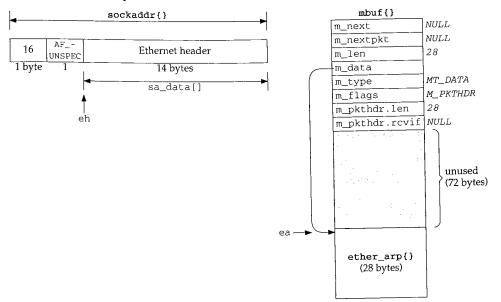


Figure 21.12 sockaddr and mbuf built by arprequest.

Figure 21.13 shows the arprequest function.

1

日本のないで

194

arprequest Function

685

Section 21.6

er 21

∴ It

the

vice

the

ther

ally

put

the

:on-

pre-

юр.

low

ain,

ion.

l es)

if ether.c 209 static void 210 arprequest(ac, sip, tip, enaddr) 211 struct arpcom \*ac; 212 u\_long \*sip, \*tip; 213 u\_char \*enaddr; 214 { struct mbuf \*m; 215 216 struct ether\_header \*eh; 217 struct ether\_arp \*ea; struct sockaddr sa; 218 if ((m = m\_gethdr(M\_DONTWAIT, MT\_DATA)) == NULL) 219 220 return; 221 m->m\_len = sizeof(\*ea); m->m\_pkthdr.len = sizeof(\*ea); 222 MH\_ALIGN(m, sizeof(\*ea)); 223 224 ea = mtod(m, struct ether\_arp \*); eh = (struct ether\_header \*) sa.sa\_data; 225 bzero((caddr\_t) ea, sizeof(\*ea)); 226 bcopy((caddr\_t) etherbroadcastaddr, (caddr\_t) eh->ether\_dhost, 227 228 sizeof(eh->ether\_dhost)); /\* if\_output() will swap \*/ 229 eh->ether\_type = ETHERTYPE\_ARP; ea->arp\_hrd = htons(ARPHRD\_ETHER); 230 ea->arp\_pro = htons(ETHERTYPE\_IP); 231 ea->arp\_hln = sizeof(ea->arp\_sha); /\* hardware address length \*/ 232 /\* protocol address length \*/ 233 ea->arp\_pln = sizeof(ea->arp\_spa); ea->arp\_op = htons(ARPOP\_REQUEST); 234 bcopy((caddr\_t) enaddr, (caddr\_t) ea->arp\_sha, sizeof(ea->arp\_sha)); 235 bcopy((caddr\_t) sip, (caddr\_t) ea->arp\_spa, sizeof(ea->arp\_spa)); 236 bcopy((caddr\_t) tip, (caddr\_t) ea->arp\_tpa, sizeof(ea->arp\_tpa)); 237 sa.sa\_family = AF\_UNSPEC; 238 239 sa.sa\_len = sizeof(sa); (\*ac->ac\_if.if\_output) (&ac->ac\_if, m, &sa, (struct rtentry \*) 0); 240 241 } - if ether.c Figure 21.13 arprequest function: build an ARP request packet and send it.

# Allocate and initialize mbuf

209-223

A packet header mbuf is allocated and the two length fields are set. MH\_ALIGN allows room for a 28-byte ether\_arp structure at the end of the mbuf, and sets the m\_data pointer accordingly. The reason for moving this structure to the end of the mbuf is to allow ether\_output to prepend the 14-byte Ethernet header in the same mbuf.

Secti

21.

į.

#### Initialize pointers

The two pointers ea and eh are set and the ether\_arp structure is set to 0. The 224-226 only purpose of the call to bzero is to set the target hardware address to 0, because the other eight fields in this structure are explicitly set to their respective value.

#### Fill in Ethernet header

- The destination Ethernet address is set to the Ethernet broadcast address and the 227-229 Ethernet type field is set to ETHERTYPE\_ARP. Note the comment that this 2-byte field will be converted from host byte order to network byte order by the interface output function. This function also fills in the Ethernet source address field. Figure 21.14 shows the different values for the Ethernet type field.

Constant	Value	Description
ETHERTYPE_IP	0x0800	IP frames
ETHERTYPE_ARP	0x0806	ARP frames
ETHERTYPE_REVARP	0x8035	reverse ARP (RARP) frames
ETHERTYPE_IPTRAILERS	0x1000	trailer encapsulation (deprecated)

RARP maps an Ethernet address to an IP address and is used when a diskless system bootstraps. RARP is normally not part of the kernel's implementation of TCP/IP, so it is not covered in this text. Chapter 5 of Volume 1 describes RARP.

#### Fill in ARP fields

230-237

All fields in the ether\_arp structure are filled in, except the target hardware address, which is what the ARP request is looking for. The constant ARPHRD\_ETHER, which has a value of 1, specifies the format of the hardware addresses as 6-byte Ethernet addresses. To identify the protocol addresses as 4-byte IP addresses, arp\_pro is set to the Ethernet type field for IP from Figure 21.14. Figure 21.15 shows the various ARP operation codes. We encounter the first two in this chapter. The last two are used with RARP.

Constant	Value	Description
ARPOP_REQUEST	1	ARP request to resolve protocol address
ARPOP_REPLY	2	reply to ARP request
ARPOP_REVREQUEST	3	RARP request to resolve hardware address
ARPOP_REVREPLY	4	reply to RARP request

Figure 21.15 Al	P operation codes.
-----------------	--------------------

# Fill in sockaddr and call interface output function

238-241 The sa\_family member of the socket address structure is set to AF\_UNSPEC and the sa\_len member is set to 16. The interface output function is called, which we said is ether\_output.

319-

arpintr Function

Section 21.7

а

e

e

t

ιt

4

m

is

re

R,

et

to ЗЪ

th

#### arpintr Function 21.7

In Figure 4.13 we saw that when ether\_input receives an Ethernet frame with a type field of ETHERTYPE\_ARP, it schedules a software interrupt of priority NETISR\_ARP and appends the frame to ARP's input queue: arpintrg. When the kernel processes the software interrupt, the function arpintr, shown in Figure 21.16, is called.

- if \_ether.c

687

```
319 void
320 arpintr()
321 {
        struct mbuf *m;
322
        struct arphdr *ar;
323
                 s;
        int
324
        while (arpintrq.ifq_head) {
325
             s = splimp();
326
            IF_DEQUEUE(&arpintrq, m);
327
             splx(s);
328
             if (m == 0 || (m->m_flags & M_PKTHDR) == 0)
329
                 panic("arpintr");
330
             if (m->m_len >= sizeof(struct arphdr) &&
331
                 (ar = mtod(m, struct arphdr *)) &&
332
                 ntohs(ar->ar_hrd) == ARPHRD_ETHER &&
333
                 m->m_len >= sizeof(struct arphdr) + 2*ar->ar_hln + 2*ar->ar_pln)
 334
                      switch (ntohs(ar->ar_pro)) {
 335
                      case ETHERTYPE_IP:
 336
                      case ETHERTYPE_IPTRAILERS:
 337
                          in_arpinput(m);
 338
                          continue;
 339
                      }
 340
             m_freem(m);
 341
          }
 342
                                                                              - if_ether.c
 343 }
```

Figure 21.16 arpintr function: process Ethernet frames containing ARP requests or replies.

The while loop processes one frame at a time, as long as there are frames on the queue. The frame is processed if the hardware type specifies Ethernet addresses, and if 319-343 the size of the frame is greater than or equal to the size of an arphdr structure plus the sizes of two hardware addresses and two protocol addresses. If the type of protocol addresses is either ETHERTYPE\_IP or ETHERTYPE\_IPTRAILERS, the in\_arpinput function, shown in the next section, is called. Otherwise the frame is discarded.

Notice the order of the tests within the if statement. The length is checked twice. First, if the length is at least the size of an arphdr structure, then the fields in that structure can be examined. The length is checked again, using the two length fields in the arphdr structure.

nd ıid AND ALT ALL

Section :

# 21.8 in\_arpinput Function

This function is called by arpintr to process each received ARP request or ARP reply. While ARP is conceptually simple, numerous rules add complexity to the implementation. The following two scenarios are typical:

- 1. If a request is received for one of the host's IP addresses, a reply is sent. This is the normal case of some other host on the Ethernet wanting to send this host a packet. Also, since we're about to receive a packet from that other host, and we'll probably send a reply, an ARP entry is created for that host (if one doesn't already exist) because we have its IP address and hardware address. This optimization avoids another ARP exchange when the packet is received from the other host.
- 2. If a reply is received in response to a request sent by this host, the corresponding ARP entry is now complete (the hardware address is known). The other host's hardware address is stored in the sockaddr\_dl structure and any queued packet for that host can now be sent. Again, this is the normal case.

ARP requests are normally broadcast so each host sees *all* ARP requests on the Ethernet, even those requests for which it is not the target. Recall from arprequest that when a request is sent, it contains the *sender's* IP address and hardware address. This allows the following tests also to occur.

- 3. If some other host sends a request or reply with a sender IP address that equals this host's IP address, one of the two hosts is misconfigured. Net/3 detects this error and logs a message for the administrator. (We say "request or reply" here because in\_arpinput doesn't examine the operation type. But ARP replies are normally unicast, in which case only the target host of the reply receives the reply.)
- 4. If this host receives a request or reply from some other host for which an ARP entry already exists, and if the other host's hardware address has changed, the hardware address in the ARP entry is updated accordingly. This can happen if the other host is shut down and then rebooted with a different Ethernet interface (hence a different hardware address) before its ARP entry times out. The use of this technique, along with the other host sending a gratuitous ARP request when it reboots, prevents this host from being unable to communicate with the other host after the reboot because of an ARP entry that is no longer valid.
- 5. This host can be configured as a *proxy ARP server*. This means it responds to ARP requests for some other host, supplying the other host's hardware address in the reply. The host whose hardware address is supplied in the proxy ARP reply must be one that is able to forward IP datagrams to the host that is the target of the ARP request. Section 4.6 of Volume 1 discusses proxy ARP.

A Net/3 system can be configured as a proxy ARP server. These ARP entries are added with the arp command, specifying the IP address, hardware address,

358-37

376-38

V.

1-

ίs

а

d

.′t

i-

۱e

ıg

's

!d

эt,

. a

ne

ils

uis

re

ire

he

ЗP

he

ı if

er-

he

RP

ate

zer

to

ass

RP

ar-

ies

2SS,

– if\_ether.c

and the keyword pub. We'll see the support for this in Figure 21.20 and we describe it in Section 21.12.

We examine in\_arpinput in four parts. Figure 21.17 shows the first part.

358 static void 359 in\_arpinput(m) 360 struct mbuf \*m; 361 { struct ether\_arp \*ea; 362 struct arpcom \*ac = (struct arpcom \*) m->m\_pkthdr.rcvif; 363 struct ether\_header \*eh; 364 struct llinfo\_arp \*la = 0; 365 struct rtentry \*rt; 366 struct in\_ifaddr \*ia, \*maybe\_ia = 0; 367 struct sockaddr\_dl \*sdl; 368 struct sockaddr sa; 369 struct in\_addr isaddr, itaddr, myaddr; 370 int op; 371 ea = mtod(m, struct ether\_arp \*); 372 op = ntohs(ea->arp\_op); bcopy((caddr\_t) ea->arp\_spa, (caddr\_t) & isaddr, sizeof(isaddr)); 373 bcopy((caddr\_t) ea->arp\_tpa, (caddr\_t) & itaddr, sizeof(itaddr)); 374 375 for (ia = in\_ifaddr; ia; ia = ia->ia\_next) 376 if (ia->ia\_ifp == &ac->ac\_if) { 377 maybe\_ia = ia; 378 if ((itaddr.s\_addr == ia->ia\_addr.sin\_addr.s\_addr) || 379 (isaddr.s\_addr == ia->ia\_addr.sin\_addr.s\_addr)) 380 break; 381 } 382 if (maybe\_ia == 0) 383 goto out; myaddr = ia ? ia->ia\_addr.sin\_addr : maybe\_ia->ia\_addr.sin\_addr; 384 - if\_ether.c 385

Figure 21.17 in\_arpinput function: look for matching interface.

<sup>358–375</sup> The length of the ether\_arp structure was verified by the caller, so ea is set to point to the received packet. The ARP operation (request or reply) is copied into op but it isn't examined until later in the function. The sender's IP address and target IP address are copied into isaddr and itaddr.

# Look for matching interface and IP address

The linked list of Internet addresses for the host is scanned (the list of in\_ifaddr structures, Figure 6.5). Remember that a given interface can have multiple IP addresses. Since the received packet contains a pointer (in the mbuf packet header) to the receiving interface's ifnet structure, the only IP addresses considered in the for loop are those associated with the receiving interface. If either the target IP address or the sender's IP address matches one of the IP addresses for the receiving interface, the break terminates the loop.

**INTEL Ex.1013.715** 

If the loop terminates with the variable maybe\_ia equal to 0, the entire list of configured IP addresses was searched and not one was associated with the received interface. The function jumps to out (Figure 21.19), where the mbuf is discarded and the function returns. This should only happen if an ARP request is received on an interface that has been initialized but has not been assigned an IP address.

If the for loop terminates having located a receiving interface (maybe\_ia is nonnull) but none of its IP addresses matched the sender or target IP address, myaddr is set to the final IP address assigned to the interface. Otherwise (the normal case) myaddr contains the local IP address that matched either the sender or target IP address.

Figure 21.18 shows the next part of the in\_arpinput function, which performs some validation of the packet.

		if ether.c
386	if	(!bcmp((caddr_t) ea->arp_sha, (caddr_t) ac->ac_enaddr,
387		<pre>sizeof(ea-&gt;arp_sha)))</pre>
388		goto out; /* it's from me, ignore it. */
389	if	(!bcmp((caddr_t) ea->arp_sha, (caddr_t) etherbroadcastaddr,
390		<pre>sizeof(ea-&gt;arp_sha))) {</pre>
391		log(LOG_ERR,
392		"arp: ether address is broadcast for IP address %x!\n",
393		<pre>ntohl(isaddr.s_addr));</pre>
394		goto out;
395	}	
396	if	(isaddr.s_addr == myaddr.s_addr) {
397		log(LOG_ERR,
398		"duplicate IP address %x!! sent from ethernet address: %s\n",
399		ntohl(isaddr.s_addr), ether_sprintf(ea->arp_sha));
400		itaddr = myaddr;
401		goto reply;
402	}	if ether.c

Figure 21.18 in\_arpinput function: validate received packet.

#### Validate sender's hardware address

<sup>386–388</sup> If the sender's hardware address equals the hardware address of the interface, the host received a copy of its own request, which is ignored.

<sup>389–395</sup> If the sender's hardware address is the Ethernet broadcast address, this is an error. The error is logged and the packet is discarded.

#### Check sender's IP address

396-402

385

If the sender's IP address equals myaddr, then the sender is using the same IP address as this host. This is also an error—probably a configuration error by the system administrator on either this host or the sending host. The error is logged and the function jumps to reply (Figure 21.19), after setting the target IP address to myaddr (the duplicate address). Notice that this ARP packet could have been destined for some other host on the Ethernet—it need not have been sent to this host. Nevertheless, if this form of IP address spoofing is detected, the error is logged and a reply generated.

Figure 21.19 shows the next part of in\_arpinput.

Chapter 21

40

Sec

pter 21

ne IP

ystem

func-

r (the some

if this

404

- if\_ether.c )f conla = arplookup(isaddr.s\_addr, itaddr.s\_addr == myaddr.s\_addr, 0); 403 interif (la && (rt = la->la\_rt) && (sdl = SDL(rt->rt\_gateway))) {  $4 \cap 4$ nd the 405 if (sdl->sdl\_alen && :erface bcmp((caddr\_t) ea->arp\_sha, LLADDR(sdl), sdl->sdl\_alen)) 406 log(LOG\_INFO, "arp info overwritten for %x by %s\n", 407 isaddr.s\_addr, ether\_sprintf(ea->arp\_sha)); s non-408 bcopy((caddr\_t) ea->arp\_sha, LLADDR(sdl), 409 : is set sdl->sdl\_alen = sizeof(ea->arp\_sha)); 410 /addr if (rt->rt\_expire) 411 rt->rt\_expire = time.tv\_sec + arpt\_keep; 412 rt->rt\_flags &= ~RTF\_REJECT; 413 forms la->la\_asked = 0; 414 if (la->la\_hold) { 415 (\*ac->ac\_if.if\_output) (&ac->ac\_if, la->la\_hold, 416 \_ether.c rt\_key(rt), rt); 417 la->la\_hold = 0; 418 } 419 420 } 421 reply: 422 if (op != ARPOP\_REQUEST) { 423 out: m\_freem(m); 424 425 return; 426 } - if\_ether.c Figure 21.19 in\_arpinput function: create a new ARP entry or update existing entry. ב" . Search routing table for match with sender's IP address arplookup searches the ARP cache for the sender's IP address (isaddr). The sec-403 ether.c ond argument is 1 if the target IP address equals myaddr (meaning create a new entry if an entry doesn't exist), or 0 otherwise (do not create a new entry). An entry is always created for the sender if this host is the target; otherwise the host is processing a broadcast intended for some other target, so it just looks for an existing entry for the sender. As mentioned earlier, this means that if a host receives an ARP request for itself from :e, the another host, an ARP entry is created for that other host on the assumption that, since that host is about to send us a packet, we'll probably send a reply. error.

The third argument is 0, which means do not look for a proxy ARP entry (described later). The return value is a pointer to an llinfo\_arp structure, or a null pointer if an entry is not found or created.

#### Update existing entry or fill in new entry

The code associated with the *if* statement is executed only if the following three conditions are all true:

- 1. an ARP entry was found or a new ARP entry was successfully created (la is nonnull),
- 2. the ARP entry points to a routing table entry (rt), and

Chapter 21

Sf

3. the rt\_gateway field of the routing table entry points to a sockaddr\_dl structure.

The first condition is false for every broadcast ARP request not directed to this host, from some other host whose IP address is not currently in the routing table.

# Check if sender's hardware addresses changed

405-408

If the link-level address length (sdl\_alen) is nonzero (meaning that an existing entry is being referenced and not a new entry that was just created), the link-level address is compared to the sender's hardware address. If they are different, the sender's Ethernet address has changed. This can happen if the sending host is shut down, its Ethernet interface card replaced, and it reboots before the ARP entry times out. While not common, this is a possibility that must be handled. An informational message is logged and the code continues, which will update the hardware address with its new value.

The sender's IP address in the log message should be converted to host byte order. This is a bug.

#### Record sender's hardware address

The sender's hardware address is copied into the sockaddr\_dl structure pointed to by the rt\_gateway member of the routing table entry. The link-level address length (sdl\_alen) in the sockaddr\_dl structure is also set to 6. This assignment of the length field is required if this is a newly created entry (Exercise 21.3).

# Update newly resolved ARP entry

When the sender's hardware address is resolved, the following steps occur. If the expiration time is nonzero, it is reset to 20 minutes (arpt\_keep) in the future. This test exists because the arp command can create permanent entries: entries that never time out. These entries are marked with an expiration time of 0. We'll also see in Figure 21.24 that when an ARP request is sent (i.e., for a nonpermanent ARP entry) the expiration time is set to the current time, which is nonzero.

413-414 The RTF\_REJECT flag is cleared and the la\_asked counter is set to 0. We'll see that these last two steps are used in arpresolve to avoid ARP flooding.

If ARP is holding onto an mbuf awaiting ARP resolution of that host's hardware address (the la\_hold pointer), the mbuf is passed to the interface output function. (We show this in Figure 21.3.) Since this mbuf was being held by ARP, the destination address must be on a local Ethernet so the interface output function is ether\_output. This function again calls arpresolve, but the hardware address was just filled in, allowing the mbuf to be queued on the actual device's output queue.

#### Finished with ARP reply packets

421-426 If the ARP operation is not a request, the received packet is discarded and the function returns.

The remainder of the function, shown in Figure 21.20, generates a reply to an ARP request. A reply is generated in only two instances:

の構成で

議論に対応であるまで

- 1. this host is the target of a request for its hardware address, or
- 2. this host receives a request for another host's hardware address for which this host has been configured to act as an ARP proxy server.

At this point in the function, an ARP request has been received, but since ARP requests are normally broadcast, the request could be for any system on the Ethernet.

```
- if_ether.c
        if (itaddr.s_addr == myaddr.s_addr) {
427
            /* I am the target */
428
            bcopy((caddr_t) ea->arp_sha, (caddr_t) ea->arp_tha,
429
                  sizeof(ea->arp_sha));
430
            bcopy((caddr_t) ac->ac_enaddr, (caddr_t) ea->arp_sha,
431
                  sizeof(ea->arp_sha));
432
433
        } else {
            la = arplookup(itaddr.s_addr, 0, SIN_PROXY);
434
            if (la == NULL)
435
                goto out;
436
           rt = la->la_rt;
437
           bcopy((caddr_t) ea->arp_sha, (caddr_t) ea->arp_tha,
438
                  sizeof(ea->arp_sha));
439
            sdl = SDL(rt->rt_gateway);
440
            bcopy(LLADDR(sdl), (caddr_t) ea->arp_sha, sizeof(ea->arp_sha));
441
        }
442
        bcopy((caddr_t) ea->arp_spa, (caddr_t) ea->arp_tpa, sizeof(ea->arp_spa));
443
        bcopy((caddr_t) & itaddr, (caddr_t) ea->arp_spa, sizeof(ea->arp_spa));
444
445
        ea->arp_op = htons(ARPOP_REPLY);
        ea->arp_pro = htons(ETHERTYPE_IP); /* let's be sure! */
446
        eh = (struct ether_header *) sa.sa_data;
447
        bcopy((caddr_t) ea->arp_tha, (caddr_t) eh->ether_dhost,
448
              sizeof(eh->ether_dhost));
449
        eh->ether_type = ETHERTYPE_ARP;
450
        sa.sa_family = AF_UNSPEC;
451
        sa.sa_len = sizeof(sa);
452
        (*ac->ac_if.if_output) (&ac->ac_if, m, &sa, (struct rtentry *) 0);
453
454
        return;
455 }
                                                                           - if ether.c
```

Figure 21.20 in\_arpinput function: form ARP reply and send it.

# This host is the target

427-432 If the target IP address equals myaddr, this host is the target of the request. The source hardware address is copied into the target hardware address (i.e., whoever sent it becomes the target) and the Ethernet address of the interface is copied from the arpcom structure into the source hardware address. The remainder of the ARP reply is constructed after the else clause.

# Check if this host is a proxy server for target

433-436

Even if this host is not the target, this host can be configured to be a proxy server for the specified target. arplookup is called again with the create flag set to 0 (the second

argument) and the third argument set to SIN\_PROXY. This finds an entry in the routing table only if that entry's SIN\_PROXY flag is set. If an entry is not found (the typical case where this host receives a copy of some other ARP request on the Ethernet), the code at out discards the mbuf and returns.

#### Form proxy reply

437-442

To handle a proxy ARP request, the sender's hardware address becomes the target hardware address and the Ethernet address from the ARP entry is copied into the sender hardware address field. This value from the ARP entry can be the Ethernet address of any host on the Ethernet capable of sending IP datagrams to the target IP address. Normally the host providing the proxy ARP service supplies its own Ethernet address, but that's not required. Proxy entries are created by the system administrator using the arp command, with the keyword pub, specifying the target IP address (which becomes the key of the routing table entry) and an Ethernet address to return in the ARP reply.

#### Complete construction of ARP reply packet

443-444

- The remainder of the function completes the construction of the ARP reply. The sender and target hardware addresses have been filled in. The sender and target IP addresses are now swapped. The target IP address is contained in itaddr, which might have been changed if another host was found using this host's IP address (Figure 21.18).
- The ARP operation is set to ARPOP\_REPLY and the type of protocol address is set to 445-446 ETHERTYPE\_IP. The comment "let's be sure!" is because arpintr also calls this function when the type of protocol address is ETHERTYPE\_IPTRAILERS, but the use of trailer encapsulation is no longer supported.

#### Fill in sockaddr with Ethernet header

A sockaddr structure is filled in with the 14-byte Ethernet header, as shown in Fig-447-452 ure 21.12. The target hardware address also becomes the Ethernet destination address.

453-455

The ARP reply is passed to the interface's output routine and the function returns.

#### **ARP** Timer Functions 21.9

ARP entries are normally dynamic-they are created when needed and time out automatically. It is also possible for the system administrator to create permanent entries (i.e., no timeout), and the proxy entries we discussed in the previous section are always permanent. Recall from Figure 21.1 and the #define at the end of Figure 21.10 that the rmx\_expire member of the routing metrics structure is used by ARP as a timer.

#### arptimer Function

This function, shown in Figure 21.21, is called every 5 minutes. It goes through all the ARP entries to see if any have expired.

467-

E.

arp

86

r 21 r

ting

Case

e at

rget

the

net

t IP

inet

ator

uich

the

Гhe

: IP

iich

₹ig-

t to

.nc-

e of

-ĩg-

ito-

ries ays

the

the

s. }. - if ether.c

if\_ether.c

74 static void 75 arptimer(ignored\_arg) 76 void \*ignored\_arg; 77 { s = splnet();78 int struct llinfo\_arp \*la = llinfo\_arp.la\_next; 79 timeout(arptimer, (caddr\_t) 0, arpt\_prune \* hz); 80 while (la != &llinfo\_arp) { 81 82 struct rtentry \*rt = la->la\_rt; la = la->la\_next; 83 if (rt->rt\_expire && rt->rt\_expire <= time.tv\_sec) 84 arptfree(la->la\_prev); /\* timer has expired, clear \*/ 85 86 } 87 splx(s); 88 }

Figure 21.21 arptimer function: check all ARP timers every 5 minutes.

#### Set next timeout

We'll see that the arp\_rtrequest function causes arptimer to be called the first time, and from that point arptimer causes itself to be called 5 minutes (arpt\_prune) in the future.

#### Check all ARP entries

Each entry in the linked list is processed. If the timer is nonzero (it is not a permanent entry) and if the timer has expired, arptfree releases the entry. If rt\_expire is nonzero, it contains a count of the number of seconds since the Unix Epoch when the entry expires.

# arptfree Function

80

This function, shown in Figure 21.22, is called by arptimer to delete a single entry from the linked list of llinfo\_arp entries.

#### Invalidate (don't delete) entries in use

- If the routing table reference count is greater than 0 and the rt\_gateway member points to a sockaddr\_dl structure, arptfree takes the following steps:
  - 1. the link-layer address length is set to 0,
  - 2. the la\_asked counter is reset to 0, and
  - 3. the RTF\_REJECT flag is cleared.

The function then returns. Since the reference count is nonzero, the routing table entry is not deleted. But setting sdl\_alen to 0 invalidates the entry, so the next time the entry is used, an ARP request will be generated.

696 ARP: Address Resolution Protocol

Chapter 21

```
if_ether.c
```

Sec

459 sta	tic void
460 arp	tfree(la)
461 str	uct llinfo_arp *la;
462 {	
463	struct rtentry *rt = la->la_rt;
464	struct sockaddr_dl *sdl;
465	if (rt == 0)
466	<pre>panic("arptfree");</pre>
467	if (rt->rt_refcnt > 0 && (sdl = SDL(rt->rt_gateway)) &&
468	sdl->sdl_family == AF_LINK) {
469	<pre>sdl-&gt;sdl_alen = 0;</pre>
470	la->la_asked = 0;
471	rt->rt_flags &= ~RTF_REJECT;
472	return;
473	}
474	<pre>rtrequest(RTM_DELETE, rt_key(rt), (struct sockaddr *) 0, rt_mask(rt),</pre>
475	0, (struct rtentry **) 0);
476 }	if ether.c

Figure 21.22 arpt free function: delete or invalidate an ARP entry.

#### **Delete unreferenced entries**

474-475 rtrequest deletes the routing table entry, and we'll see in Section 21.13 that it calls arp\_rtrequest. This latter function frees any mbuf chain held by the ARP entry (the la\_hold pointer) and deletes the corresponding llinfo\_arp entry.

## **21.10** arpresolve Function

We saw in Figure 4.16 that ether\_output calls arpresolve to obtain the Ethernet address for an IP address. arpresolve returns 1 if the destination Ethernet address is known, allowing ether\_output to queue the IP datagram on the interface's output queue. A return value of 0 means arpresolve does not know the Ethernet address. The datagram is "held" by arpresolve (using the la\_hold member of the llinfo\_arp structure) and an ARP request is sent. If and when an ARP reply is received, in\_arpinput completes the ARP entry and sends the held datagram.

arpresolve must also avoid *ARP flooding*, that is, it must not repeatedly send ARP requests at a high rate when an ARP reply is not received. This can happen when several datagrams are sent to the same unresolved IP address before an ARP reply is received, or when a datagram destined for an unresolved address is fragmented, since each fragment is sent to ether\_output as a separate packet. Section 11.9 of Volume 1 contains an example of ARP flooding caused by fragmentation, and discusses the associated problems. Figure 21.23 shows the first half of arpresolve.

252-261

dst is a pointer to a sockaddr\_in containing the destination IP address and desten is an array of 6 bytes that is filled in with the corresponding Ethernet address, if known.

はいないないないない

if ether.c

if\_ether.c

```
252 int
253 arpresolve(ac, rt, m, dst, desten)
254 struct arpcom *ac;
255 struct rtentry *rt;
256 struct mbuf *m;
257 struct sockaddr *dst;
258 u_char *desten;
259 {
260
        struct llinfo_arp *la;
261
        struct sockaddr_dl *sdl;
262
        if (m->m_flags & M_BCAST) { /* broadcast */
263
            bcopy((caddr_t) etherbroadcastaddr, (caddr_t) desten,
264
                  sizeof(etherbroadcastaddr));
265
            return (1);
266
        }
267
        if (m->m_flags & M_MCAST) { /* multicast */
            ETHER_MAP_IP_MULTICAST(&SIN(dst)->sin_addr, desten);
268
269
            return (1);
270
        }
271
        if (rt)
272
            la = (struct llinfo_arp *) rt->rt_llinfo;
273
        else {
274
            if (la = arplookup(SIN(dst)->sin_addr.s_addr, 1, 0))
275
                rt = la->la_rt;
276
        }
277
        if (la == 0 || rt == 0) {
278
            log(LOG_DEBUG, "arpresolve: can't allocate llinfo");
279
            m_freem(m);
            return (0);
280
281
        }
```

Figure 21.23 arpresolve function: find ARP entry if required.

#### Handle broadcast and multicast destinations

262-270 If the M\_BCAST flag of the mbuf is set, the destination is filled in with the Ethernet broadcast address and the function returns 1. If the M\_MCAST flag is set, the ETHER\_MAP\_IP\_MULTICAST macro (Figure 12.6) converts the class D address into the corresponding Ethernet address.

#### Get pointer to llinfo\_arp structure

- <sup>271–276</sup> The destination address is a unicast address. If a pointer to a routing table entry is passed by the caller, la is set to the corresponding llinfo\_arp structure. Otherwise arplookup searches the routing table for the specified IP address. The second argument is 1, telling arplookup to create the entry if it doesn't already exist; the third argument is 0, which means don't look for a proxy ARP entry.
- If either rt or la are null pointers, one of the allocations failed, since arplookup should have created an entry if one didn't exist. An error message is logged, the packet released, and the function returns 0.

Chapter 21

Figure 21.24 contains the last half of arpresolve. It checks whether the ARP entry is still valid, and, if not, sends an ARP request.

		— if ether.c
282	<pre>sdl = SDL(rt-&gt;rt_gateway);</pre>	)_
283	/*	
284	* Check the address family and length is valid, the address	
285	* is resolved; otherwise, try to resolve.	
286	*/	
287	if ((rt->rt_expire == 0    rt->rt_expire > time.tv_sec) &&	
288	sdl->sdl_family == AF_LINK && sdl->sdl_alen != 0) {	
289	<pre>bcopy(LLADDR(sdl), desten, sdl-&gt;sdl_alen);</pre>	
290	return 1;	
291	}	
292	/*	
293	* There is an arptab entry, but no ethernet address	
294	* response yet. Replace the held mbuf with this	
295	* latest one.	
296	*/	
297	if (la->la_hold)	
298	m_freem(la->la_hold);	
299	la->la_hold = m;	
300	if (rt->rt_expire) {	
301	rt->rt_flags &= ~RTF_REJECT;	
302	if (la->la_asked == 0    rt->rt_expire != time.tv_sec) {	
303	rt->rt_expire = time.tv_sec;	
304	if (la->la_asked++ < arp_maxtries)	
305	arpwhohas(ac, &(SIN(dst)->sin_addr));	
306	else {	
307	rt->rt_flags  = RTF_REJECT;	
308	rt->rt_expire += arpt_down;	
309	<pre>la-&gt;la_asked = 0;</pre>	
310	}	
311	}	
312	}	
313	return (0);	
314 }		if other c

Figure 21.24 arpresolve function: check if ARP entry valid, send ARP request if not.

#### Check ARP entry for validity

282–291 Even though an ARP entry is located, it must be checked for validity. The entry is valid if the following conditions are all true:

- 1. the entry is permanent (the expiration time is 0) or the expiration time is greater than the current time, and
- 2. the family of the socket address structure pointed to by rt\_gateway is AF\_LINK, and
- 3. the link-level address length (sdl\_alen) is nonzero.

292-299

Section 2

300-314

- if\_ether.c

Recall that arptfree invalidated an ARP entry that was still referenced by setting sdl\_alen to 0. If the entry is valid, the Ethernet address contained in the sockaddr\_dl is copied into desten and the function returns 1.

#### Hold only most recent IP datagram

292-299

At this point an ARP entry exists but it does not contain a valid Ethernet address. An ARP request must be sent. First the pointer to the mbuf chain is saved in la\_hold, after releasing any mbuf chain that was already pointed to by la\_hold. This means that if multiple IP datagrams are sent quickly to a given destination, and an ARP entry does not already exist for the destination, during the time it takes to send an ARP request and receive a reply only the *last* datagram is held, and all prior ones are discarded. An example that generates this condition is NFS. If NFS sends an 8500-byte IP datagram that is fragmented into six IP fragments, and if all six fragments are sent by ip\_output to ether\_output in the time it takes to send an ARP request and receive a reply, the first five fragments are discarded and only the final fragment is sent when the reply is received. This in turn causes an NFS timeout, and a retransmission of all six fragments.

#### Send ARP request but avoid ARP flooding

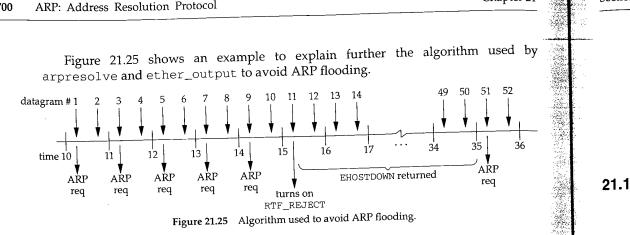
300-314

RFC 1122 requires ARP to avoid sending ARP requests to a given destination at a high rate when a reply is not received. The technique used by Net/3 to avoid ARP flooding is as follows.

- Net/3 never sends more than one ARP request in any given second to a destination.
- If a reply is not received after five ARP requests (i.e., after about 5 seconds), the RTF\_REJECT flag in the routing table is set and the expiration time is set for 20 seconds in the future. This causes ether\_output to refuse to send IP datagrams to this destination for 20 seconds, returning EHOSTDOWN or EHOSTUNREACH instead (Figure 4.15).
- After the 20-second pause in ARP requests, arpresolve will send ARP requests to that destination again.

If the expiration time is nonzero (i.e., this is not a permanent entry) the RTF\_REJECT flag is cleared, in case it had been set earlier to avoid flooding. The counter la\_asked counts the number of consecutive times an ARP request has been sent to this destination. If the counter is 0 or if the expiration time does not equal the current time (looking only at the seconds portion of the current time), an ARP request might be sent. This comparison avoids sending more than one ARP request during any second. The expiration time is then set to the current time in seconds (i.e., the microseconds portion, time.tv\_usec is ignored).

The counter is compared to the limit of 5 (arp\_maxtries) and then incremented. If the value was less than 5, arpwhohas sends the request. If the request equals 5, however, ARP has reached its limit: the RTF\_REJECT flag is set, the expiration time is set to 20 seconds in the future, and the counter la\_asked is reset to 0.



700

We show 26 seconds of time, labeled 10 through 36. We assume a process is sending an IP datagram every one-half second, causing two datagrams to be sent every second. The datagrams are numbered 1 through 52. We also assume that the destination host is down, so there are no replies to the ARP requests. The following actions take place:

- We assume la\_asked is 0 when datagram 1 is written by the process. la\_hold is set to point to datagram 1, rt\_expire is set to the current time (10), la\_asked becomes 1, and an ARP request is sent. The function returns 0.
- When datagram 2 is written by the process, datagram 1 is discarded and la\_hold is set to point to datagram 2. Since rt\_expire equals the current time (10), nothing else happens (an ARP request is not sent) and the function returns 0.
- When datagram 3 is written, datagram 2 is discarded and la\_hold is set to point to datagram 3. The current time (11) does not equal rt\_expire (10), so rt\_expire is set to 11. la\_asked is less than 5, so la\_asked becomes 2 and an ARP request is sent.
- When datagram 4 is written, datagram 3 is discarded and la\_hold is set to point to datagram 4. Since rt\_expire equals the current time (11), nothing else happens and the function returns 0.
- Similar actions occur for datagrams 5 through 10. After datagram 9 causes an ARP request to be sent, la\_asked is 5.
- When datagram 11 is written, datagram 10 is discarded and la\_hold is set to point to datagram 11. The current time (15) does not equal rt\_expire (14), so rt\_expire is set to 15. la\_asked is no longer less than 5, so the ARP flooding avoidance algorithm takes place: RTF\_REJECT flag is set, rt\_expire is set to 35 (20 seconds in the future), and la\_asked is reset to 0. The function returns 0.
- When datagram 12 is written, ether\_output notices that the RTF\_REJECT flag is set and that the current time is less than rt\_expire (35) causing EHOSTDOWN to be returned to the sender (normally ip\_output).
- The EHOSTDOWN error is returned for datagrams 13 through 50.

111-

Chapter 21

Sectio

INTEL Ex.1013.727

Section 21.11

When datagram 51 is written, even though the RTF\_REJECT flag is set ether\_output does not return the error because the current time (35) is no longer less than rt\_expire (35). arpresolve is called and the entire process starts over again: five ARP requests are sent in 5 seconds, followed by a 20-second pause. This continues until the sending process gives up or the destination host responds to an ARP request.

# 21.11 arplookup Function

arplookup calls the routing function rtalloc1 to look up an ARP entry in the Internet routing table. We've seen three calls to arplookup:

- 1. from in\_arpinput to look up and possibly create an entry corresponding to the source IP address of a received ARP packet,
- from in\_arpinput to see if a proxy ARP entry exists for the destination IP address of a received ARP request, and
- 3. from arpresolve to look up or create an entry corresponding to the destination IP address of a datagram that is about to be sent.

If arplookup succeeds, a pointer is returned to the corresponding llinfo\_arp structure; otherwise a null pointer is returned.

arplookup has three arguments. The first is the IP address to search for, the second is a flag that is true if a new entry should be created if the entry is not found, and the third is a flag that is true if a proxy ARP entry should be searched for and possibly created.

Proxy ARP entries are handled by defining a different form of the Internet socket address structure, a sockaddr\_inarp structure, shown in Figure 21.26 This structure is used only by ARP.

```
- if_ether.h
111 struct sockaddr_inarp {
                                     /* sizeof(struct sockaddr_inarp) = 16 */
112
        u_char sin_len;
        u_char sin_family;
                                     /* AF_INET */
113
114
        u_short sin_port;
                                     /* IP address */
        struct in_addr sin_addr;
115
        struct in_addr sin_srcaddr; /* not used */
116
                                     /* not used */
117
        u_short sin_tos;
                                     /* 0 or SIN_PROXY */
118
        u_short sin_other;
119 \};
                                                                            if_ether.h
```

Figure 21.26 sockaddr\_inarp structure.

111-119 The first 8 bytes are the same as a sockaddr\_in structure and the sin\_family is also set to AF\_INET. The final 8 bytes, however, are different: the sin\_srcaddr, sin\_tos, and sin\_other members. Of these three, only the final one is used, being set to SIN\_PROXY (1) if the entry is a proxy entry.

21

УУ

ın d.

is

d

1),

.d

٦t

'n

:0 :0

d

0

;e

n

0

0

g

0 ).

g

Ν

Chapter 21

Figure 21.27 shows the arplookup function.

```
if_ether.c
480 static struct llinfo_arp *
481 arplookup(addr, create, proxy)
482 u_long addr;
            create, proxy;
483 int
484 {
        struct rtentry *rt;
485
        static struct sockaddr_inarp sin =
486
        {sizeof(sin), AF_INET};
487
        sin.sin_addr.s_addr = addr;
488
        sin.sin_other = proxy ? SIN_PROXY : 0;
489
        rt = rtalloc1((struct sockaddr *) &sin, create);
490
        if (rt == 0)
491
            return (0);
492
        rt->rt_refcnt--;
        if ((rt->rt_flags & RTF_GATEWAY) || (rt->rt_flags & RTF_LLINFO) == 0 ||
493
494
            rt->rt_gateway->sa_family != AF_LINK) {
495
496
             if (create)
                 log(LOG_DEBUG, "arptnew failed on %x\n", ntohl(addr));
497
             return (0);
498
499
         }
         return ((struct llinfo_arp *) rt->rt_llinfo);
500
                                                                             if ether.c
 501 }
```

Figure 21.27 arplookup function: look up an ARP entry in the routing table.

# Initialize sockaddr\_inarp to look up

480-489 The sin\_addr member is set to the IP address that is being looked up. The sin\_other member is set to SIN\_PROXY if the proxy argument is nonzero, or 0 otherwise.

#### Look up entry in routing table

490-492 rtalloc1 looks up the IP address in the Internet routing table, creating a new entry if the create argument is nonzero. If the entry is not found, the function returns 0 (a null pointer).

#### Decrement routing table reference count

- If the entry is found, the reference count for the routing table entry is decremented. This is because ARP is not considered to "hold onto" a routing table entry like the transport layers, so the increment of rt\_refcnt that was done by the routing table lookup is undone here by ARP.
- If the RTF\_GATEWAY flag is set, or the RTF\_LLINFO flag is not set, or the address family of the socket address structure pointed to by rt\_gateway is not AF\_LINK, something is wrong and a null pointer is returned. If the entry was created this way, a log message is created.

The log message with the function name arptnew refers to the older Net/2 function that created ARP entries.

If rtalloc1 creates a new entry because the matching entry had the RTF\_CLONING flag set, the function arp\_rtrequest (which we describe in Section 21.13) is also called by rtrequest.

## 21.12 Proxy ARP

Net/3 supports proxy ARP, as we saw in the previous section. Two different types of proxy ARP entries can be added to the routing table. Both are added with the arp command, specifying the pub option. Adding a proxy ARP entry always causes a gratuitous ARP request to be issued by arp\_rtrequest (Figure 21.28) because the RTF\_ANNOUNCE flag is set when the entry is created.

The first type of proxy ARP entry allows an IP address for a host on an attached network to be entered into the ARP cache. Any Ethernet address can be assigned to the entry. These entries are added to the routing table with an explicit mask of  $0 \times fffffffff$ . The purpose of this mask is to allow the call to rtalloc1 in Figure 21.27 to match this entry, even if the SIN\_PROXY flag is set in the socket address structure of the search key. This in turn allows the call to arplookup from Figure 21.20 to match this entry when a search is made for the target address with the SIN\_PROXY flag set.

This type of entry can be used if a host H1 that doesn't implement ARP is on an attached network. The host with the proxy entry answers all ARP requests for H1's hardware address, supplying the Ethernet address that was specified when the proxy entry was created (i.e., the Ethernet address of H1). These entries are output with the notation "published" by the arp -a command.

The second type of proxy ARP entry is for a host for which a routing table entry already exists. The kernel creates another routing table entry for the destination, with this new entry containing the link-layer information (i.e., the Ethernet address). The SIN\_PROXY flag is set in the sin\_other member of the sockaddr\_inarp structure (Figure 21.26) in the new routing table entry. Recall that routing table searches compare 12 bytes of the Internet socket address structure (Figure 18.39). This use of the SIN\_PROXY flag is the only time the final 8 bytes of the structure are nonzero. When arplookup specifies the SIN\_PROXY value in the sin\_other member of the structure passed to rtalloc1, the only entries in the routing table that will match are ones that also have the SIN\_PROXY flag set.

This type of entry normally specifies the Ethernet address of the host acting as the proxy server. If the proxy entry was created for a host HD, the sequence of steps is as follows.

- 1. The proxy server receives a broadcast ARP request for HD's hardware address from some other host HS. The host HS thinks HD is on the local network.
- 2. The proxy server responds, supplying its own Ethernet address.
- 3. HS sends the datagram with a destination IP address of HD to the proxy server's Ethernet address.

ter.c

r 21

11

irns

The

ner-

ted. anskup

ress

:y, а

: cre-

e.

4. The proxy server receives the datagram for HD and forwards it, using the normal routing table entry for HD.

This type of entry was used on the router netb in the example in Section 4.6 of Volume 1. These entries are output by the arp -a command with the notation "published (proxy only)."

# 21.13 arp\_rtrequest Function

Figure 21.3 provides an overview of the relationship between the ARP functions and the routing functions. We've encountered two calls to the routing table functions from the ARP functions.

1. arplookup calls rtalloc1 to look up an ARP entry and possibly create a new entry if a match isn't found.

If a matching entry is found in the routing table and the RTF\_CLONING flag is not set (i.e., it is a matching entry for the destination host), the pointer to the matching entry is returned. But if the RTF\_CLONING bit is set, rtalloc1 calls rtrequest with a command of RTM\_RESOLVE. This is how the entries for 140.252.13.33 and 140.252.13.34 in Figure 18.2 were created—they were cloned from the entry for 140.252.13.32.

2. arptfree calls rtrequest with a command of RTM\_DELETE to delete an entry from the routing table that corresponds to an ARP entry.

Additionally, the arp command manipulates the ARP cache by sending and receiving routing messages on a routing socket. The arp command issues routing messages with commands of RTM\_ADD, RTM\_DELETE, and RTM\_GET. The first two commands cause rtrequest to be called and the third causes rtalloc1 to be called.

Finally, when an Ethernet device driver has an IP address assigned to the interface, rtinit adds a route to the network. This causes rtrequest to be called with a command of RTM\_ADD and with the flags of RTF\_UP and RTF\_CLONING. This is how the entry for 140.252.13.32 in Figure 18.2 was created.

As described in Chapter 19, each ifaddr structure can contain a pointer to a function (the ifa\_rtrequest member) that is automatically called when a routing table entry is added or deleted for that interface. We saw in Figure 6.17 that in\_ifinit sets this pointer to the function arp\_rtrequest for all Ethernet devices. Therefore, whenever the routing functions are called to add or delete a routing table entry for ARP, arp\_rtrequest is also called. The purpose of this function is to do whatever type of initialization or cleanup is required above and beyond what the generic routing table functions perform. For example, this is where a new llinfo\_arp structure is allocated and initialized whenever a new ARP entry is created. In a similar way, the llinfo\_arp structure is deleted by this function after the generic routing routines have completed processing an RTM\_DELETE command.

Figure 21.28 shows the first part of the arp\_rtrequest function.

- if\_ether.c 92 void 93 arp\_rtrequest(req, rt, sa) 94 int req; 95 struct rtentry \*rt; 96 struct sockaddr \*sa; 97 { struct sockaddr \*gate = rt->rt\_gateway; 98 99 struct llinfo\_arp \*la = (struct llinfo\_arp \*) rt->rt\_llinfo; 100 static struct sockaddr\_dl null\_sdl = {sizeof(null\_sdl), AF\_LINK}; 101 if (!arpinit\_done) { 102 103 arpinit\_done = 1; 104 timeout(arptimer, (caddr\_t) 0, hz); 105 } 106 if (rt->rt\_flags & RTF\_GATEWAY) 107 return; 108 switch (req) { case RTM\_ADD: 109 110 /\* \* XXX: If this is a manually added route to interface 111 \* such as older version of routed or gated might provide, 112 113 \* restore cloning bit. \*/ 114if ((rt->rt\_flags & RTF\_HOST) == 0 && 115 SIN(rt\_mask(rt))->sin\_addr.s\_addr != 0xffffffff) 116 117 rt->rt\_flags |= RTF\_CLONING; if (rt->rt\_flags & RTF\_CLONING) { 118 119 /\* \* Case 1: This route should come from a route to iface. 120 \*/ 121 rt\_setgate(rt, rt\_key(rt), 122 (struct sockaddr \*) &null\_sdl); 123 gate = rt->rt\_gateway; 124 SDL(gate)->sdl\_type = rt->rt\_ifp->if\_type; 125 SDL(gate)->sdl\_index = rt->rt\_ifp->if\_index; 126 rt->rt\_expire = time.tv\_sec; 127 break; 128 129 } /\* Announce a new entry if requested. \*/ 130 if (rt->rt\_flags & RTF\_ANNOUNCE) 131 arprequest((struct arpcom \*) rt->rt\_ifp, 132 &SIN(rt\_key(rt))->sin\_addr.s\_addr, 133 &SIN(rt\_key(rt))->sin\_addr.s\_addr, 134 135 (u\_char \*) LLADDR(SDL(gate))); /\* FALLTHROUGH \*/ 136 if\_ether.c

Figure 21.28 arp\_rtrequest function: RTM\_ADD command.

## Initialize ARP timeout function

The first time arp\_rtrequest is called (when the first Ethernet interface is 92-105 assigned an IP address during system initialization), the timeout function schedules the function arptimer to be called in 1 second. This starts the ARP timer code running every 5 minutes, since arptimer always calls timeout.

### Ignore indirect routes

- If the RTF\_GATEWAY flag is set, the function returns. This flag indicates an indirect 106-107 routing table entry and all ARP entries are direct routes.
- The remainder of the function is a switch with three cases: RTM\_ADD, RTM\_RESOLVE, and RTM\_DELETE. (The latter two are shown in figures that follow.) 108

## RTM\_ADD command

109

The first case for RTM\_ADD is invoked by either the arp command manually creating an ARP entry or by an Ethernet interface being assigned an IP address by rtinit (Figure 21.3).

## Backward compatibility

110-117

If the RTF\_HOST flag is cleared, this routing table entry has an associated mask (i.e., it is a network route, not a host route). If that mask is not all one bits, then the entry is really a route to an interface, so the RTF\_CLONING flag is set. As the comment indicates, this is for backward compatibility with older versions of some routing daemons. Also, the command

route add -net 224.0.0.0 -interface bsdi

that is in the file /etc/netstart creates the entry for this network shown in Figure 18.2 that has the RTF\_CLONING flag set.

## Initialize entry for network route to interface

118-126

If the RTF\_CLONING flag is set (which in\_ifinit sets for all Ethernet interfaces), this entry is probably being added by rtinit. rt\_setgate allocates space for a sockaddr\_dl structure, which is pointed to by the rt\_gateway member. This datalink socket address structure is the one associated with the routing table entry for 140.252.13.32 in Figure 21.1. The sdl\_len and sdl\_family members are initialized from the static definition of null\_sdl at the beginning of the function, and the sdl\_type (probably IFT\_ETHER) and sdl\_index members are copied from the interface's ifnet structure. This structure never contains an Ethernet address and the

sdl\_alen member remains 0. Finally, the expiration time is set to the current time, which is simply the time the 127-128 entry was created, and the break causes the function to return. For entries created at system initialization, their rmx\_expire value is the time at which the system was bootstrapped. Notice in Figure 21.1 that this routing table entry does not have an associated llinfo\_arp structure, so it is never processed by arptimer. Nevertheless this sockaddr\_dl structure is used: since it is the rt\_gateway structure for the entry that is cloned for host-specific entries on this Ethernet, it is copied by rtrequest when the newly cloned entries are created with the RTM\_RESOLVE command. Also, the netstat program prints the sdl\_index value as link#n, as we see in Figure 18.2.

130-135

Section 2

Chapter 21

136

137-144

145-146

147-158

159-16

136

#### Send gratuitous ARP request

130-135 If the RTF\_ANNOUNCE flag is set, this entry is being created by the arp command with the pub option. This option has two ramifications: (1) the SIN\_PROXY flag will be set in the sin\_other member of the sockaddr\_inarp structure, and (2) the RTF\_ANNOUNCE flag will be set. Since the RTF\_ANNOUNCE flag is set, arprequest broadcasts a gratuitous ARP request. Notice that the second and third arguments are the same, which causes the sender IP address to equal the target IP address in the ARP request.

ter 21

le is

lules

ning

irect

reat-

.nit

(i.e.,

ry is

ndions.

Fig-

ces),

or a

:ata-

for

ized

the

the

the

d at

oot-

ited

this

that

the

tat

iter-

ADD,

The code falls through to the case for the RTM\_RESOLVE command.

Figure 21.29 shows the next part of the arp\_rtrequest function, which handles the RTM\_RESOLVE command. This command is issued when rtalloc1 matches an entry with the RTF\_CLONING flag set and its second argument is nonzero (the create argument to arplookup). A new llinfo\_arp structure must be allocated and initialized.

#### Verify sockaddr\_d1 structure

137-144 The family and length of the sockaddr\_dl structure pointed to by the rt\_gateway pointer are verified. The interface type (probably IFT\_ETHER) and index are then copied into the new sockaddr\_dl structure.

#### Handle route changes

145-146 Normally the routing table entry is new and does not point to an llinfo\_arp structure. If the la pointer is nonnull, however, arp\_rtrequest was called when a route changed for an existing routing table entry. Since the llinfo\_arp structure is already allocated, the break causes the function to return.

### Initialize llinfo\_arp structure

147-158 An llinfo\_arp structure is allocated and its pointer is stored in the rt\_llinfo pointer of the routing table entry. The two statistics arp\_inuse and arp\_allocated are incremented and the llinfo\_arp structure is set to 0. This sets la\_hold to a null pointer and la\_asked to 0.

The rt pointer is stored in the llinfo\_arp structure and the RTF\_LLINFO flag is set. In Figure 18.2 we see that the three routing table entries created by ARP, 140.252.13.33, 140.252.13.34, and 140.252.13.35, all have the L flag enabled, as does the entry for 224.0.0.1. Recall that the arp program looks only for entries with this flag (Figure 19.36). Finally the new structure is added to the front of the linked list of llinfo\_arp structures by insque.

The ARP entry has been created: rtrequest creates the routing table entry (often cloning a network-specific entry for the Ethernet) and arp\_rtrequest allocates and initializes an llinfo\_arp structure. All that remains is for an ARP request to be broadcast so that an ARP reply can fill in the host's Ethernet address. In the common sequence of events, arp\_rtrequest is called because arpresolve called arplookup (the intermediate sequence of function calls can be followed in Figure 21.3). When control returns to arpresolve, it broadcasts the ARP request.

Chapter 21

3

Se

16

175

177

137	case RTM_RESOLVE:
138	if (gate->sa_family != AF_LINK
139	<pre>gate-&gt;sa_len &lt; sizeof(null sdl)) {</pre>
140	<pre>log(LOG_DEBUG, "arp_rtrequest: bad gateway value");</pre>
141	break;
142	
143	<pre>SDL(gate)-&gt;sdl_type = rt-&gt;rt_ifp-&gt;if_type;</pre>
144	<pre>SDL(gate)-&gt;sdl_index = rt-&gt;rt_ifp-&gt;if_index;</pre>
145	if (la != 0)
146	break; /* This happens on a route change */
147	/*
148	* Case 2: This route may come from cloning, or a manual route
149	* add with a LL address.
150	*/
151	<pre>R_Malloc(la, struct llinfo_arp *, sizeof(*la));</pre>
152	$rt > rt_1linfo \approx (caddr_t) la;$
153	if (la == 0) {
154	<pre>log(LOG_DEBUG, "arp_rtrequest: malloc failed\n");</pre>
155	break;
156	}
157	arp_inuse++, arp_allocated++;
158	<pre>Bzero(la, sizeof(*la));</pre>
159	la->la_rt = rt;
160	rt->rt_flags  = RTF_LLINFO;
161	<pre>insque(la, &amp;llinfo_arp);</pre>
162	if (SIN(rt_key(rt))->sin_addr.s addr ==
163	(IA_SIN(rt->rt_ifa))->sin_addr.s_addr) {
164	/*
165	* This test used to be
166	<pre>* if (loif.if_flags &amp; IFF UP)</pre>
167	* It allowed local traffic to be forced
168	* through the hardware by configuring the loopback down.
169	* However, it causes problems during network configuration
170	* for boards that can't receive packets they send.
171	* It is now necessary to clear "useloopback" and remove
172	* the route to force traffic out to the hardware.
173	*/
174	rt->rt_expire = 0;
175	<pre>Bcopy(((struct arpcom *) rt-&gt;rt_ifp)-&gt;ac_enaddr,</pre>
176	LLADDR(SDL(gate)), SDL(gate)->sdl_alen = 6);
177	if (useloopback)
178	<pre>rt-&gt;rt_ifp = &amp;loif</pre>
179	}
180	break:

 $Figure \ 21.29 \quad \texttt{arp\_rtrequest function: RTM\_RESOLVE command}.$ 

#### Handle local host specially

This portion of code is a special test that is new with 4.4BSD (although the comment 162-173 is left over from earlier releases). It creates the rightmost routing table entry in Figure 21.1 with a key consisting of the local host's IP address (140.252.13.35). The if test checks whether the routing table key equals the IP address of the interface. If so, the entry that was just created (probably as a clone of the interface entry) refers to the local host.

## Make entry permanent and set Ethernet address

The expiration time is set to 0, making the entry permanent—it will never time out. 174-176 The Ethernet address is copied from the arpcom structure of the interface into the sockaddr\_dl structure pointed to by the rt\_gateway member.

## Set interface pointer to loopback interface

If the global useloopback is nonzero (it defaults to 1), the interface pointer in the 177--178 routing table entry is changed to point to the loopback interface. This means that any datagrams sent to the host's own IP address are sent to the loopback interface instead. Prior to 4.4BSD, the route from the host's own IP address to the loopback interface was established using a command of the form

route add 140.252.13.35 127.0.0.1

in the /etc/netstart file. Although this still works with 4.4BSD, it is unnecessary because the code we just looked at creates an equivalent route automatically, the first time an IP datagram is sent to the host's own IP address. Also realize that this piece of code is executed only once per interface. Once the routing table entry and the permanent ARP entry are created, they don't expire, so another RTM\_RESOLVE for this IP address won't occur.

The final part of arp\_rtrequest, shown in Figure 21.30, handles the RTM\_DELETE request. From Figure 21.3 we see that this command can be generated from the arp command, to delete an entry manually, and from the arptfree function, when an ARP entry times out.

		if_ether.c
181	case RTM_DELETE:	
182	if (la == 0)	
183	break;	
184	arp_inuse;	
185	<pre>remque(la);</pre>	
186	<pre>rt-&gt;rt_llinfo = 0;</pre>	
187	rt->rt_flags &= ~RTF_LLINFO;	
188	if (la->la_hold)	
189	<pre>m_freem(la-&gt;la_hold);</pre>	
190	<pre>Free((caddr_t) la);</pre>	
191	}	
192 }		if_ether.c

Figure 21.30 arp\_rtrequest function: RTM\_DELETE command.

f\_ether.c

pter 21

:e



- if\_ether.c

£.

14

#### Verify 1a pointer

182-183 The la pointer should always be nonnull (that is, the routing table entry should always point to an llinfo\_arp structure); otherwise the break causes the function to return.

### Delete llinfo\_arp structure

184--190

The arp\_inuse statistic is decremented and the llinfo\_arp structure is removed from the doubly linked list by remque. The rt\_llinfo pointer is set to 0 and the RTF\_LLINFO flag is cleared. If an mbuf is held by the ARP entry (i.e., an ARP request is outstanding), that mbuf is released. Finally the llinfo\_arp structure is released.

Notice that the switch statement does not provide a default case and does not provide a case for the RTM\_GET command. This is because the RTM\_GET command issued by the arp program is handled entirely by the route\_output function, and rtrequest is not called. Also, the call to rtalloc1 that we show in Figure 21.3, which is caused by an RTM\_GET command, specifies a second argument of 0; therefore rtalloc1 does not call rtrequest in this case.

## 21.14 ARP and Multicasting

If an IP datagram is destined for a multicast group, ip\_output checks whether the process has assigned a specific interface to the socket (Figure 12.40), and if so, the datagram is sent out that interface. Otherwise, ip\_output selects the outgoing interface using the normal IP routing table (Figure 8.24). Therefore, on a system with more than one multicast-capable interface, the IP routing table specifies the default interface for each multicast group.

We saw in Figure 18.2 that an entry was created in our routing table for the 224.0.0.0 network and since that entry has its "clone" flag set, all multicast groups starting with 224 had the associated interface (1e0) as its default. Additional routing table entries can be created for the other multicast groups (the ones beginning with 225–239), or specific entries can be created for particular multicast groups to assign an explicit default. For example, a routing table entry could be created for 224.0.1.1 (the network time protocol) with an interface that differs from the interface for 224.0.0.0. If an entry for a multicast group does not exist in the routing table, and the process doesn't specify an interface with the IP\_MULTICAST\_IF socket option, the default interface for the group becomes the interface associated with the "default" route in the table. In Figure 18.2 the entry for 224.0.0.0 isn't really needed, since both it and the default route use the interface 1e0.

Once the interface is selected, if the interface is an Ethernet, arpresolve is called to convert the multicast group address into its corresponding Ethernet address. In Figure 21.23 this was done by invoking the macro ETHER\_MAP\_IP\_MULTICAST. Since this simple macro logically ORs the low-order 23 bits of the multicast group with a constant (Figure 12.6), an ARP request-reply is not required and the mapping does not need to go into the ARP cache. The macro is just invoked each time the conversion is required.

Multicast group addresses appear in the Net/3 ARP cache if the multicast group is cloned from another entry, as we saw in Figure 21.5. This is because these entries have

Chapter 21

pter 21

should tion to

noved nd the uest is

es not mand n, and e 21.3, erefore

ver the

e dataterface :e than ace for '4.0.0.0 g with entries or spelefault. k time y for a cify an group 8.2 the ≥ inter-

called In Figice this instant ieed to iired. roup is is have the RTF\_LLINFO flag set. These are not true ARP entries because they do not require an ARP request-reply, and they do not have an associated link-layer address, since the mapping is done when needed by the ETHER\_MAP\_IP\_MULTICAST macro.

The timeout of the ARP entries for these multicast group addresses is different from normal ARP entries. When a routing table entry is created for a multicast group, such as the entry for 224.0.0.1 in Figure 18.2, rtrequest copies the rt\_metrics structure from the entry being cloned (Figure 19.9). We mentioned with Figure 21.28 that the network entry has an rmx\_expire value of the time the RTM\_ADD command was executed, normally the time the system was initialized. The new entry for 224.0.0.1 has this same expiration time.

This means the ARP entry for a multicast group such as 224.0.0.1 expires the next time arptimer executes, because its expiration time is always in the past. The entry is created again the next time it is looked up in the routing table.

## 21.15 Summary

ARP provides the dynamic mapping between IP addresses and hardware addresses. This chapter has examined an implementation of ARP that maps IP addresses to Ethernet addresses.

The Net/3 implementation is a major change from previous BSD releases. The ARP information is now stored in various structures: the routing table, a data-link socket address structure, and an llinfo\_arp structure. Figure 21.1 shows the relationships between all the structures.

Sending an ARP request is simple: the appropriate fields are filled in and the request is sent as a broadcast. Processing a received request is more complicated because each host receives *all* broadcast ARP requests. Besides responding to requests for one of the host's IP addresses, in\_arpinput also checks that some other host isn't using the host's IP address. Since all ARP requests contain the sender's IP and hardware addresses, any host on the Ethernet can use this information to update an existing ARP entry for the sender.

ARP flooding can be a problem on a LAN and Net/3 is the first BSD release to handle this. A maximum of one ARP request per second is sent to any given destination, and after five consecutive requests without a reply, a 20-second pause occurs before another ARP request is sent to that destination.

## Exercises

- 21.1 What assumption is made in the assignment of the local variable ac in Figure 21.17?
- **21.2** If we ping the broadcast address of the local Ethernet and then execute arp -a, we see that this causes the ARP cache to be filled with entries for almost every other host on the local Ethernet. Why?
- 21.3 Follow through the code and explain why the assignment of 6 to sdl\_alen is required in Figure 21.19.

Chapter 21

- 21.4 With the separate ARP table in Net/2, independent of the routing table, each time arpresolve was called, a search was made of the ARP table. Compare this to the Net/3 approach. Which is more efficient?
- **21.5** The ARP code in Net/2 explicitly set a timeout of 3 minutes for an incomplete entry in the ARP cache, that is, for an entry that is awaiting an ARP reply. We've never explicitly said how Net/3 handles this timeout. When does Net/3 time out an incomplete ARP entry?
- 21.6 What changes in the avoidance of ARP flooding when a Net/3 system is acting as a router and the packets that cause the flooding are from some other host?
- 21.7 What are the values of the four rmx\_expire variables shown in Figure 21.1? Where in the code are the values set?
- 21.8 What change would be required to the code in this chapter to cause an ARP entry to be created for every host that broadcasts an ARP request?
- **21.9** To verify the example in Figure 21.25 the authors ran the sock program from Appendix C of Volume 1, writing a UDP datagram every 500 ms to a nonexistent host on the local Ethernet. (The -p option of the program was modified to allow millisecond waits.) But only 10 UDP datagrams were sent without an error, instead of the 11 shown in Figure 21.25, before the first EHOSTDOWN error was returned. Why?
- **21.10** Modify ARP to hold onto *all* packets for a destination, awaiting an ARP reply, instead of just the most recent one. What are the implications of this change? Should there be a limit, as there is for each interface's output queue? Are any changes required to the data structures?

22.1

r 21

22

# **Protocol Control Blocks**

## 22.1 Introduction

Protocol control blocks (PCBs) are used at the protocol layer to hold the various pieces of information required for each UDP or TCP socket. The Internet protocols maintain *Internet protocol control blocks* and *TCP control blocks*. Since UDP is connectionless, everything it needs for an end point is found in the Internet PCB; there are no UDP control blocks.

The Internet PCB contains the information common to all UDP and TCP end points: foreign and local IP addresses, foreign and local port numbers, IP header prototype, IP options to use for this end point, and a pointer to the routing table entry for the destination of this end point. The TCP control block contains all of the state information that TCP maintains for each connection: sequence numbers in both directions, window sizes, retransmission timers, and the like.

In this chapter we describe the Internet PCBs used in Net/3, saving TCP's control blocks until we describe TCP in detail. We examine the numerous functions that operate on Internet PCBs, since we'll encounter them when we describe UDP and TCP. Most of the functions begin with the six characters in\_pcb.

Figure 22.1 summarizes the protocol control blocks that we describe and their relationship to the file and socket structures. There are numerous points to consider in this figure.

• When a socket is created by either socket or accept, the socket layer creates a file structure and a socket structure. The file type is DTYPE\_SOCKET and the socket type is SOCK\_DGRAM for UDP end points or SOCK\_STREAM for TCP end points.

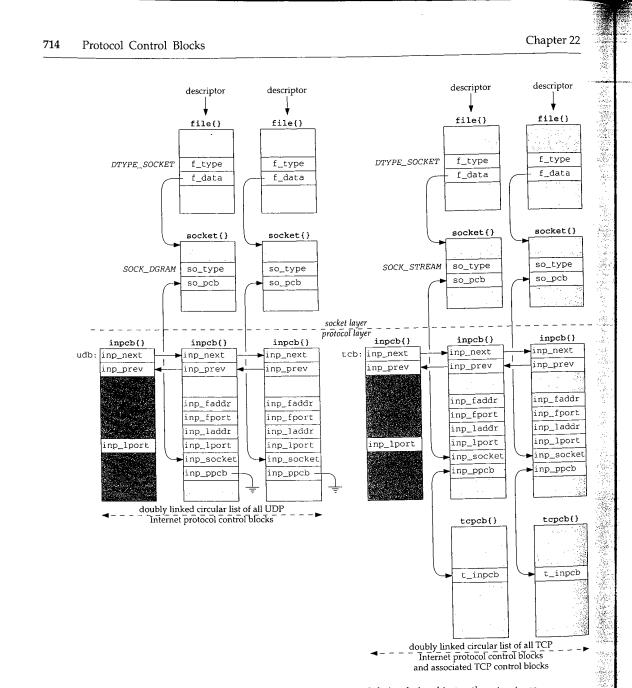


Figure 22.1 Internet protocol control blocks and their relationship to other structures.

- The protocol layer is then called. UDP creates an Internet PCB (an inpcb structure) and links it to the socket structure: the so\_pcb member points to the inpcb structure and the inp\_socket member points to the socket structure.
- TCP does the same and also creates its own control block (a tcpcb structure) and links it to the inpcb using the inp\_ppcb and t\_inpcb pointers. In the

INTEL Ex.1013.741

two UDP inpcbs the inp\_ppcb member is a null pointer, since UDP does not maintain its own control block.

- The four other members of the inpcb structure that we show, inp\_faddr through inp\_lport, form the socket pair for this end point: the foreign IP address and port number along with the local IP address and port number.
- Both UDP and TCP maintain a doubly linked list of all their Internet PCBs, using the inp\_next and inp\_prev pointers. They allocate a global inpcb structure as the head of their list (named udb and tcb) and only use three members in the structure: the next and previous pointers, and the local port number. This latter member contains the next ephemeral port number to use for this protocol.

The Internet PCB is a transport layer data structure. It is used by TCP, UDP, and raw IP, but not by IP, ICMP, or IGMP.

We haven't described raw IP yet, but it too uses Internet PCBs. Unlike TCP and UDP, raw IP does not use the port number members in the PCB, and raw IP uses only two of the functions that we describe in this chapter: in\_pcballoc to allocate a PCB, and in\_pcbdetach to release a PCB. We return to raw IP in Chapter 32.

## 22.2 Code Introduction

All the PCB functions are in a single C file and a single header contains the definitions, as shown in Figure 22.2.

File	Description		
netinet/in_pcb.h	inpcb structure definition		
netinet/in_pcb.c	PCB functions		

Figure 22.2 Files discussed in this chapter.

### **Global Variables**

One global variable is introduced in this chapter, which is shown in Figure 22.3.

Variable	Datatype	Description	
zeroin_addr	struct in_addr	32-bit IP address of all zero bits	

Figure 22.3 Global variable introduced in this chapter.

#### Statistics

Internet PCBs and TCP PCBs are both allocated by the kernel's malloc function with a type of M\_PCB. This is just one of the approximately 60 different types of memory

ıche

2.

re)

he

allocated by the kernel. Mbufs, for example, are allocated with a type of M\_BUF, and socket structures are allocated with a type of M\_SOCKET.

Since the kernel can keep counters of the different types of memory buffers that are allocated, various statistics on the number of PCBs can be maintained. The command vmstat -m shows the kernel's memory allocation statistics and the netstat -m command shows the mbuf allocation statistics.

## 22.3 inpcb Structure

Figure 22.4 shows the definition of the inpub structure. It is not a big structure, and occupies only 84 bytes.

42 str	uct inpcb {
43	<pre>struct inpcb *inp_next, *inp_prev; /* doubly linked list */</pre>
4.4	struct inpcb *inp_head; /* pointer back to chain of inpcb's for
45	this protocol */
46	struct in_addr inp_faddr; /* foreign IP address */
47	<pre>u_short inp_fport; /* foreign port# */</pre>
48	struct in_addr inp_laddr; /* local IP address */
49	u_short inp_lport; /* local port# */
50	struct socket *inp_socket; /* back pointer to socket */
51	<pre>caddr_t inp_ppcb; /* pointer to per-protocol PCB */</pre>
52	<pre>struct route inp_route;  /* placeholder for routing entry */</pre>
53	<pre>int inp_flags;</pre>
54	<pre>struct ip inp_ip; /* header prototype; should have more */</pre>
55	struct mbuf *inp_options; /* IP options */
56	<pre>struct ip_moptions *inp_moptions; /* IP multicast options */</pre>
57 };	in and h

Figure 22.4 inpcb structure.

43-45 inp\_next and inp\_prev form the doubly linked list of all PCBs for UDP and TCP. Additionally, each PCB has a pointer to the head of the protocol's linked list (inp\_head). For PCBs on the UDP list, inp\_head always points to udb (Figure 22.1); for PCBs on the TCP list, this pointer always points to tcb.

The next four members, inp\_faddr, inp\_fport, inp\_laddr, and inp\_lport, contain the socket pair for this IP end point: the foreign IP address and port number and the local IP address and port number. These four values are maintained in the PCB in network byte order, not host byte order.

The Internet PCB is used by both transport layers, TCP and UDP. While it makes sense to store the local and foreign IP addresses in this structure, the port numbers really don't belong here. The definition of a port number and its size are specified by each transport layer and could differ between different transport layers. This problem was identified in [Partridge 1987], where 8-bit port numbers were used in version 1 of RDP, which required reimplementing several standard kernel routines to use 8-bit port numbers. Version 2 of RDP [Partridge and Hinden 1990] uses 16-bit port numbers. The port numbers really belong in a transport-specific control block, such as TCP's tcpcb. A new UDP-specific PCB would then be required. While doable, this would complicate some of the routines we'll examine shortly.

50-51

52

53

54

inp\_socket is a pointer to the socket structure for this PCB and inp\_ppcb is a pointer to an optional transport-specific control block for this PCB. We saw in Figure 22.1 that the inp\_ppcb pointer is used with TCP to point to the corresponding tcpcb, but is not used by UDP. The link between the socket and inpcb is two way because sometimes the kernel starts at the socket layer and needs to find the corresponding Internet PCB (e.g., user output), and sometimes the kernel starts at the PCB and needs to locate the corresponding socket structure (e.g., processing a received IP datagram).

If IP has a route to the foreign address, it is stored in the inp\_route entry. We'll see that when an ICMP redirect message is received, all Internet PCBs are scanned and all those with a foreign IP address that matches the redirected IP address have their inp\_route entry marked as invalid. This forces IP to find a new route to the foreign address the next time the PCB is used for output.

Various flags are stored in the inp\_flags member. Figure 22.5 lists the individual flags.

inp_flags	Description
INP_HDRINCL INP_RECVOPTS INP_RECVRETOPTS INP_RECVDSTADDR	process supplies entire IP header (raw socket only) receive incoming IP options as control information (UDP only, not implemented) receive IP options for reply as control information (UDP only, not implemented) receive IP destination address as control information (UDP only)
INP_CONTROLOPTS	INP_RECVOPTS   INP_RECVRETOPTS   INP_RECVDSTADDR

Figure 22.5 inp\_flags values.

A copy of an IP header is maintained in the PCB but only two members are used, the TOS and TTL. The TOS is initialized to 0 (normal service) and the TTL is initialized by the transport layer. We'll see that TCP and UDP both default the TTL to 64. A process can change these defaults using the IP\_TOS or IP\_TTL socket options, and the new value is recorded in the inpcb.inp\_ip structure. This structure is then used by TCP and UDP as the prototype IP header when sending IP datagrams.

A process can set the IP options for outgoing datagrams with the IP\_OPTIONS socket option. A copy of the caller's options are stored in an mbuf by the function ip\_pcbopts and a pointer to that mbuf is stored in the inp\_options member. Each time TCP or UDP calls the ip\_output function, a pointer to these IP options is passed for IP to insert into the outgoing IP datagram. Similarly, a pointer to a copy of the user's IP multicast options is maintained in the inp\_moptions member.

## 22.4 in\_pcballoc and in\_pcbdetach Functions

An Internet PCB is allocated by TCP, UDP, and raw IP when a socket is created. A PRU\_ATTACH request is issued by the socket system call. In the case of UDP, we'll see in Figure 23.33 that the resulting call is

ter 22

it are nand com-

and

and

\_pcb.h

\_pcb.h

TCP. [ list 22.1); ort, c and

CB in

) store ; here. ld difwhere everal linden ontrol oable, struct socket \*so; int error;

error = in\_pcballoc(so, &udb);

Figure 22.6 shows the in\_pcballoc function.

36 int		
37 in_	pcballoc(so, head)	
38 str	uct socket *so;	
39 str	uct inpcb *head;	
40 {		
41	struct inpcb *inp;	
42	MALLOC(inp, struct inpcb *, sizeof(*inp), M_PCB, M_WAITOK);	
43	if (inp == NULL)	
44	return (ENOBUFS);	
45	<pre>bzero((caddr_t) inp, sizeof(*inp));</pre>	
46	inp->inp_head = head;	
47	inp->inp_socket = so;	
48	<pre>insque(inp, head);</pre>	
49	so->so_pcb = (caddr_t) inp;	
50	return (0);	
51 }		— in_pcb.c

Figure 22.6 in\_pcballoc function: allocate an Internet PCB.

#### Allocate PCB and initialize to zero

<sup>36-45</sup> in\_pcballoc calls the kernel's memory allocator using the macro MALLOC. Since these PCBs are always allocated as the result of a system call, it is OK to wait for one.

Net/2 and earlier Berkeley releases stored both Internet PCBs and TCP PCBs in mbufs. Their sizes were 80 and 108 bytes, respectively. With the Net/3 release, the sizes went to 84 and 140 bytes, so TCP control blocks no longer fit into an mbuf. Net/3 uses the kernel's memory allocator instead of mbufs for both types of control blocks.

Careful readers may note that the example in Figure 2.6 shows 17 mbufs allocated for PCBs, yet we just said that Net/3 no longer uses mbufs for Internet PCBs or TCP PCBs. Net/3 does, however, use mbufs for Unix domain PCBs, and that is what this counter refers to. The mbuf statistics output by netstat are for all mbufs in the kernel across all protocol suites, not just the Internet protocols.

bzero sets the PCB to 0. This is important because the IP addresses and port numbers in the PCB must be initialized to 0.

#### Link structures together

46-49

The inp\_head member points to the head of the protocol's PCB list (either udb or tcb), the inp\_socket member points to the socket structure, the new PCB is added to the protocol's doubly linked list (insque), and the socket structure points to the PCB. The insque function puts the new PCB at the head of the protocol's list.

Chapter 22

in\_pcb.c

t

5

r

f

е

An Internet PCB is deallocated when a PRU\_DETACH request is issued. This happens when the socket is closed. The function in\_pcbdetach, shown in Figure 22.7, is eventually called.

- in\_pcb.c

```
252 int
253 in_pcbdetach(inp)
254 struct inpcb *inp;
255 {
        struct socket *so = inp->inp_socket;
256
        so->so_pcb = 0;
257
258
        sofree(so);
        if (inp->inp_options)
259
             (void) m_free(inp->inp_options);
260
        if (inp->inp_route.ro_rt)
261
            rtfree(inp->inp_route.ro_rt);
262
        ip_freemoptions(inp->inp_moptions);
263
        remque(inp);
264
265
        FREE(inp, M_PCB);
266 }
```

– in\_pcb.c

Figure 22.7 in\_pcbdetach function: deallocate an Internet PCB.

The PCB pointer in the socket structure is set to 0 and that structure is released by sofree. If an mbuf with IP options was allocated for this PCB, it is released by m\_free. If a route is held by this PCB, it is released by rtfree. Any multicast options are also released by ip\_freemoptions.

The PCB is removed from the protocol's doubly linked list by remque and the memory used by the PCB is returned to the kernel.

## 22.5 Binding, Connecting, and Demultiplexing

Before examining the kernel functions that bind sockets, connect sockets, and demultiplex incoming datagrams, we describe the rules imposed by the kernel on these actions.

## Binding of Local IP Address and Port Number

Figure 22.8 shows the six different combinations of a local IP address and local port number that a process can specify in a call to bind.

The first three lines are typical for servers—they bind a specific port, termed the server's *well-known port*, whose value is known by the client. The last three lines are typical for clients—they don't care what the local port, termed an *ephemeral port*, is, as long as it is unique on the client host.

Most servers and most clients specify the wildcard IP address in the call to bind. This is indicated in Figure 22.8 by the notation \* on lines 3 and 6.

De

Chapter 22

Local IP address	Local port	Description		
unicast or broadcast	nonzero	one local interface, specific port		
multicast	nonzero	one local multicast group, specific port		
*	nonzero	any local interface or multicast group, specific port		
unicast or broadcast	0	one local interface, kernel chooses port		
multicast	0	one multicast group, kernel chooses port		
*	0	any local interface, kernel chooses port		

Figure 22.8 Combination of local IP address and local port number for bind.

If a server binds a specific IP address to a socket (i.e., not the wildcard address), then only IP datagrams arriving with that specific IP address as the destination IP address—be it unicast, broadcast, or multicast—are delivered to the process. Naturally, when the process binds a specific unicast or broadcast IP address to a socket, the kernel verifies that the IP address corresponds to a local interface.

It is rare, though possible, for a client to bind a specific IP address (lines 4 and 5 in Figure 22.8). Normally a client binds the wildcard IP address (the final line in Figure 22.8), which lets the kernel choose the outgoing interface based on the route chosen to reach the server.

What we don't show in Figure 22.8 is what happens if the client tries to bind a local port that is already in use with another socket. By default a process cannot bind a port number if that port is already in use. The error EADDRINUSE (address already in use) is returned if this occurs. The definition of *in use* is simply whether a PCB exists with that port as its local port. This notion of "in use" is relative to a given protocol: TCP or UDP, since TCP port numbers are independent of UDP port numbers.

Net/3 allows a process to change this default behavior by specifying one of following two socket options:

SO\_REUSEADDR Allows the process to bind a port number that is already in use, but the IP address being bound (including the wildcard) must not already be bound to that same port.

For example, if an attached interface has the IP address 140.252.1.29 then one socket can be bound to 140.252.1.29, port 5555; another socket can be bound to 127.0.0.1, port 5555; and another socket can be bound to the wildcard IP address, port 5555. The call to bind for the second and third cases must be preceded by a call to setsockopt, setting the SO\_REUSEADDR option.

SO\_REUSEPORT Allows a process to reuse both the IP address and port number, but each binding of the IP address and port number, including the first, must specify this socket option. With SO\_REUSEADDR, the first binding of the port number need not specify the socket option.

For example, if an attached interface has the IP address 140.252.1.29 and a socket is bound to 140.252.1.29, port 6666 specifying the

Со

Sec

SO\_REUSEPORT socket option, then another socket can also specify this same socket option and bind 140.252.1.29, port 6666.

Later in this section we describe what happens in this final example when an IP datagram arrives with a destination address of 140.252.1.29 and a destination port of 6666, since two sockets are bound to that end point.

> The SO\_REUSEPORT option is new with Net/3 and was introduced with the support for multicasting in 4.4BSD. Before this release it was never possible for two sockets to be bound to the same IP address and same port number.

> Unfortunately the SO\_REUSEPORT option was not part of the original Stanford multicast sources and is therefore not widely supported. Other systems that support multicasting, such as Solaris 2.x, let a process specify SO\_REUSEADDR to specify that it is OK to bind multiple sockets to the same IP address and same port number.

## **Connecting a UDP Socket**

We normally associate the connect system call with TCP clients, but it is also possible for a UDP client or a UDP server to call connect and specify the foreign IP address and foreign port number for the socket. This restricts the socket to exchanging UDP datagrams with that one particular peer.

There is a side effect when a UDP socket is connected: the local IP address, if not already specified by a call to bind, is automatically set by connect. It is set to the local interface address chosen by IP routing to reach the specified peer.

Figure 22.9 shows the three different states of a UDP socket along with the pseudocode of the function calls to end up in that state.

Local socket	Foreign socket	Description		
localIP.lport	foreignIP.fport	<pre>restricted to one peer: socket(), bind(*, lport), connect(foreignIP, fport) socket(), bind(localIP, lport), connect(foreignIP, fport)</pre>		
localIP.lport	*.*	restricted to datagrams arriving on one local interface: localIP socket(), bind(localIP, lport)		
*.lport	*.*	<pre>receives all datagrams sent to lport: socket(), bind(*, lport)</pre>		

Figure 22.9 Specification of local and foreign IP addresses and port numbers for UDP sockets.

The first of the three states is called a *connected UDP socket* and the next two states are called *unconnected UDP sockets*. The difference between the two unconnected sockets is that the first has a fully specified local address and the second has a wildcarded local IP address.

## Demultiplexing of Received IP Datagrams by TCP

Figure 22.10 shows the state of three Telnet server sockets on the host sun. The first two sockets are in the LISTEN state, waiting for incoming connection requests, and the third

ard address), estination IP ss. Naturally, et, the kernel

: port

les 4 and 5 in l line in Figroute chosen

o bind a local of bind a port ady in use) is ists with that TCP or UDP,

me of follow-

ly in use, but st not already

3 140.252.1.29 i555; another socket can be bind for the setsockopt,

number, but ling the first, the first bindı.

3 140.252.1.29 Decifying the

is connected to a client at port 1500 on the host with an IP address of 140.252.1.11. The first listening socket will handle connection requests that arrive on the 140.252.1.29 interface and the second listening socket will handle all other interfaces (since its local IP address is the wildcard).

1	Local address	Local port	Foreign address	Foreign port	TCP state
	140 050 1 20		*	*	LISTEN
	140.252.1.29	23	*	*	LISTEN
	*		140.252.1.11	1500	ESTABLISHED
	140.252.1.29	23	140.202.1.11	1000	L

Figure 22.10	Three TCP sockets with a local	port of	23.
--------------	--------------------------------	---------	-----

We show both of the listening sockets with unspecified foreign IP addresses and port numbers because the sockets API doesn't allow a TCP server to restrict either of these values. A TCP server must accept the client's connection and is then told of the client's IP address and port number after the connection establishment is complete (i.e., when TCP's three-way handshake is complete). Only then can the server close the connection if it doesn't like the client's IP address and port number. This isn't a required TCP feature, it is just the way the sockets API has always worked.

When TCP receives a segment with a destination port of 23 it searches through its list of Internet PCBs looking for a match by calling in\_pcblookup. When we examine this function shortly we'll see that it has a preference for the smallest number of *wildcard matches*. To determine the number of wildcard matches we consider only the local and foreign IP addresses. We do not consider the foreign port number. The local port number must match, or we don't even consider the PCB. The number of wildcard matches can be 0, 1 (local IP address or foreign IP address), or 2 (both local and foreign IP addresses).

For example, assume the incoming segment is from 140.252.1.11, port 1500, destined for 140.252.1.29, port 23. Figure 22.11 shows the number of wildcard matches for the three sockets from Figure 22.10.

Local address	Local port	Foreign address	Foreign port	TCP state	#wildcard matches
140.252.1.29	23	*	*	LISTEN	1
*	23	*	*	LISTEN	2
140.252.1.29	23	140.252.1.11	1500	ESTABLISHED	0

Figure 22.11 Incoming segment from (140.252.1.11, 1500) to (140.252.1.29, 23).

The first socket matches these four values, but with one wildcard match (the foreign IP address). The second socket also matches the incoming segment, but with two wildcard matches (the local and foreign IP addresses). The third socket is a complete match with no wildcards. Net/3 uses the third socket, the one with the smallest number of wildcard matches.

Continuing this example, assume the incoming segment is from 140.252.1.11, port 1501, destined for 140.252.1.29, port 23. Figure 22.12 shows the number of wildcard matches.

Local address	Local port	Foreign address	Foreign port	TCP state	#wildcard matches
140.252.1.29	23	*	*	LISTEN	1
*	23	*	*	LISTEN	2
140.252.1.29	23	140.252.1.11	1500	ESTABLISHED	

Figure 22.12 Incoming segment from {140.252.1.11, 1501} to {140.252.1.29, 23}.

The first socket matches with one wildcard match; the second socket matches with two wildcard matches; and the third socket doesn't match at all, since the foreign port numbers are unequal. (The foreign port numbers are compared only if the foreign IP address in the PCB is not a wildcard.) The first socket is chosen.

In these two examples we never said what type of TCP segment arrived: we assume that the segment in Figure 22.11 contains data or an acknowledgment for an established connection since it is delivered to an established socket. We also assume that the segment in Figure 22.12 is an incoming connection request (a SYN) since it is delivered to a listening socket. But the demultiplexing code in in\_pcblookup doesn't care. If the TCP segment is the wrong type for the socket that it is delivered to, we'll see later how TCP handles this. For now the important fact is that the demultiplexing code only compares the source and destination socket pair from the IP datagram against the values in the PCB.

#### Demultiplexing of Received IP Datagrams by UDP

The delivery of UDP datagrams is more complicated than the TCP example we just examined, since UDP datagrams can be sent to a broadcast or multicast address. Since Net/3 (and most systems with multicast support) allow multiple sockets to have identical local IP addresses and ports, how are multiple recipients handled? The Net/3 rules are:

- 1. An incoming UDP datagram destined for either a broadcast IP address or a multicast IP address is delivered to *all* matching sockets. There is no concept of a "best" match here (i.e., the one with the smallest number of wildcard matches).
- 2. An incoming UDP datagram destined for a unicast IP address is delivered only to *one* matching socket, the one with the smallest number of wildcard matches. If there are multiple sockets with the same "smallest" number of wildcard matches, which socket receives the incoming datagram is implementation-dependent.

Figure 22.13 shows four UDP sockets that we'll use for some examples. Having four UDP sockets with the same local port number requires using either SO\_REUSEADDR or SO\_REUSEPORT. The first two sockets have been connected to a foreign IP address and port number, and the last two are unconnected.

'he .29 cal

d

Chapter 22

Local address	Local port	Foreign address	Foreign port	Comment
140.252.1.29 140.252.13.63 140.252.13.63 *	577 577 577 577 577	140.252.1.11 140.252.13.35 * *	1500	connected, local IP = unicast connected, local IP = broadcast unconnected, local IP = broadcast unconnected, local IP = wildcard

Figure 22.13 Four UDP sockets with a local port of 577.

Consider an incoming UDP datagram destined for 140.252.13.63 (the broadcast address on the 140.252.13 subnet), port 577, from 140.252.13.34, port 1500. Figure 22.14 shows that it is delivered to the third and fourth sockets.

Local address	Local port	Foreign address	Foreign port	Delivered?
140.252.1.29 140.252.13.63 140.252.13.63 *	577 577 577 577 577 577	140.252.1.11 140.252.13.35 * *	1500 1500 *	no, local and foreign IP mismatch no, foreign IP mismatch yes yes

Figure 22.14 Received datagram from {140.252.13.34, 1500} to (140.252.13.63, 577).

The broadcast datagram is not delivered to the first socket because the local IP address doesn't match the destination IP address and the foreign IP address doesn't match the source IP address. It isn't delivered to the second socket because the foreign IP address doesn't match the source IP address.

As the next example, consider an incoming UDP datagram destined for 140.252.1.29 (a unicast address), port 577, from 140.252.1.11, port 1500. Figure 22.15 shows to which sockets the datagram is delivered.

Local address	Local port	Foreign address	Foreign port	Delivered?
140.252.1.29 140.252.13.63 140.252.13.63 *	577 577 577 577 577	140.252.1.11 140.252.13.35 * *	1500 1500 * *	yes, 0 wildcard matches no, local and foreign IP mismatch no, local IP mismatch no, 2 wildcard matches

Figure 22.15 Received datagram from {140.252.1.11, 1500} to {140.252.1.29, 577}.

The datagram matches the first socket with no wildcard matches and also matches the fourth socket with two wildcard matches. It is delivered to the first socket, the best match.

## 22.6 in\_pcblookup Function

The function in\_pcblookup serves four different purposes.

1. When either TCP or UDP receives an IP datagram, in\_pcblookup scans the protocol's list of Internet PCBs looking for a matching PCB to receive the

datagram. This is transport layer demultiplexing of a received datagram.

- 2. When a process executes the bind system call, to assign a local IP address and local port number to a socket, in\_pcbbind is called by the protocol to verify that the requested local address pair is not already in use.
- 3. When a process executes the bind system call, requesting an ephemeral port be assigned to its socket, the kernel picks an ephemeral port and calls in\_pcbbind to check if the port is in use. If it is in use, the next ephemeral port number is tried, and so on, until an unused port is located.
- 4. When a process executes the connect system call, either explicitly or implicitly, in\_pcbbind verifies that the requested socket pair is unique. (An implicit call to connect happens when a UDP datagram is sent on an unconnected socket. We'll see this scenario in Chapter 23.)

In cases 2, 3, and 4 in\_pcbbind calls in\_pcblookup. Two options confuse the logic of the function. First, a process can specify either the SO\_REUSEADDR or SO\_REUSEPORT socket option to say that a duplicate local address is OK.

Second, sometimes a wildcard match is OK (e.g., an incoming UDP datagram can match a PCB that has a wildcard for its local IP address, meaning that the socket will accept UDP datagrams that arrive on any local interface), while other times a wildcard match is forbidden (e.g., when connecting to a foreign IP address and port number).

In the original Stanford IP multicast code appears the comment that "The logic of in\_pcblookup is rather opaque and there is not a single comment, ...." The adjective *opaque* is an understatement.

The publicly available IP multicast code available for BSD/386, which is derived from the port to 4.4BSD done by Craig Leres, fixed the overloaded semantics of this function by using in\_pcblookup only for case 1 above. Cases 2 and 4 are handled by a new function named in\_pcbconflict, and case 3 is handled by a new function named in\_uniqueport. Dividing the original functionality into separate functions is much clearer, but in the Net/3 release, which we're describing in this text, the logic is still combined into the single function in\_pcblookup.

## Figure 22.16 shows the in\_pcblookup function.

The function starts at the head of the protocol's PCB list and potentially goes through every PCB on the list. The variable match remembers the pointer to the entry with the best match so far, and matchwild remembers the number of wildcards in that match. The latter is initialized to 3, which is a value greater than the maximum number of wildcard matches that can be encountered. (Any value greater than 2 would work.) Each time around the loop, the variable wildcard starts at 0 and counts the number of wildcard matches for each PCB.

#### Compare local port number

<sup>416–417</sup> The first comparison is the local port number. If the PCB's local port doesn't match the lport argument, the PCB is ignored.

Chapter 22

342

14

```
- in_pcb.c
405 struct inpcb *
406 in_pcblookup(head, faddr, fport_arg, laddr, lport_arg, flags)
407 struct inpcb *head;
408 struct in_addr faddr, laddr;
409 u_int fport_arg, lport_arg;
410 int
            flags;
411 {
        struct inpcb *inp, *match = 0;
412
413
        int
               matchwild = 3, wildcard;
414
        u_short fport = fport_arg, lport = lport_arg;
        for (inp = head->inp_next; inp != head; inp = inp->inp_next) {
415
416
            if (inp->inp_lport != lport)
417
                continue;
                                     /* ignore if local ports are unequal */
418
            wildcard = 0;
419
            if (inp->inp_laddr.s_addr != INADDR_ANY) {
                if (laddr.s_addr == INADDR_ANY)
420
421
                    wildcard++;
422
                else if (inp->inp_laddr.s_addr != laddr.s_addr)
423
                    continue;
424
            } else {
425
                if (laddr.s_addr != INADDR_ANY)
426
                    wildcard++;
427
            }
428
            if (inp->inp_faddr.s_addr != INADDR_ANY) {
429
                 if (faddr.s_addr == INADDR_ANY)
430
                    wildcard++;
431
                 else if (inp->inp_faddr.s_addr != faddr.s_addr ||
432
                          inp->inp_fport != fport)
433
                     continue;
434
            } else {
435
                 if (faddr.s_addr != INADDR_ANY)
436
                    wildcard++;
437
             }
438
            if (wildcard && (flags & INPLOOKUP_WILDCARD) == 0)
439
                 continue;
                                     /* wildcard match not allowed */
440
            if (wildcard < matchwild) {
441
                match = inp;
                 matchwild = wildcard;
442
443
                 if (matchwild == 0)
                                      /* exact match, all done */
444
                    break:
445
             }
446
         }
447
        return (match);
448 }
                                                                            - in_pcb.c
```

Figure 22.16 in\_pcblookup function: search all the PCBs for a match.

#### Compare local address

419-427 in\_pcblookup compares the local address in the PCB with the laddr argument. If one is a wildcard and the other is not a wildcard, the wildcard counter is incremented. If both are not wildcards, then they must be the same, or this PCB is ignored. If both are wildcards, nothing changes: they can't be compared and the wildcard counter isn't incremented. Figure 22.17 summarizes the four different conditions.

PCB local IP	laddr argument	Description
not *	*	wildcard++
not *	not *	compare IP addresses, skip PCB if not equal
*	*	can't compare
*	not *	wildcard++

Figure 22.17 Four scenarios for the local IP address comparison done by in\_pcblookup.

### Compare foreign address and foreign port number

<sup>428-437</sup> These lines perform the same test that we just described, but using the foreign addresses instead of the local addresses. Also, if both foreign addresses are not wildcards then not only must the two IP addresses be equal, but the two foreign ports must also be equal. Figure 22.18 summarizes the foreign IP comparisons.

PCB foreign IP	faddr argument	Description
not *	*	wildcard++
not *	not *	compare IP addresses and ports, skip PCB if not equal
*	*	can't compare
*	not *	wildcard++

Figure 22.18 Four scenarios for the foreign IP address comparison done by in\_pcblookup.

The additional comparison of the foreign port numbers can be performed for the second line of Figure 22.18 because it is not possible to have a PCB with a nonwildcard foreign address and a foreign port number of 0. This restriction is enforced by connect, which we'll see shortly requires a nonwildcard foreign IP address and a nonzero foreign port. It is possible, however, and common, to have a wildcard local address with a nonzero local port. We saw this in Figures 22.10 and 22.13.

### Check if wildcard match allowed

438-439

The flags argument can be set to INPLOOKUP\_WILDCARD, which means a match containing wildcards is OK. If a match is found containing wildcards (wildcard is nonzero) and this flag was not specified by the caller, this PCB is ignored. When TCP and UDP call this function to demultiplex an incoming datagram, INPLOOKUP\_WILDCARD is always set, since a wildcard match is OK. (Recall our examples using Figures 22.10 and 22.13.) But when this function is called as part of the connect system call, in order to verify that a socket pair is not already in use, the flags argument is set to 0.

- 14

## \_\_\_

## Remember best match, return if exact match found

440-447 These statements remember the best match found so far. Again, the best match is considered the one with the fewest number of wildcard matches. If a match is found with one or two wildcards, that match is remembered and the loop continues. But if an exact match is found (wildcard is 0), the loop terminates, and a pointer to the PCB with that exact match is returned.

## Example—Demultiplexing of Received TCP Segment

Figure 22.19 is from the TCP example we discussed with Figure 22.11. Assume in\_pcblookup is demultiplexing a received datagram from 140.252.1.11, port 1500, destined for 140.252.1.29, port 23. Also assume that the order of the PCBs is the order of the rows in the figure. laddr is the destination IP address, lport is the destination TCP port, faddr is the source IP address, and fport is the source TCP port.

	wildcard			
Local address	wilucalu			
140.252.1.29	23	*	*	1
*	23	*	*	2
140.252.1.29	23	140.252.1.11	1500	0

Figure 22.19	laddr = 140.252.1.29, lport = 23, faddr = 140.252.1.11, fport = 1500.
--------------	---

When the first row is compared to the incoming segment, wildcard is 1 (the foreign IP address), flags is set to INPLOOKUP\_WILDCARD, so match is set to point to this PCB and matchwild is set to 1. The loop continues since an exact match has not been found yet. The next time around the loop, wildcard is 2 (the local and foreign IP addresses) and since this is greater than matchwild, the entry is not remembered, and the loop continues. The next time around the loop, wildcard is 0, which is less than matchwild (1), so this entry is remembered in match. The loop also terminates since an exact match has been found and the pointer to this PCB is returned to the caller.

If in\_pcblookup were used by TCP and UDP only to demultiplex incoming datagrams, it could be simplified. First, there's no need to check whether the faddr or laddr arguments are wildcards, since these are the source and destination IP addresses from the received datagram. Also the flags argument could be removed, along with its corresponding test, since wildcard matches are always OK.

This section has covered the mechanics of the in\_pcblookup function. We'll return to this function and discuss its meaning after seeing how it is called from the in\_pcbbind and in\_pcbconnect functions.

## 22.7 in\_pobbind Function

The next function, in\_pcbbind, binds a local address and port number to a socket. It is called from five functions:

- 1. from bind for a TCP socket (normally to bind a server's well-known port);
- 2. from bind for a UDP socket (either to bind a server's well-known port or to bind an ephemeral port to a client's socket);
- from connect for a TCP socket, if the socket has not yet been bound to a nonzero port (this is typical for TCP clients);
- 4. from listen for a TCP socket, if the socket has not yet been bound to a nonzero port (this is rare, since listen is called by a TCP server, which normally binds a well-known port, not an ephemeral port); and
- 5. from in\_pebconnect (Section 22.8), if the local IP address and local port number have not been set (typical for a call to connect for a UDP socket or for each call to sendto for an unconnected UDP socket).

In cases 3, 4, and 5, an ephemeral port number is bound to the socket and the local IP address is not changed (in case it is already set).

We call cases 1 and 2 *explicit binds* and cases 3, 4, and 5 *implicit binds*. We also note that although it is normal in case 2 for a server to bind a well-known port, servers invoked using remote procedure calls (RPC) often bind ephemeral ports and then register their ephemeral port with another program that maintains a mapping between the server's RPC program number and its ephemeral port (e.g., the Sun port mapper described in Section 29.4 of Volume 1).

We'll show the in\_pobbind function in three sections. Figure 22.20 is the first section.

```
in_pcb.c
52 int
53 in_pcbbind(inp, nam)
54 struct inpcb *inp;
55 struct mbuf *nam;
56 {
       struct socket *so = inp->inp_socket;
57
       struct inpcb *head = inp->inp_head;
58
       struct sockaddr_in *sin;
59
                                   /* XXX */
       struct proc *p = curproc;
60
       u_short lport = 0;
61
               wild = 0, reuseport = (so->so_options & SO_REUSEPORT);
62
       int
               error;
63
       int
       if (in_ifaddr == 0)
64
           return (EADDRNOTAVAIL);
65
       if (inp->inp_lport || inp->inp_laddr.s_addr != INADDR_ANY)
66
67
           return (EINVAL);
       if ((so->so_options & (SO_REUSEADDR | SO_REUSEPORT)) == 0 &&
68
           ((so->so_proto->pr_flags & PR_CONNREQUIRED) == 0 ||
69
             (so->so_options & SO_ACCEPTCONN) == 0))
70
            wild = INPLOOKUP_WILDCARD;
71
                                                                           -in_pcb.c
```

Figure 22.20 in\_pcbbind function: bind a local address and port number.

art a transmission and a transmission of the second second second second second second second second second se Second 68-71

72-75

<sup>64–67</sup> The first two tests verify that at least one interface has been assigned an IP address and that the socket is not already bound. You can't bind a socket twice.

This if statement is confusing. The net result sets the variable wild to INPLOOKUP\_WILDCARD if neither SO\_REUSEADDR or SO\_REUSEPORT are set.

The second test is true for UDP sockets since PR\_CONNREQUIRED is false for connectionless sockets and true for connection-oriented sockets.

The third test is where the confusion lies [Torek 1992]. The socket flag SO\_ACCEPTCONN is set only by the listen system call (Section 15.9), which is valid only for a connection-oriented server. In the normal scenario, a TCP server calls socket, bind, and then listen. Therefore, when in\_pcbbind is called by bind, this socket flag is cleared. Even if the process calls socket and then listen, without calling bind, TCP's PRU\_LISTEN request calls in\_pcbbind to assign an ephemeral port to the socket *before* the socket layer sets the SO\_ACCEPTCONN flag. This means the third test in the if statement, testing whether SO\_ACCEPTCONN is not set, is always true. The if statement is therefore equivalent to

if ((so->so\_options & (SO\_REUSEADDR|SO\_REUSEPORT)) == 0 &&
 ((so->so\_proto->pr\_flags & PR\_CONNREQUIRED) == 0 || 1)
 wild = INPLOOKUP\_WILDCARD;

Since anything logically ORed with 1 is always true, this is equivalent to

if ((so->so\_options & (SO\_REUSEADDR|SO\_REUSEPORT)) == 0)
wild = INPLOOKUP\_WILDCARD;

which is simpler to understand: if either of the REUSE socket options is set, wild is left as 0. If neither of the REUSE socket options are set, wild is set to INPLOOKUP\_WILDCARD. In other words, when in\_pcblookup is called later in the function, a wildcard match is allowed only if *neither* of the REUSE socket options are on.

The next section of the in\_pobbind, shown in Figure 22.22, function processes the optional nam argument.

The nam argument is a nonnull pointer only when the process calls bind explicitly.

For an implicit bind (a side effect of connect, listen, or in\_pcbconnect, cases 3, 4, and 5 from the beginning of this section), nam is a null pointer. When the argument is specified, it is an mbuf containing a sockaddr\_in structure. Figure 22.21 shows the four cases for the nonnull nam argument.

nam argument:		PCB meml	per gets set to:	Comment	
localIP	lport	inp_laddr	inp_lport	Continent	
not *	0	localIP	ephemeral port	localIP must be local interface	
not *	nonzero	localIP	lport	subject to in_pcblookup	
*	0	*	ephemeral port		
*	nonzero	*	lport	subject to in_pcblookup	

Figure 22.21 Four cases for nam argument to in\_pcbbind.

<sup>76-83</sup> The test for the correct address family is commented out, yet the identical test in the in\_pcbconnect function (Figure 22.25) is performed. We expect either both to be in or both to be out.

72 if (nam) { in\_pcb.c 73 sin = mtod(nam, struct sockaddr\_in \*); 74 if (nam->m\_len != sizeof(\*sin)) 75 return (EINVAL); 76 #ifdef notdef 77 /\* 78 \* We should check the family, but old programs \* incorrectly fail to initialize it. 79 80 \*/ 81 if (sin->sin\_family != AF\_INET) 82 return (EAFNOSUPPORT); 83 #endif 84 lport = sin->sin\_port; /\* might be 0 \*/ if (IN\_MULTICAST(ntohl(sin->sin\_addr.s\_addr))) { 85 86 1\* 87 \* Treat SO\_REUSEADDR as SO\_REUSEPORT for multicast; \* allow complete duplication of binding if 88 \* SO\_REUSEPORT is set, or if SO\_REUSEADDR is set 89 90 \* and a multicast address is bound on both 91 \* new and duplicated sockets. 92 \*/ 93 if (so->so\_options & SO\_REUSEADDR) 94 reuseport = SO\_REUSEADDR | SO\_REUSEPORT; 95 } else if (sin->sin\_addr.s\_addr != INADDR\_ANY) { 96 sin->sin\_port = 0; /\* yech... \*/ 97 if (ifa\_ifwithaddr((struct sockaddr \*) sin) == 0) 98 return (EADDRNOTAVAIL); 99 } 100 if (lport) { 101 struct inpcb \*t; 102 /\* GROSS \*/ 103 if (ntohs(lport) < IPPORT\_RESERVED && 104 (error = suser(p->p\_ucred, &p->p\_acflag))) 105 return (error); 106 t = in\_pcblookup(head, zeroin\_addr, 0, 107 sin->sin\_addr, lport, wild); 108 if (t && (reuseport & t->inp\_socket->so\_options) == 0) 109 return (EADDRINUSE); 110 } 111 inp->inp\_laddr = sin->sin\_addr; /\* might be wildcard \*/ 112 -in pcb.c

Figure 22.22 in\_pobbind function: process optional nam argument.

Net/3 tests whether the IP address being bound is a multicast group. If so, the SO\_REUSEADDR option is considered identical to SO\_REUSEPORT.
 Otherwise, if the local address being bound by the state of the local address being bound by the state.

Otherwise, if the local address being bound by the caller is not the wildcard, ifa\_ifwithaddr verifies that the address corresponds to a local interface.

The comment "yech" is probably because the port number in the socket address structure must be 0 because ifa\_ifwithaddr does a binary comparison of the entire structure, not just a comparison of the IP addresses.

This is one of the few instances where the process *must* zero the socket address structure before issuing the system call. If bind is called and the final 8 bytes of the socket address structure (sin\_zero[8]) are nonzero, ifa\_ifwithaddr will not find the requested interface, and in\_pcbbind will return an error.

100-105

The next if statement is executed when the caller is binding a nonzero port, that is, the process wants to bind one particular port number (the second and fourth scenarios from Figure 22.21). If the requested port is less than 1024 (IPPORT\_RESERVED) the process must have superuser privilege. This is not part of the Internet protocols, but a Berkeley convention. A port number less than 1024 is called a *reserved port* and is used, for example, by the rcmd function [Stevens 1990], which in turn is used by the rlogin and rsh client programs as part of their authentication with their servers.

106-109

The function in\_pcblookup (Figure 22.16) is then called to check whether a PCB already exists with the same local IP address and local port number. The second argument is the wildcard IP address (the foreign IP address) and the third argument is a port number of 0 (the foreign port). The wildcard value for the second argument causes in\_pcblookup to ignore the foreign IP address and foreign port in the PCB—only the local IP address and local port are compared to sin->sin\_addr and lport, respectively. We mentioned earlier that wild is set to INPLOOKUP\_WILDCARD only if neither of the REUSE socket options are set.

111

The caller's value for the local IP address is stored in the PCB. This can be the wildcard address, if that's the value specified by the caller. In this case the local IP address is chosen by the kernel, but not until the socket is connected at some later time. This is because the local IP address is determined by IP routing, based on foreign IP address.

The final section of in\_pcbbind handles the assignment of an ephemeral port when the caller explicitly binds a port of 0, or when the nam argument is a null pointer (an implicit bind).

in pcb.c 113 if (lport == 0) 114 do { 115 if (head->inp\_lport++ < IPPORT\_RESERVED || 116 head~>inp\_lport > IPPORT\_USERRESERVED) 117 head->inp\_lport = IPPORT\_RESERVED; 118 lport = htons(head->inp\_lport); 119 } while (in\_pcblookup(head, 120 zeroin\_addr, 0, inp->inp\_laddr, lport, wild)); 121 inp->inp\_lport = lport; 122 return (0); 123 } in pcb.c

Figure 22.23 in\_pcbbind function: choose an ephemeral port.

113-122

The next ephemeral port number to use for this protocol (TCP or UDP) is maintained in the head of the protocol's PCB list: tcb or udb. Other than the inp\_next and inp\_back pointers in the protocol's head PCB, the only other element of the inpcb structure that is used is the local port number. Confusingly, this local port number is maintained in host byte order in the head PCB, but in network byte order in all the other PCBs on the list! The ephemeral port numbers start at 1024 が確認を行いていた。

(IPPORT\_RESERVED) and get incremented by 1 until port 5000 is used (IPPORT\_USERRESERVED), then cycle back to 1024. The loop is executed until in\_pcbbind does not find a match.

### SO\_REUSEADDR Examples

Let's look at some common examples to see the interaction of in\_pcbbind with in\_pcblookup and the two REUSE socket options.

1. A TCP or UDP server normally starts by calling socket and bind. Assume a TCP server that calls bind, specifying the wildcard IP address and its nonzero well-known port, say 23 (the Telnet server). Also assume that the server is not already running and that the process does not set the SO\_REUSEADDR socket option.

in\_pcbbind calls in\_pcblookup with INPLOOKUP\_WILDCARD as the final argument. The loop in in\_pcblookup won't find a matching PCB, assuming no other process is using the server's well-known TCP port, causing a null pointer to be returned. This is OK and in\_pcbbind returns 0.

2. Assume the same scenario as above, but with the server already running when someone tries to start the server a second time.

When in\_pcblookup is called it finds the PCB with a local socket of {\*, 23}. Since the wildcard counter is 0, in\_pcblookup returns the pointer to this entry. Since reuseport is 0, in\_pcbbind returns EADDRINUSE.

3. Assume the same scenario as the previous example, but when the attempt is made to start the server a second time, the SO\_REUSEADDR socket option is specified.

Since this socket option is specified, in\_pcbbind calls in\_pcblookup with a final argument of 0. But the PCB with a local socket of {\*, 23} is still matched and returned because wildcard is 0, since in\_pcblookup cannot compare the two wildcard addresses (Figure 22.17). in\_pcbbind again returns EADDRINUSE, preventing us from starting two instances of the server with identical local sockets, regardless of whether we specify SO\_REUSEADDR or not.

4. Assume that a Telnet server is already running with a local socket of {\*, 23} and we try to start another with a local socket of {140.252.13.35, 23}.

Assuming SO\_REUSEADDR is not specified, in\_pcblookup is called with a final argument of INPLOOKUP\_WILDCARD. When it compares the PCB containing \*.23, the counter wildcard is set to 1. Since a wildcard match is allowed, this match is remembered as the best match and a pointer to it is returned after all the TCP PCBs are scanned. in\_pcbbind returns EADDRINUSE.

5. This example is the same as the previous one, but we specify the SO\_REUSEADDR socket option for the second server that tries to bind the local socket {140.252.13.35, 23}.

The final argument to in\_pcblookup is now 0, since the socket option is specified. When the PCB with the local socket {\*, 23} is compared, the wildcard counter is 1, but since the final flags argument is 0, this entry is skipped and is not remembered as a match. After comparing all the TCP PCBs, the function returns a null pointer and in publind returns 0.

6. Assume the first Telnet server is started with a local socket of {140.252.13.35, 23} when we try to start a second server with a local socket of {\*, 23}. This is the same as the previous example, except we're starting the servers in reverse order this time.

The first server is started without a problem, assuming no other socket has already bound port 23. When we start the second server, the final argument to in\_pcblookup is INPLOOKUP\_WILDCARD, assuming the SO\_REUSEADDR socket option is not specified. When the PCB with the local socket of {140.252.13.35, 23} is compared, the wildcard counter is set to 1 and this entry is remembered. After all the TCP PCBs are compared, the pointer to this entry is returned, causing in\_pcbbind to return EADDRINUSE.

7. What if we start two instances of a server, both with a nonwildcard local IP address? Assume we start the first Telnet server with a local socket of {140.252.13.35, 23} and then try to start a second with a local socket of {127.0.0.1, 23}, without specifying SO\_REUSEADDR.

When the second server calls in\_pcbbind, it calls in\_pcblookup with a final argument of INPLOOKUP\_WILDCARD. When the PCB with the local socket of {140.252.13.35, 23} is compared, it is skipped because the local IP addresses are not equal. in\_pcblookup returns a null pointer, and in\_pcbbind returns 0.

From this example we see that the SO\_REUSEADDR socket option has no effect on nonwildcard IP addresses. Indeed the test on the flags value INPLOOKUP\_WILDCARD in in\_pcblookup is made only when wildcard is greater than 0, that is, when either the PCB entry has a wildcard IP address or the IP address being bound is the wildcard.

8. As a final example, assume we try to start two instances of the same server, both with the same nonwildcard local IP address, say 127.0.0.1.

When the second server is started, in\_pcblookup always returns a pointer to the matching PCB with the same local socket. This happens regardless of the SO\_REUSEADDR socket option, because the wildcard counter is always 0 for this comparison. Since in\_pcblookup returns a nonnull pointer, in\_pcbbind returns EADDRINUSE.

From these examples we can state the rules about the binding of local IP addresses and the SO\_REUSEADDR socket option. These rules are shown in Figure 22.24. We assume that *localIP1* and *localIP2* are two different unicast or broadcast IP addresses valid on the local host, and that *localmcastIP* is a multicast group. We also assume that the process is trying to bind the same nonzero port number that is already bound to the existing PCB.

We need to differentiate between a unicast or broadcast address and a multicast address, because we saw that in\_pcbbind considers SO\_REUSEADDR to be the same as SO\_REUSEPORT for a multicast address.

#### in\_pcbconnect Function 735

Existing PCB	Try to bind	SO_REUSEADDR		Description	
Existing I CD	ITY to billa	off	on	Description	
localIP1	localIP1	error	error	one server per IP address and port	
localIP1	localIP2	OK	OK	one server for each local interface	
localIP1	*	error	OK	one server for one interface, other server for remaining interfaces	
*	localIP1	error	OK	one server for one interface, other server for remaining interfaces	
*	*	error	error	can't duplicate local sockets (same as first example)	
localmcastIP	localmcastIP	error	ОК	multiple multicast recipients	

Figure 22.24 Effect of SO\_REUSEADDR socket option on binding of local IP address.

#### **SO\_REUSEPORT** Socket Option

The handling of SO\_REUSEPORT in Net/3 changes the logic of in\_pcbbind to allow duplicate local sockets as long as both sockets specify SO\_REUSEPORT. In other words, all the servers must agree to share the same local port.

## 22.8 in\_pcbconnect Function

The function in\_pcbconnect specifies the foreign IP address and foreign port number for a socket. It is called from four functions:

- 1. from connect for a TCP socket (required for a TCP client);
- 2. from connect for a UDP socket (optional for a UDP client, rare for a UDP server);
- 3. from sendto when a datagram is output on an unconnected UDP socket (common); and
- 4. from tcp\_input when a connection request (a SYN segment) arrives on a TCP socket that is in the LISTEN state (standard for a TCP server).

In all four cases it is common, though not required, for the local IP address and local port be unspecified when in\_pcbconnect is called. Therefore one function of in\_pcbconnect is to assign the local values when they are unspecified.

We'll discuss the in\_pcbconnect function in four sections. Figure 22.25 shows the first section.

in\_pcb.c

```
130 int
131 in_pcbconnect(inp, nam)
132 struct inpcb *inp;
133 struct mbuf *nam;
134 {
135 struct in_ifaddr *ia;
136 struct sockaddr_in *ifaddr;
137 struct sockaddr_in *sin = mtod(nam, struct sockaddr_in *);
```

**INTEL Ex.1013.761** 

Chapter 22

138	if	(nam->m_len != sizeof(*sin))
139		return (EINVAL);
140	if	(sin->sin_family != AF_INET)
141		return (EAFNOSUPPORT);
142	if	(sin->sin_port == 0)
143		return (EADDRNOTAVAIL);
144	if	(in_ifaddr) {
145		/*
146		* If the destination address is INADDR_ANY,
147		* use the primary local address.
148		* If the supplied address is INADDR_BROADCAST,
149		* and the primary interface supports broadcast,
150		* choose the broadcast address for that interface.
151		*/
152	#define	<pre>satosin(sa) ((struct sockaddr_in *)(sa))</pre>
153	#define	sintosa(sin) ((struct sockaddr *)(sin))
154	#define	ifatoia(ifa) ((struct in_ifaddr *)(ifa))
155		if (sin->sin_addr.s_addr == INADDR_ANY)
156		<pre>sin-&gt;sin_addr = IA_SIN(in_ifaddr)-&gt;sin_addr;</pre>
157		else if (sin->sin_addr.s_addr == (u_long) INADDR_BROADCAST &&
158		(in_ifaddr->ia_ifp->if_flags & IFF BROADCAST))
159		<pre>sin-&gt;sin_addr = satosin(∈_ifaddr-&gt;ia_broadaddr)-&gt;sin_addr;</pre>
160	}	

Figure 22.25 in\_pcbconnect function: verify arguments, check foreign IP address.

#### Validate argument

130-143

The nam argument points to an mbuf containing a sockaddr\_in structure with the foreign IP address and port number. These lines validate the argument and verify that the caller is not trying to connect to a port number of 0.

## Handle connection to 0.0.0.0 and 255.255.255.255 specially

144-160

The test of the global in\_ifaddr verifies that an IP interface has been configured. If the foreign IP address is 0.0.0.0 (INADDR\_ANY), then 0.0.0.0 is replaced with the IP address of the primary IP interface. This means the calling process is connecting to a peer on this host. If the foreign IP address is 255.255.255.255.255.255 (INADDR\_BROADCAST) and the primary interface supports broadcasting, then 255.255.255.255 is replaced with the broadcast address of the primary interface. This allows a UDP application to broadcast on the primary interface without having to figure out its IP address—it can simply send datagrams to 255.255.255.255, and the kernel converts this to the appropriate IP address for the interface.

The next section of code, Figure 22.26, handles the case of an unspecified local address. This is the common scenario for TCP and UDP clients, cases 1, 2, and 3 from the list at the beginning of this section.

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179 180

181

182

183

184

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199 200

201

202

203

204

205

in\_pcb.c

/\* XXX \*/

-in\_pcb.c

「「「「「「「」」

22

.с

е

t

)

1

)

1

~

if ((inp->inp\_socket->so\_options & SO\_DONTROUTE) == 0 && (ro->ro\_rt == (struct rtentry \*) 0 || ro->ro\_rt->rt\_ifp == (struct ifnet \*) 0)) { /\* No route yet, so try to acquire one \*/ ro->ro\_dst.sa\_family = AF\_INET; ro->ro\_dst.sa\_len = sizeof(struct sockaddr\_in); ((struct sockaddr\_in \*) &ro->ro\_dst)->sin\_addr = sin->sin\_addr; rtalloc(ro); } /\* \* If we found a route, use the address corresponding to the outgoing interface \* unless it is the loopback (in case a route \* to our address on another net goes to loopback). \*/ if (ro->ro\_rt && !(ro->ro\_rt->rt\_ifp->if\_flags & IFF\_LOOPBACK)) ia = ifatoia(ro->ro\_rt->rt\_ifa); if (ia == 0) { u\_short fport = sin->sin\_port; sin->sin\_port = 0; ia = ifatoia(ifa\_ifwithdstaddr(sintosa(sin))); if (ia == 0)ia = ifatoia(ifa\_ifwithnet(sintosa(sin))); sin->sin\_port = fport;

if (inp->inp\_laddr.s\_addr == INADDR\_ANY) {

sin->sin\_addr.s\_addr ||

ro->ro\_rt = (struct rtentry \*) 0;

RTFREE(ro->ro\_rt);

\* If route is known or can be allocated now,

\* our src addr is taken from the i/f, else punt.

(satosin(&ro->ro\_dst)->sin\_addr.s\_addr !=

inp->inp\_socket->so\_options & SO\_DONTROUTE)) {

ia = (struct in\_ifaddr \*) 0;

struct route \*ro;

ro = &inp->inp\_route;

if (ro->ro\_rt &&

if (ia == 0)

if (ia == 0)

}

ia = in\_ifaddr;

return (EADDRNOTAVAIL);

/\*

\*/

Figure 22.26 in\_pcbconnect function: local IP address not yet specified.

調調

### Release route if no longer valid

164-175

<sup>75</sup> If a route is held by the PCB but the destination of that route differs from the foreign address being connected to, or the SO\_DONTROUTE socket option is set, that route is released.

To understand why a PCB may have an associated route, consider case 3 from the list at the beginning of this section: in\_pcbconnect is called *every time* a UDP datagram is sent on an unconnected socket. Each time a process calls sendto, the UDP output function calls in\_pcbconnect, ip\_output, and in\_pcbdisconnect. If all the datagrams sent on the socket go to the same destination IP address, then the first time through in\_pcbconnect the route is allocated and it can be used from that point on. But since a UDP application can send datagrams to a different IP address with each call to sendto, the destination changes. This same test is done in ip\_output, which seems to be redundant.

The SO\_DONTROUTE socket option tells the kernel to bypass the normal routing decisions and send the IP datagram to the locally attached interface whose IP network address matches the network portion of the destination address.

#### Acquire route

176-185 If the SO\_DONTROUTE socket option is not set, and a route to the destination is not held by the PCB, try to acquire one by calling rtalloc.

#### Determine outgoing interface

186-205

The goal in this section of code is to have ia point to an interface address structure (in\_ifaddr, Section 6.5), which contains the IP address of the interface. If the PCB holds a route that is still valid, or if rtalloc found a route, and the route is not to the loopback interface, the corresponding interface is used. Otherwise ifa\_withdstaddr and ifa\_withnet are called to check if the foreign IP address is on the other end of a point-to-point link or on an attached network. Both of these functions require that the port number in the socket address structure be 0, so it is saved in fport across the calls. If this fails, the primary IP address is used (in\_ifaddr), and if no interfaces are configured (in\_ifaddr is zero), an error is returned.

Figure 22.27 shows the next section of in\_pcbconnect, which handles a destination address that is a multicast address.

206-223

If the destination address is a multicast address and the process has specified the outgoing interface to use for multicast packets (using the IP\_MULTICAST\_IF socket option), then the IP address of that interface is used as the local address. A search is made of all IP interfaces for the one matching the interface that was specified with the socket option. An error is returned if that interface is no longer up.

224-225

The code that started at the beginning of Figure 22.26 to handle the case of a wildcard local address is complete. The pointer to the sockaddr\_in structure for the local interface is saved in ifaddr.

The final section of in\_pcblookup is shown in Figure 22.28.

.

.

206	/*	— in_pcb
207	* If the destination address is multicast and an outgoing	•
208	* interface has been set as a multicast option, use the	
209	* address of that interface as our source address.	
210	*/	
211	if (IN_MULTICAST(ntohl(sin->sin_addr.s_addr)) &&	
212	inp->inp_moptions != NULL) {	
213	struct ip_moptions *imo;	
214	struct ifnet *ifp;	
215	<pre>imo = inp-&gt;inp_moptions;</pre>	
216	if (imo->imo_multicast_ifp != NULL) {	
217	<pre>ifp = imo-&gt;imo_multicast_ifp;</pre>	
218	for (ia = in_ifaddr; ia; ia = ia->ia_next)	
219	if (ia->ia_ifp == ifp)	
220	break;	
221	if (ia == 0)	
222	return (EADDRNOTAVAIL);	
223	}	
224	}	
225	ifaddr = (struct sockaddr_in *) &ia->ia_addr;	
226	)	

Figure 22.27 in\_pebconnect function: destination address is a multicast address.

```
227
        if (in_pcblookup(inp->inp_head,
                                                                              in_pcb.c
228
                          sin->sin_addr,
229
                          sin->sin_port,
230
                     inp->inp_laddr.s_addr ? inp->inp_laddr : ifaddr->sin_addr,
231
                          inp->inp_lport,
232
                          0))
233
            return (EADDRINUSE);
234
        if (inp->inp_laddr.s_addr == INADDR_ANY) {
            if (inp->inp_lport == 0)
235
236
                 (void) in_pcbbind(inp, (struct mbuf *) 0);
237
            inp->inp_laddr = ifaddr->sin_addr;
238
        }
239
        inp->inp_faddr = sin->sin_addr;
240
        inp->inp_fport = sin->sin_port;
241
        return (0);
242 }
                                                                             - in_pcb.c
```

Figure 22.28 in\_pcbconnect function: verify that socket pair is unique.

#### Verify that socket pair is unique

in\_pcblookup verifies that the socket pair is unique. The foreign address and for-227-233 eign port are the values specified as arguments to in\_pcbconnect. The local address is either the value that was already bound to the socket or the value in ifaddr that was

n\_pcb.c

Secti

22.

calculated in the code we just described. The local port can be 0, which is typical for a TCP client, and we'll see that later in this section of code an ephemeral port is chosen

This test prevents two TCP connections to the same foreign address and foreign for the local port. port from the same local address and local port. For example, if we establish a TCP connection with the echo server on the host sun and then try to establish another connection to the same server from the same local port (8888, specified with the -b option), the call to in\_pcblookup returns a match, causing connect to return the error EADDRINUSE. (We use the sock program from Appendix C of Volume 1.)

start first one in the background bsdi \$ sock -b 8888 sun echo & bsdi \$ sock -A -b 8888 sun echo then try again connect() error: Address already in use

We specify the -A option to set the SO\_REUSEADDR socket option, which lets the bind succeed, but the connect cannot succeed. This is a contrived example, as we explicitly bound the same local port (8888) to both sockets. In the normal scenario of two different clients from the host bsdi to the echo server on the host sun, the local port will be 0 when the second client calls in\_pcblookup from Figure 22.28.

This test also prevents two UDP sockets from being connected to the same foreign

address from the same local port. This test does not prevent two UDP sockets from alternately sending datagrams to the same foreign address from the same local port, as long as neither calls connect, since a UDP socket is only temporarily connected to a peer for the duration of a sendto system call.

## Implicit bind and assignment of ephemeral port

If the local address is still wildcarded for the socket, it is set to the value saved in ifaddr. This is an implicit bind: cases 3, 4, and 5 from the beginning of Section 22.7. 234--238 First a check is made as to whether the local port has been bound yet, and if not, in\_pcbbind binds an ephemeral port to the socket. The order of the call to in\_pcbbind and the assignment to inp\_laddr is important, since in\_pcbbind fails if the local address is not the wildcard address.

## Store foreign address and foreign port in PCB

239-240

The final step of this function sets the foreign IP address and foreign port number in the PCB. We are guaranteed, on successful return from this function, that both socket pairs in the PCB—the local and foreign—are filled in with specific values.

## IP Source Address Versus Outgoing Interface Address

There is a subtle difference between the source address in the IP datagram versus the IP

address of the interface used to send the datagram. The PCB member inp\_laddr is used by TCP and UDP as the source address of the

IP datagram. It can be set by the process to the IP address of any configured interface by bind. (The call to ifa\_ifwithaddr in in\_pobbind verifies the local address desired by the application.) in\_pcbconnect assigns the local address only if it is a wildcard, and when this happens the local address is based on the outgoing interface (since the destination address is known).

2:

in\_setsockaddr and in\_setpeeraddr Functions 741

Section 22.10

はためないです。

247

248

249

250

251 }

The outgoing interface, however, is also determined by ip\_output based on the destination IP address. On a multihomed host it is possible for the source address to be a local interface that is not the outgoing interface, when the process explicitly binds a local address that differs from the outgoing interface. This is allowed because Net/3 chooses the weak end system model (Section 8.4).

#### 22.9 in\_pcbdisconnect Function

A UDP socket is disconnected by in\_pcbdisconnect. This removes the foreign association by setting the foreign IP address to all 0s (INADDR\_ANY) and foreign port number to 0.

This is done after a datagram has been sent on an unconnected UDP socket and when connect is called on a connected UDP socket. In the first case the sequence of steps when the process calls sendto is: UDP calls in\_pcbconnect to connect the socket temporarily to the destination, udp\_output sends the datagram, and then in\_pcbdisconnect removes the temporary connection.

in\_pcbdisconnect is not called when a socket is closed since in\_pcbdetach handles the release of the PCB. A disconnect is required only when the PCB needs to be reused for a different foreign address or port number.

Figure 22.29 shows the function in\_pcbdisconnect.

### 243 int 244 in\_pcbdisconnect(inp) 245 struct inpcb \*inp; 246 { inp->inp\_faddr.s\_addr = INADDR\_ANY; inp->inp\_fport = 0; if (inp->inp\_socket->so\_state & SS\_NOFDREF) in\_pcbdetach(inp);

-in vcb.c

in\_pcb.c

Figure 22.29 in\_pcbdisconnect function: disconnect from foreign address and port number.

If there is no longer a file table reference for this PCB (SS\_NOFDREF is set) then in\_pcbdetach (Figure 22.7) releases the PCB.

# 22.10 in\_setsockaddr and in\_setpeeraddr Functions

The getsockname system call returns the local protocol address of a socket (e.g., the IP address and port number for an Internet socket) and the getpeername system call returns the foreign protocol address. Both system calls end up issuing a PRU\_SOCKADDR request or a PRU\_PEERADDR request. The protocol then calls either in\_setsockaddr or in\_setpeeraddr. We show the first of these in Figure 22.30.

742 Protocol Control Blocks

Chapter 22

in\_pcb.c

in\_pcb.c

in\_pcb.c

\$£) 22

267 int 268 in\_setsockaddr(inp, nam) 269 struct inpcb \*inp; 270 struct mbuf \*nam; 271 { struct sockaddr\_in \*sin; 272 273 nam->m\_len = sizeof(\*sin); sin = mtod(nam, struct sockaddr\_in \*); 274 bzero((caddr\_t) sin, sizeof(\*sin)); 275 sin->sin\_family = AF\_INET; 276 277 sin->sin\_len = sizeof(\*sin); sin->sin\_port = inp->inp\_lport; 278 sin->sin\_addr = inp->inp\_laddr; 279 280 }

Figure 22.30 in\_setsockaddr function: return local address and port number.

The argument nam is a pointer to an mbuf that will hold the result: a sockaddr\_in structure that the system call copies back to the process. The code fills in the socket address structure and copies the IP address and port number from the Internet PCB into the sin\_addr and sin\_port members.

Figure 22.31 shows the in\_setpeeraddr function. It is nearly identical to Figure 22.30, but copies the foreign IP address and port number from the PCB.

```
281 int
282 in_setpeeraddr(inp, nam)
283 struct inpcb *inp;
284 struct mbuf *nam;
285 {
286
        struct sockaddr_in *sin;
287
        nam->m_len = sizeof(*sin);
288
        sin = mtod(nam, struct sockaddr_in *);
289
        bzero((caddr_t) sin, sizeof(*sin));
290
        sin->sin_family = AF_INET;
291
        sin->sin_len = sizeof(*sin);
292
        sin->sin_port = inp->inp_fport;
293
        sin->sin_addr = inp->inp_faddr;
294 }
```

in\_pcb.c

#### Figure 22.31 in\_setpeeraddr function: return foreign address and port number.

# 22.11 in\_pcbnotify, in\_rtchange, and in\_losing Functions

The function in\_pcbnotify is called when an ICMP error is received, in order to notify the appropriate process of the error. The "appropriate process" is found by searching all the PCBs for one of the protocols (TCP or UDP) and comparing the local

.C

с

n t

Э

С

С

)

1

1

and foreign IP addresses and port numbers with the values returned in the ICMP error. For example, when an ICMP source quench error is received in response to a TCP segment that some router discarded, TCP must locate the PCB for the connection that caused the error and slow down the transmission on that connection.

Before showing the function we must review how it is called. Figure 22.32 summarizes the functions called to process an ICMP error. The two shaded ellipses are the functions described in this section.

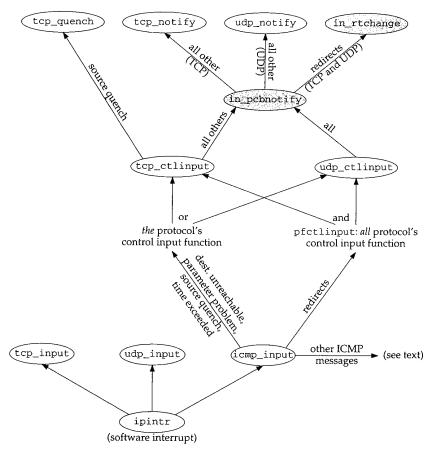


Figure 22.32 Summary of processing of ICMP errors.

When an ICMP message is received, icmp\_input is called. Five of the ICMP messages are classified as errors (Figures 11.1 and 11.2):

- destination unreachable,
- parameter problem,
- redirect,
- source quench, and
- time exceeded.

INTEL Ex.1013.769

#### Chapter 22

Redirects are handled differently from the other four errors. All other ICMP messages (the queries) are handled as described in Chapter 11.

Each protocol defines its control input function, the pr\_ctlinput entry in the protosw structure (Section 7.4). The ones for TCP and UDP are named tcp\_ctlinput and udp\_ctlinput, and we'll show their code in later chapters. Since the ICMP error that is received contains the IP header of the datagram that caused the error, the protocol that caused the error (TCP or UDP) is known. Four of the five ICMP errors cause that protocol's control input function to be called. Redirects are handled differently: the function pfctlinput is called, and it in turn calls the control input functions for all the protocols in the family (Internet). TCP and UDP are the only protocols in the Internet family with control input functions.

Redirects are handled specially because they affect all IP datagrams going to that destination, not just the one that caused the redirect. On the other hand, the other four errors need only be processed by the protocol that caused the error.

The final points we need to make about Figure 22.32 are that TCP handles source quenches differently from the other errors, and redirects are handled specially by in\_pcbnotify: the function in\_rtchange is called, regardless of the protocol that caused the error.

Figure 22.33 shows the in\_pcbnotify function. When it is called by TCP, the first argument is the address of tcb and the final argument is the address of the function tcp\_notify. For UDP, these two arguments are the address of udb and the address of the function udp\_notify.

#### Verify arguments

The cmd argument and the address family of the destination are verified. The for-306-324 eign address is checked to ensure it is not 0.0.0.

# Handle redirects specially

325-338

If the error is a redirect it is handled specially. (The error PRC\_HOSTDEAD is an old error that was generated by the IMPs. Current systems should never see this error-it is a historical artifact.) The foreign port, local port, and local address are all set to 0 so that the for loop that follows won't compare them. For a redirect we want that loop to select the PCBs to receive notification based only on the foreign IP address, because that is the IP address for which our host received a redirect. Also, the function that is called for a redirect is in\_rtchange (Figure 22.34) instead of the notify argument specified by the caller.

339

The global array inetctlerrmap maps one of the protocol-independent error codes (the PRC\_xxx values from Figure 11.19) into its corresponding Unix errno value (the final column in Figure 11.1).

306 int

2

S

e

t

e

e

Р

t

ιt

ιt

r

e

y

it

st

n

١f

t

t

Э

Э

t

ł

£

r

е

in\_pcbnotify, in\_rtchange, and in\_losing Functions 745

– in\_pcb.c

307 in\_pcbnotify(head, dst, fport\_arg, laddr, lport\_arg, cmd, notify) 308 struct inpcb \*head; 309 struct sockaddr \*dst; fport\_arg, lport\_arg; 310 u\_int 311 struct in\_addr laddr; 312 int cmd; 313 void (\*notify) (struct inpcb \*, int); 314 { 315 extern u\_char inetctlerrmap[]; struct inpcb \*inp, \*oinp; 316 317 struct in\_addr faddr; 318 u\_short fport = fport\_arg, lport = lport\_arg; 319 int errno; 320 if ((unsigned) cmd > PRC\_NCMDS || dst->sa\_family != AF\_INET) 321 return; 322 faddr = ((struct sockaddr\_in \*) dst)->sin\_addr; 323 if (faddr.s\_addr == INADDR\_ANY) 324 return: 325 /\*  $\star$  Redirects go to all references to the destination, 326 \* and use in\_rtchange to invalidate the route cache. 327 328 \* Dead host indications: notify all references to the destination. 329  $\star$  Otherwise, if we have knowledge of the local port and address, 330 \* deliver only to that socket. \* / 331 if (PRC\_IS\_REDIRECT(cmd) || cmd == PRC\_HOSTDEAD) { 332 333 fport = 0;334 lport = 0;335 laddr.s\_addr = 0; 336 if (cmd != PRC\_HOSTDEAD) 337 notify = in\_rtchange; 338 } 339 errno = inetctlerrmap[cmd]; 340 for (inp = head->inp\_next; inp != head;) { 341 if (inp->inp\_faddr.s\_addr != faddr.s\_addr !| 342 inp->inp\_socket == 0 || 343 (lport && inp->inp\_lport != lport) || 344 (laddr.s\_addr && inp->inp\_laddr.s\_addr != laddr.s\_addr) || 345 (fport && inp->inp\_fport != fport)) { 346 inp = inp->inp\_next; 347 /\* skip this PCB \*/ continue; 348 } 349 oinp = inp; 350 inp = inp->inp\_next; 351 if (notify) 352 (\*notify) (oinp, errno); 353 } 354 } in pcb.c

Figure 22.33 in\_pcbnot ify function: pass error notification to processes.

in\_pcb.c

in\_pcb.c

Sect

#### Call notify function for selected PCBs

<sup>340–353</sup> This loop selects the PCBs to be notified. Multiple PCBs can be notified—the loop keeps going even after a match is located. The first if statement combines five tests, and if any one of the five is true, the PCB is skipped: (1) if the foreign addresses are unequal, (2) if the PCB does not have a corresponding socket structure, (3) if the local ports are unequal, (4) if the local addresses are unequal, or (5) if the foreign ports are unequal. The foreign addresses *must* match, while the other three foreign and local elements are compared only if the corresponding argument is nonzero. When a match is found, the notify function is called.

#### in\_rtchange Function

We saw that in\_pcbnotify calls the function in\_rtchange when the ICMP error is a redirect. This function is called for all PCBs with a foreign address that matches the IP address that has been redirected. Figure 22.34 shows the in\_rtchange function.

```
391 void
392 in_rtchange(inp, errno)
393 struct inpch *inp.
```

393	struct	<pre>inpcb *inp;</pre>
394	int	errno;
395	{	
396	if	(inp->inp_route.ro_rt) {
397		rtfree(inp->inp_route.ro_rt);
398		<pre>inp-&gt;inp_route.ro_rt = 0;</pre>
399		/*
400		* A new route can be allocated the next time
401		* output is attempted.
402		*/
403	}	
404	}	·

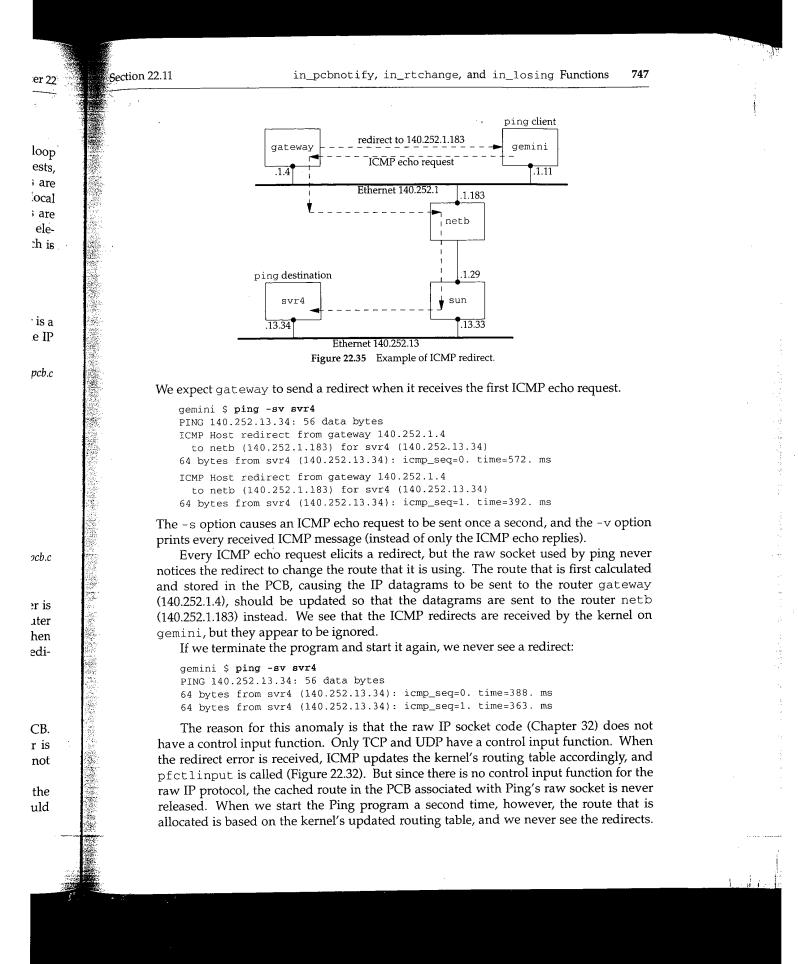
Figure 22.34 in\_rtchange function: invalidate route.

If the PCB holds a route, that route is released by rtfree, and the PCB member is marked as empty. We don't try to update the route at this time, using the new router address returned in the redirect. The new route will be allocated by ip\_output when this PCB is used next, based on the kernel's routing table, which is updated by the redirect, before pfctlinput is called.

#### Redirects and Raw Sockets

Let's examine the interaction of redirects, raw sockets, and the cached route in the PCB. If we run the Ping program, which uses a raw socket, and an ICMP redirect error is received for the IP address being pinged, Ping continues using the original route, not the redirected route. We can see this as follows.

We ping the host svr4 on the 140.252.13 network from the host gemini on the 140.252.1 network. The default router for gemini is gateway, but the packets should be sent to the router netb instead. Figure 22.35 shows the arrangement.



**INTEL Ex.1013.773** 

# Chapter 22

i

#### ICMP Errors and UDP Sockets

One confusing part of the sockets API is that ICMP errors received on a UDP socket are not passed to the application unless the application has issued a connect on the socket, restricting the foreign IP address and port number for the socket. We now see where this limitation is enforced by in\_pcbnotify.

Consider an ICMP port unreachable, probably the most common ICMP error on a UDP socket. The foreign IP address and the foreign port number in the dst argument to in\_pcbnotify are the IP address and port number that caused the ICMP error. But if the process has not issued a connect on the socket, the inp\_faddr and inp\_fport members of the PCB are both 0, preventing in\_pcbnotify from ever calling the notify function for this socket. The for loop in Figure 22.33 will skip every UDP PCB.

This limitation arises for two reasons. First, if the sending process has an unconnected UDP socket, the only nonzero element in the socket pair is the local port. (This assumes the process did not call bind.) This is the only value available to in\_pcbnotify to demultiplex the incoming ICMP error and pass it to the correct process. Although unlikely, there could be multiple processes bound to the same local port, making it ambiguous which process should receive the error. There's also the possibility that the process that sent the datagram that caused the ICMP error has terminated, with another process then starting and using the same local port. This is also unlikely since ephemeral ports are assigned in sequential order from 1024 to 5000 and reused only after cycling around (Figure 22.23).

The second reason for this limitation is because the error notification from the kernel to the process—an errno value—is inadequate. Consider a process that calls sendto on an unconnected UDP socket three times in a row, sending a UDP datagram to three different destinations, and then waits for the replies with recvfrom. If one of the datagrams generates an ICMP port unreachable error, and if the kernel were to return the corresponding error (ECONNREFUSED) to the recvfrom that the process issued, the errno value doesn't tell the process which of the three datagrams caused the error. The kernel has all the information required in the ICMP error, but the sockets API doesn't provide a way to return this to the process.

Therefore the design decision was made that if a process wants to be notified of these ICMP errors on a UDP socket, that socket must be connected to a single peer. If the error ECONNREFUSED is returned on that connected socket, there's no question which peer generated the error.

There is still a remote possibility of an ICMP error being delivered to the wrong process. One process sends the UDP datagram that elicits the ICMP error, but it terminates before the error is received. Another process then starts up before the error is received, binds the same local port, and connects to the same foreign address and foreign port, causing this new process to receive the error. There's no way to prevent this from occurring, given UDP's lack of memory. We'll see that TCP handles this with its TIME\_WAIT state.

oter 22

et are

эсket, vhere

' on a

ment

: But port

g the

UDP

ncon-(This

le to

: proport,

sibil-

ated,

ikely

used

? ker-

calls

gram

ne of

re to

ocess

used

ckets

ed of

er. If

stion

rong

ermi-

or is

l for-

: this

h its

Section 22.11

In our preceding example, one way for the application to get around this limitation is to use three connected UDP sockets instead of one unconnected socket, and call select to determine when any one of the three has a received datagram or an error to be read.

> Here we have a scenario where the kernel has the information but the API (sockets) is inadequate. With most implementations of Unix System V and the other popular API (TLI), the reverse is true: the TLI function t\_rcvuderr can return the peer's IP address, port number, and an error value, but most SVR4 streams implementations of TCP/IP don't provide a way for ICMP to pass the error to an unconnected UDP end point.

> In an ideal world, in\_pcbnotify delivers the ICMP error to all UDP sockets that match, even if the only nonwildcard match is the local port. The error returned to the process would include the destination IP address and destination UDP port that caused the error, allowing the process to determine if the error corresponds to a datagram sent by the process.

#### in\_losing Function

The final function dealing with PCBs is in\_losing, shown in Figure 22.36. It is called by TCP when its retransmission timer has expired four or more times in a row for a given connection (Figure 25.26).

	in pcb.c
361	int
362	in_losing(inp)
363	struct inpcb *inp;
364	{
365	struct rtentry *rt;
366	struct rt_addrinfo info;
367	<pre>if ((rt = inp-&gt;inp_route.ro_rt)) {</pre>
368	<pre>inp_route.ro_rt = 0;</pre>
369	<pre>bzero((caddr_t) &amp; info, sizeof(info));</pre>
370	info.rti_info[RTAX_DST] =
371	(struct sockaddr *) &inp->inp_route.ro_dst;
372	info.rti_info[RTAX_GATEWAY] = rt->rt_gateway;
373	info.rti_info[RTAX_NETMASK] = rt_mask(rt);
374	rt_missmsg(RTM_LOSING, &info, rt->rt_flags, 0);
375	if (rt->rt_flags & RTF_DYNAMIC)
376	<pre>(void) rtrequest(RTM_DELETE, rt_key(rt),</pre>
377	rt->rt_gateway, rt_mask(rt), rt->rt_flags,
378	<pre>(struct rtentry **) 0);</pre>
379	else
380	/*
381	* A new route can be allocated
382	* the next time output is attempted.
383	*/
384	<pre>rtfree(rt);</pre>
385	}
386	in pcb.c

Figure 22.36 in\_losing function: invalidate cached route information.

# 2

Se

# Generate routing message

If the PCB holds a route, that route is discarded. An rt\_addrinfo structure is filled in with information about the cached route that appears to be failing. The func-361-374 tion rt\_missmsg is then called to generate a message from the routing socket of type RTM\_LOSING, indicating a problem with the route.

# Delete or release route

If the cached route was generated by a redirect (RTF\_DYNAMIC is set), the route is deleted by calling rtrequest with a request of RTM\_DELETE. Otherwise the cached 375-384 route is released, causing the next output on the socket to allocate another route to the destination-hopefully a better route.

# 22.12 Implementation Refinements

Undoubtedly the most time-consuming algorithm we've encountered in this chapter is the linear searching of the PCBs done by in\_pcblookup. At the beginning of Section 22.6 we noted four instances when this function is called. We can ignore the calls to bind and connect, as they occur much less frequently than the calls to in\_pcblookup from TCP and UDP, to demultiplex every received IP datagram.

In later chapters we'll see that TCP and UDP both try to help this linear search by maintaining a pointer to the last PCB that the protocol referenced: a one-entry cache. If the local address, local port, foreign address, and foreign port in the cached PCB match the values in the received datagram, the protocol doesn't even call in\_pcblookup. If the protocol's data fits the packet train model [Jain and Routhier 1986], this simple cache works well. But if the data does not fit this model and, for example, looks like data entry into an on-line transaction processing system, the one-entry cache performs poorly [McKenney and Dove 1992].

One proposal for a better PCB arrangement is to move a PCB to the front of the PCB list when the PCB is referenced. ([McKenney and Dove 1992] attribute this idea to Jon Crowcroft; [Partridge and Pink 1993] attribute it to Gary Delp.) This movement of the PCB is easy to do since it is a doubly linked list and a pointer to the head of the list is the first argument to in\_pcblookup.

[McKenney and Dove 1992] compare the original Net/1 implementation (no cache), an enhanced one-entry send-receive cache, the move-to-the-front heuristic, and their own algorithm that uses hash chains. They show that maintaining a linear list of PCBs on hash chains provides an order of magnitude improvement over the other algorithms. The only cost for the hash chains is the memory required for the hash chain headers and the computation of the hash function. They also consider adding the move-to-the-front heuristic to their hash-chain algorithm and conclude that it is easier simply to add more hash chains.

Another comparison of the BSD linear search to a hash table search is in [Hutchinson and Peterson 1991]. They show that the time required to demultiplex an incoming UDP datagram is constant as the number of sockets increases for a hash table, but with a linear search the time increases as the number of sockets increases.

# 22.13 Summary

An Internet PCB is associated with every Internet socket: TCP, UDP, and raw IP. It contains information common to all Internet sockets: local and foreign IP addresses, pointer to a route structure, and so on. All the PCBs for a given protocol are placed on a doubly linked list maintained by that protocol.

In this chapter we've looked at numerous functions that manipulate the PCBs, and three in detail.

1. in\_pcblookup is called by TCP and UDP to demultiplex every received datagram. It chooses which socket receives the datagram, taking into account wildcard matches.

This function is also called by in\_pcbbind to verify that the local address and local process are unique, and by in\_pcbconnect to verify that the combination of a local address, local process, foreign address, and foreign process are unique.

- 2. in\_pebbind explicitly or implicitly binds a local address and local port to a socket. An explicit bind occurs when the process calls bind, and an implicit bind occurs when a TCP client calls connect without calling bind, or when a UDP process calls sendto or connect without calling bind.
- 3. in\_pcbconnect sets the foreign address and foreign process. If the local address has not been set by the process, a route to the foreign address is calculated and the resulting local interface becomes the local address. If the local port has not been set by the process, in\_pcbbind chooses an ephemeral port for the socket.

Figure 22.37 summarizes the common scenarios for various TCP and UDP applications and the values stored in the PCB for the local address and port and the foreign address and port. We have not yet covered all the actions shown in Figure 22.37 for TCP and UDP processes, but will examine the code in later chapters.

752 Protocol Control Blocks

Application	local address: inp_laddr	iocai port.	foreign address	s: foreign port:
TCP client: connect (foreignIP, fpor	in_pcbconnect calls rtalloc to allocate route to foreignIP. Local address is local	inp_lport in_pcbconnect calls in_pcbbind choose ephemeral port.	inp_faddr foreignIP	fport
TCP client: bind (localIP, lport) connect (foreignIP, fport TCP client:	interface. localIP	lport	foreignIP	fport
<pre>bind(*, lport) connect(foreignIP, fport) TCP client:</pre>	foreignIP. Local address is local interface.	lport	foreignIP	fport
<pre>bind(localIP, 0) connect(foreignIP, fport) TCP server:</pre>		in_pcbbind chooses ephemeral port.	foreignIP	fport
<pre>bind(localIP, lport) listen() accept() TCP server:</pre>	localIP	lport	Source address from IP header.	Source port from TCP header.
<pre>bind(*, lport) listen() accept() JDP client:</pre>	Destination address from IP header.	lport	Source address from IP header.	Source port from TCP header.
sendto (foreignIP, fport) DP client:	interface. Reset to 0.0.0.0 after datagram sent.	in_pcbconnect calls in_pcbbind to choose ephemeral port. Not changed on subsequent calls to sendto.	foreignIP. Reset to 0.0.0.0 after datagram sent.	fport. Reset to 0 after datagram sent.
connect(foreignIP, fport) vrite() a iu iu ci	allocate route to deforming for the formed of the formed o	in_pcbconnect calls in_pcbbind to choose ephemeral port. Not changed on subsequent calls to write.	foreignIP j	fport

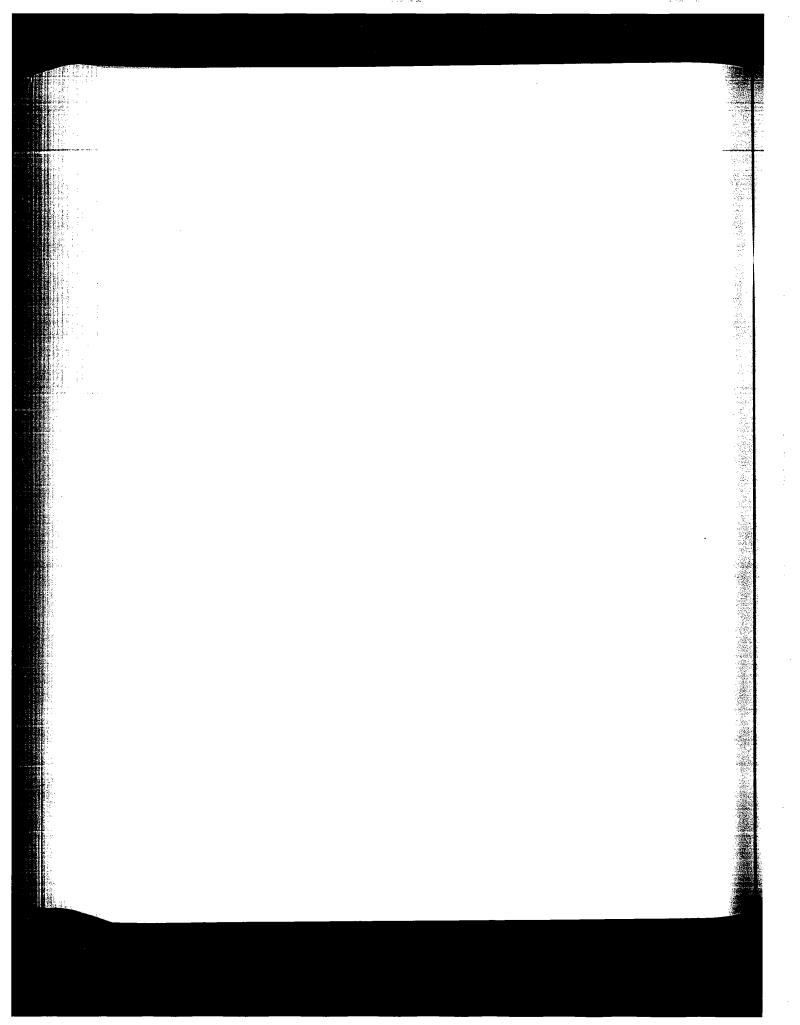
Figure 22.37 Summary of in\_pobbind and in\_pobconnect.

er

22

# **Exercises**

- **22.1** What happens in Figure 22.23 when the process asks for an ephemeral port and every ephemeral port is in use?
- 22.2 In Figure 22.10 we showed two Telnet servers with listening sockets: one with a specific local IP address and one with the wildcard for its local IP address. Does your system's Telnet daemon allow you to specify the local IP address, and if so, how?
- 22.3 Assume a socket is bound to the local socket {140.252.1.29, 8888}, and this is the only socket using local port 8888. (1) Go through the steps performed by in\_pcbbind when another socket is bound to {140.252.13.33, 8888}, without any socket options. (2) Go through the steps performed when another socket is bound to the wildcard IP address, port 8888, without any socket options. (3) Go through the steps performed when another socket is bound to the wildcard IP address, port 8888, with the SO\_REUSEADDR socket option.
- 22.4 What is the first ephemeral port number allocated by UDP?
- 22.5 When a process calls bind, which elements in the sockaddr\_in structure must be filled in?
- 22.6 What happens if a process tries to bind a local broadcast address? What happens if a process tries to bind the limited broadcast address (255.255.255.255)?



# 23

# **UDP: User Datagram Protocol**

# 23.1 Introduction

三部の記述を

The User Datagram Protocol, or UDP, is a simple, datagram-oriented, transport-layer protocol: each output operation by a process produces exactly one UDP datagram, which causes one IP datagram to be sent.

A process accesses UDP by creating a socket of type SOCK\_DGRAM in the Internet domain. By default the socket is termed *unconnected*. Each time the process sends a datagram it must specify the destination IP address and port number. Each time a datagram is received for the socket, the process can receive the source IP address and port number from the datagram.

We mentioned in Section 22.5 that a UDP socket can also be *connected* to one particular IP address and port number. This causes all datagrams written to the socket to go to that destination, and only datagrams arriving from that IP address and port number are passed to the process.

This chapter examines the implementation of UDP.

# 23.2 Code Introduction

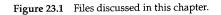
There are nine UDP functions in a single C file and various UDP definitions in two headers, as shown in Figure 23.1.

Figure 23.2 shows the relationship of the six main UDP functions to other kernel functions. The shaded ellipses are the six functions that we cover in this chapter. We also cover three additional UDP functions that are called by some of these six functions.

755

756 UDP: User Datagram Protocol

File	Description
netinet/udp.h netinet/udp_var.h	udphdr structure definition other UDP definitions
netinet/udp_usrreq.c	UDP functions



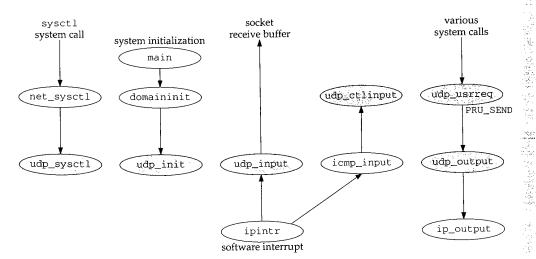


Figure 23.2 Relationship of UDP functions to rest of kernel.

## **Global Variables**

Seven global variables are introduced in this chapter, which are shown in Figure 23.3.

Variable	Datatype	Description
		head of the UDP PCB list pointer to PCB for last received datagram: one-behind cache
udpcksum int		flag for calculating and verifying UDP checksum
udp_in struct sockaddr_in		holds sender's IP address and port on input
udpstat struct udpstat		UDP statistics (Figure 23.4)
udp_recvspace	u_long	default size of socket receive buffer, 41,600 bytes
udp_sendspace	u_long	default size of socket send buffer, 9216 bytes

Figure 23.3 Global variables introduced in this chapter.

いたので、「「「「」」

# Statistics

2

Various UDP statistics are maintained in the global structure udpstat, described in Figure 23.4. We'll see where these counters are incremented as we proceed through the code.

udpstat member	Description	
udps_badlen	#received datagrams with data length larger than packet	•
udps_badsum	#received datagrams with checksum error	•
udps_fullsock	#received datagrams not delivered because input socket full	
udps_hdrops	#received datagrams with packet shorter than header	•
udps_ipackets	total #received datagrams	•
udps_noport	#received datagrams with no process on destination port	•
udps_noportbcast	#received broadcast/multicast datagrams with no process on dest. port	•
udps_opackets	total #output datagrams	•
udpps_pcbcachemiss	#received input datagrams missing pcb cache	

Figure 23.4 UDP statistics maintained in the udpstat structure.

Figure 23.5 shows some sample output of these statistics, from the  ${\tt netstat}$  -s command.

netstat -s output	udpstat member
18,575,142 datagrams received	udps_ipackets
0 with incomplete header	udps_hdrops
18 with bad data length field	udps_badlen
58 with bad checksum	udps_badsum
84,079 dropped due to no socket	udps_noport
446 broadcast/multicast datagrams dropped due to no socket	udps_noportbcast
5,356 dropped due to full socket buffers	udps_fullsock
18,485,185 delivered	(see text)
18,676,277 datagrams output	udps_opackets

Figure 23.5 Sample UDP statistics.

The number of UDP datagrams delivered (the second from last line of output) is the number of datagrams received (udps\_ipackets) minus the six variables that precede it in Figure 23.5.

# **SNMP** Variables

Figure 23.6 shows the four simple SNMP variables in the UDP group and which counters from the udpstat structure implement that variable.

Figure 23.7 shows the UDP listener table, named udpTable. The values returned by SNMP for this table are taken from a UDP PCB, not the udpstat structure.

1.1

# 758 UDP: User Datagram Protocol

Chapter 23

Sectior

23.4

「「「「「「」」」

SNMP variable	udpstat member	Description
udpInDatagrams	udps_ipackets	#received datagrams delivered to processes
udpInErrors	udps_hdrops + udps_badsum + udps_badlen	#undeliverable UDP datagrams for reasons other than no application at destination port (e.g., UDP checksum error)
udpNoPorts	udps_noport + udps_noportbcast	#received datagrams for which no application process was at the destination port
udpOutDatagrams	udps_opackets	#datagrams sent

Figure 23.6 Simple SNMP variables in udp group.

UDP listener	table, index = <	udpLocalAddress >.< udpLocalPort >		
SNMP variable PCB variable Description				
udpLocalAddress	inp_laddr	local IP address for this listener		
udpLocalPort	inp_lport	local port number for this listener		

Figure 23.7 Variables in UDP listener table: udpTable.

# 23.3 UDP protosw Structure

Figure 23.8 lists the protocol switch entry for UDP.

Member	inetsw[1]	Description
<pre>pr_type pr_domain pr_protocol pr_flags pr_input pr_output pr_ctlinput pr_ctloutput pr_usrreq pr_init pr_fasttimo pr_slowtimo</pre>	SOCK_DGRAM &inetdomain IPPROTO_UDP (17) PR_ATOMIC PR_ADDR udp_input 0 udp_ctlinput ip_ctlinput ip_ctloutput udp_usrreq udp_init 0 0	Description         UDP provides datagram packet services         UDP is part of the Internet domain         appears in the ip_p field of the IP header         socket layer flags, not used by protocol processing         receives messages from IP layer         not used by UDP         control input function for ICMP errors         respond to administrative requests from a process         respond to communication requests from a process         initialization for UDP         not used by UDP
pr_drain pr_sysctl	0 _udp_sysct1	not used by UDP for sysct1(8) system call

Figure 23.8 The UDP protosw structure.

We describe the five functions that begin with udp\_ in this chapter. We also cover a sixth function, udp\_output, which is not in the protocol switch entry but is called by udp\_usrreq when a UDP datagram is output.

apter 23

(error)

as at the

Section 23.4

UDP Header 759

# 23.4 UDP Header

The UDP header is defined as a udphdr structure. Figure 23.9 shows the C structure and Figure 23.10 shows a picture of the UDP header.

```
udp.h
39 struct udphdr {
                                       /* source port */
       u_short uh_sport;
40
                                       /* destination port */
       u_short uh_dport;
41
                                       /* udp length */
       short uh_ulen;
42
                                       /* udp checksum */
43
       u_short uh_sum;
44 };
                                                                                    udp.h
                               Figure 23.9 udphdr structure.
                                                                                31
                                       15 16
  0
                                                         uh_dport
                 uh sport
                                                 16-bit destination port number
           16-bit source port number
                                                                                  8 bytes
```

uh\_sum

16-bit UDP checksum

data (if any)

uh\_ulen

16-bit UDP length



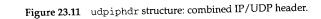
In the source code the UDP header is normally referenced as an IP header immediately followed by a UDP header. This is how udp\_input processes received IP datagrams, and how udp\_output builds outgoing IP datagrams. This combined IP/UDP header is a udpiphdr structure, shown in Figure 23.11.

– udp\_var.h 38 struct udpiphdr { /\* overlaid ip structure \*/ 39 struct ipovly ui\_i; /\* udp header \*/ struct udphdr ui\_u; 40 41 }; ui\_i.ih\_next 42 #define ui\_next ui\_i.ih\_prev 43 #define ui\_prev ui\_i.ih\_x1 44 #define ui\_x1 ui\_i.ih\_pr 45 #define ui\_pr ui\_i.ih\_len 46 #define ui\_len 47 #define ui\_src ui\_i.ih\_src ui\_i.ih\_dst 48 #define ui\_dst ui\_u.uh\_sport 49 #define ui\_sport

over a led by

3

зs



ui\_u.uh\_dport

ui\_u.uh\_ulen

ui\_u.uh\_sum

50 #define ui\_dport

51 #define ui\_ulen

52 #define ui\_sum

– udp\_var.h

Chapter 23

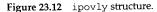
ip\_var.h

udp usrreq.c

Sectio

The 20-byte IP header is defined as an ipovly structure, shown in Figure 23.12.

			· · · · · · · · · · · · · · · · · · ·
38 struct ipovly {			
39 caddr_t ih_ne	ext, ih_prev;	/* for protocol se	quence q's */
40 u_char ih_x	1;	/* (unused) */	
41 u_char ih_p:	r;	/* protocol */	
42 short ih_l	en;	/* protocol length	*/
43 struct in_ade	dr ih_src;	/* source internet	
44 struct in_ad	dr ih_dst;	/* destination int	ernet address */
45 };			ip var.h
			ip_our.it



Unfortunately this structure is not a real IP header, as shown in Figure 8.8. The size is the same (20 bytes) but the fields are different. We'll return to this discrepancy when we discuss the calculation of the UDP checksum in Section 23.6.

# 23.5 udp\_init Function

The domaininit function calls UDP's initialization function (udp\_init, Figure 23.13) at system initialization time.

```
50 void
51 udp_init()
52 {
53     udb.inp_next = udb.inp_prev = &udb;
54 }
_______udp_usrreq.c
```

Figure 23.13 udp\_init function.

The only action performed by this function is to set the next and previous pointers in the head PCB (udb) to point to itself. This is an empty doubly linked list.

The remainder of the udb PCB is initialized to 0, although the only other field used in this head PCB is inp\_lport, the next UDP ephemeral port number to allocate. In the solution for Exercise 22.4 we mention that because this local port number is initialized to 0, the first ephemeral port number will be 1024.

# 23.6 udp\_output Function

UDP output occurs when the application calls one of the five write functions: send, sendto, sendmsg, write, or writev. If the socket is connected, any of the five functions can be called, although a destination address cannot be specified with sendto or sendmsg. If the socket is unconnected, only sendto and sendmsg can be called, and a



'3

h

h

.S

e

;)

.C

.C

S

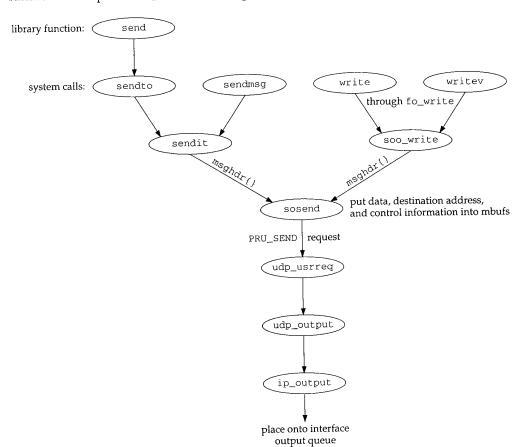
f

n

r

a

傸



destination address must be specified. Figure 23.14 summarizes how these five write functions end up with udp\_output being called, which in turn calls ip\_output.

Figure 23.14 How the five write functions end up calling udp\_output.

All five functions end up calling sosend, passing a pointer to a msghdr structure as an argument. The data to output is packaged into an mbuf chain and an optional destination address and optional control information are also put into mbufs by sosend. A PRU\_SEND request is issued.

UDP calls the function udp\_output, which we show the first half of in Figure 23.15. The four arguments are inp, a pointer to the socket Internet PCB; m, a pointer to the mbuf chain for output; addr, an optional pointer to an mbuf with the destination address packaged as a sockaddr\_in structure; and control, an optional pointer to an mbuf with control information from sendmsg.

Chapter 23

「「「「「「「「「「「「」」」」」

```
udp_usrreq.c
333 int
334 udp_output(inp, m, addr, control)
335 struct inpcb *inp;
336 struct mbuf *m;
337 struct mbuf *addr, *control;
338 {
339
        struct udpiphdr *ui;
340
        int
                len = m->m_pkthdr.len;
        struct in_addr laddr;
341
342
              s, error = 0;
        int
343
        if (control)
344
            m_freem(control);
                                     /* XXX */
345
        if (addr) {
346
            laddr = inp->inp_laddr;
347
            if (inp->inp_faddr.s_addr != INADDR_ANY) {
348
                error = EISCONN;
349
                goto release;
350
            }
351
            /*
             * Must block input while temporarily connected.
352
353
             */
354
            s = splnet();
355
            error = in_pcbconnect(inp, addr);
356
            if (error) {
357
                splx(s);
358
                goto release;
359
            }
360
        } else {
361
            if (inp->inp_faddr.s_addr == INADDR_ANY) {
362
                error = ENOTCONN;
363
                goto release;
364
            }
365
        }
        /*
366
367
         \ast Calculate data length and get an mbuf for UDP and IP headers.
368
         */
369
        M_PREPEND(m, sizeof(struct udpiphdr), M_DONTWAIT);
370
        if (m == 0) {
371
            error = ENOBUFS;
372
            goto release;
373
        }
                                                    andra Albertan
Maria
                    /* remainder of function shown in Figure 23.20 */
         \hat{\mathcal{F}}_{q} \notin \mathcal{A}_{q}
409
      release:
410
        m_freem(m);
411
        return (error);
412 }
                                                                         -udp_usrreq.c
```

Figure 23.15 udp\_output function: temporarily connect an unconnected socket.

**INTEL Ex.1013.788** 

# Discard optional control information

Any optional control information is discarded by m\_freem, without generating an error. UDP output does not use control information for any purpose.

The comment XXX is because the control information is ignored without generating an error. Other protocols, such as the routing domain and TCP, generate an error if the process passes control information.

### Temporarily connect an unconnected socket

345--359

10.00

If the caller specifies a destination address for the UDP datagram (addr is nonnull), the socket is temporarily connected to that destination address by in\_pcbconnect. The socket will be disconnected at the end of this function. Before doing this connect, a check is made as to whether the socket is already connected, and, if so, the error EISCONN is returned. This is why a sendto that specifies a destination address on a connected socket returns an error.

Before the socket is temporarily connected, IP input processing is stopped by splnet. This is done because the temporary connect changes the foreign address, foreign port, and possibly the local address in the socket's PCB. If a received UDP datagram were processed while this PCB was temporarily connected, that datagram could be delivered to the wrong process. Setting the processor priority to splnet only stops a software interrupt from causing the IP input routine to be executed (Figure 1.12), it does not prevent the interface layer from accepting incoming packets and placing them onto IP's input queue.

[Partridge and Pink 1993] note that this operation of temporarily connecting the socket is expensive and consumes nearly one-third of the cost of each UDP transmission.

The local address from the PCB is saved in laddr before temporarily connecting, because if it is the wildcard address it will be changed by in\_pcbconnect when it calls in pcbbind.

The same rules apply to the destination address that would apply if the process called connect, since in\_pcbconnect is called for both cases.

<sup>360–364</sup> If the process doesn't specify a destination address, and the socket is not connected, ENOTCONN is returned.

### Prepend IP and UDP headers

<sup>366–373</sup> M\_PREPEND allocates room for the IP and UDP headers in front of the data. Figure 1.8 showed one scenario, assuming there is not room in the first mbuf on the chain for the 28 bytes of header. Exercise 23.1 details the other possible scenarios. The flag M\_DONTWAIT is specified because if the socket is temporarily connected, IP processing is blocked, and M\_PREPEND should not block.

Earlier Berkeley releases incorrectly specified M\_WAIT here.

#### Prepending IP/UDP Headers and Mbuf Clusters

There is a subtle interaction between the M\_PREPEND macro and mbuf clusters. If the user data is placed into a cluster by sosend, then 56 bytes (max\_hdr from Figure 7.17)

q.c

are left unused at the beginning of the cluster, allowing room for the Ethernet, IP, and UDP headers. This is to prevent M\_PREPEND from allocating another mbuf just to hold these headers. M\_PREPEND calls M\_LEADINGSPACE to calculate how much space is available at the beginning of the mbuf:

```
#define M_LEADINGSPACE(m) \
    ((m)->m_flags & M_EXT ? /* (m)->m_data - (m)->m_ext.ext_buf */ 0 : \
    (m)->m_flags & M_PKTHDR ? (m)->m_data - (m)->m_pktdat : \
    (m)->m_data - (m)->m_dat)
```

The code that correctly calculates the amount of room at the front of a cluster is commented out, and the macro always returns 0 if the data is in a cluster. This means that when the user data is in a cluster, M\_PREPEND always allocates a new mbuf for the protocol headers instead of using the room allocated for this purpose by sosend.

The reason for commenting out the correct code in M\_LEADINGSPACE is that the cluster might be shared (Section 2.9), and, if it is shared, using the space before the user's data in the cluster could wipe out someone else's data.

With UDP data, clusters are not shared, since udp\_output does not save a copy of the data. TCP, however, saves a copy of the data in its send buffer (waiting for the data to be acknowledged), and if the data is in a cluster, it is shared. But tcp\_output doesn't call M\_LEADINGSPACE, because sosend leaves room for only 56 bytes at the beginning of the cluster for datagram protocols. tcp\_output always calls MGETHDR instead, to allocate an mbuf for the protocol headers.

# UDP Checksum Calculation and Pseudo-Header

Before showing the last half of udp\_output we describe how UDP fills in some of the fields in the IP/UDP headers, calculates the UDP checksum, and passes the IP/UDP headers and the data to IP for output. The way this is done with the ipovly structure is tricky.

Figure 23.16 shows the 28-byte IP/UDP headers that are built by udp\_output in the first mbuf in the chain pointed to by m. The unshaded fields are filled in by udp\_output and the shaded fields are filled in by ip\_output. This figure shows the format of the headers as they appear on the wire.

The UDP checksum is calculated over three areas: (1) a 12-byte pseudo-header containing fields from the IP header, (2) the 8-byte UDP header, and (3) the UDP data. Figure 23.17 shows the 12 bytes of pseudo-header used for the checksum computation, along with the UDP header. The UDP header used for the checksum calculation is identical to the UDP header that appears on the wire (Figure 23.16).

The following three facts are used in computing the UDP checksum. (1) The third 32-bit word in the pseudo-header (Figure 23.17) looks similar to the third 32-bit word in the IP header (Figure 23.16): two 8-bit values and a 16-bit value. (2) The order of the three 32-bit values in the pseudo-header is irrelevant. Actually, the computation of the Internet checksum does not depend on the order of the 16-bit values that are used (Section 8.7). (3) Including additional 32-bit words of 0 in the checksum computation has no effect.

100 A 100 A

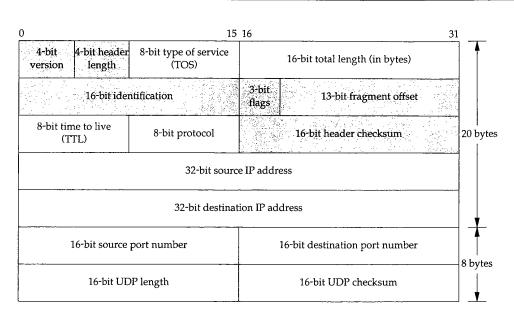


Figure 23.16 IP/UDP headers: unshaded fields filled in by UDP; shaded fields filled in by IP.

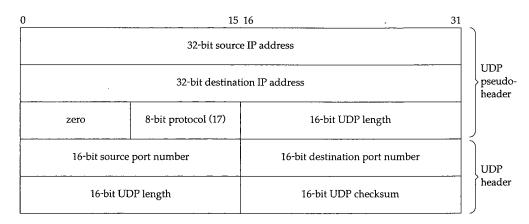


Figure 23.17 Pseudo-header used for checksum computation and UDP header.

udp\_output takes advantage of these three facts and fills in the fields in the udpiphdr structure (Figure 23.11), which we depict in Figure 23.18. This structure is contained in the first mbuf in the chain pointed to by the argument m.

The last three 32-bit words in the 20-byte IP header (the five members ui\_x1, ui\_pr, ui\_len, ui\_src, and ui\_dst) are used as the pseudo-header for the checksum computation. The first two 32-bit words in the IP header (ui\_next and ui\_prev) are also used in the checksum computation, but they're initialized to 0, and don't affect the checksum. 766 UDP: User Datagram Protocol Chapter 23

Ę,

цź

いましたが変要

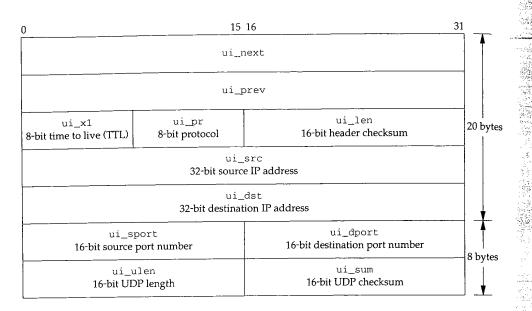


Figure 23.18 udpiphdr structure used by udp\_output.

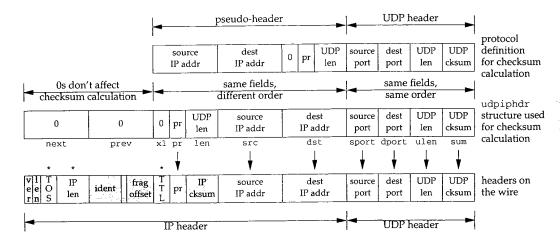


Figure 23.19 Operations to fill in IP/UDP headers and calculate UDP checksum.

Figure 23.19 summarizes the operations we've described.

1. The top picture shown in Figure 23.19 is the protocol definition of the pseudoheader, which corresponds to Figure 23.17.

er 23:

vytes

tes

un

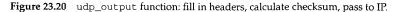
sed ım

10-

- 2. The middle picture is the udpiphdr structure that is used in the source code, which corresponds to Figure 23.11. (To make the figure readable, the prefix ui\_ has been left off all the members.) This is the structure built by udp\_output in the first mbuf and then used to calculate the UDP checksum.
- 3. The bottom picture shows the IP/UDP headers that appear on the wire, which corresponds to Figure 23.16. The seven fields with an arrow above are filled in by udp\_output before the checksum computation. The three fields with an asterisk above are filled in by udp\_output after the checksum computation. The remaining six shaded fields are filled in by ip\_output.

Figure 23.20 shows the last half of the udp\_output function.

374	/* udp_usrreq.c
375	* Fill in mbuf with extended UDP header
376	* and addresses and length put into network format.
377	*/
378	ui = mtod(m, struct udpiphdr *);
379	ui->ui_next = ui->ui_prev = 0;
380	ui->ui_x1 = 0;
381	ui->ui_pr = IPPROTO_UDP;
382	ui->ui_len = htons((u_short) len + sizeof(struct udphdr));
383	ui->ui_src = inp->inp_laddr;
384	ui->ui_dst = inp->inp_faddr;
385	ui->ui_sport = inp->inp_lport;
386	ui->ui_dport = inp->inp_fport;
387	ui->ui_ulen = ui->ui_len;
388	/*
389	* Stuff checksum and output datagram.
390	*/
391	ui->ui_sum = 0;
392	if (udpcksum) {
393	<pre>if ((ui-&gt;ui_sum = in_cksum(m, sizeof(struct udpiphdr) + len)) == 0)</pre>
394	ui->ui_sum = 0xffff;
395	)
396	((struct ip *) ui)->ip_len = sizeof(struct udpiphdr) + len;
397	((struct ip *) ui)->ip_ttl = inp->inp_ip.ip_ttl; /* XXX */
398	((struct ip *) ui)->ip_tos = inp->inp_ip.ip_tos; /* XXX */
399	udpstat.udps_opackets++;
400	error = ip_output(m, inp->inp_options, &inp->inp_route,
401	inp->inp_socket->so_options & (SO_DONTROUTE   SO_BROADCAST),
402	<pre>inp-&gt;inp_moptions);</pre>
403	if (addr) {
404	<pre>in_pcbdisconnect(inp);</pre>
405	<pre>inp-&gt;inp_laddr = laddr;</pre>
406	<pre>splx(s);</pre>
407	}
408	return (error);udp usrreq.c



#### Prepare pseudo-header for checksum computation

All the members in the udpiphdr structure (Figure 23.18) are set to their respective values. The local and foreign sockets from the PCB are already in network byte order, but the UDP length must be converted to network byte order. The UDP length is the number of bytes of data (len, which can be 0) plus the size of the UDP header (8). The UDP length field appears twice in the UDP checksum calculation: ui\_len and ui\_ulen. One of them is redundant.

#### Calculate checksum

The checksum is calculated by first setting it to 0 and then calling in\_cksum. If UDP checksums are disabled (a bad idea—see Section 11.3 of Volume 1), 0 is sent as the checksum. If the calculated checksum is 0, 16 one bits are stored in the header instead of 0. (In one's complement arithmetic, all one bits and all zero bits are both considered 0.) This allows the receiver to distinguish between a UDP packet without a checksum (the checksum field is 0) versus a UDP packet with a checksum whose value is 0 (the checksum is 16 one bits).

The variable udpcksum (Figure 23.3) normally defaults to 1, enabling UDP checksums. The kernel can be compiled for 4.2BSD compatibility, which initializes udpcksum to 0.

#### Fill in UDP length, TTL, and TOS

The pointer ui is cast to a pointer to a standard IP header (ip), and three fields in the IP header are set by UDP. The IP length field is set to the amount of data in the UDP datagram, plus 28, the size of the IP/UDP headers. Notice that this field in the IP header is stored in host byte order, not network byte order like the rest of the multibyte fields in the header. ip\_output converts it to network byte order before transmission.

The TTL and TOS fields in the IP header are then set from the values in the socket's PCB. These values are defaulted by UDP when the socket is created, but can be changed by the process using setsockopt. Since these three fields—IP length, TTL, and TOS—are not part of the pseudo-header and not used in the UDP checksum computation, they must be set after the checksum is calculated but before ip\_output is called.

#### Send datagram

400-402 ip\_output sends the datagram. The second argument, inp\_options, are IP options the process can set using setsockopt. These IP options are placed into the IP header by ip\_output. The third argument is a pointer to the cached route in the PCB, and the fourth argument is the socket options. The only socket options that are passed to ip\_output are SO\_DONTROUTE (bypass the routing tables) and SO\_BROADCAST (allow broadcasting). The final argument is a pointer to the multicast options for this socket.

#### Disconnect temporarily connected socket

403-407

If the socket was temporarily connected, in\_pcbdisconnect disconnects the socket, the local IP address is restored in the PCB, and the interrupt level is restored to its saved value.

ſ

udp\_input Function 769

Section 23.7

# 23.7 udp\_input Function

UDP output is driven by a process calling one of the five write functions. The functions shown in Figure 23.14 are all called directly as part of the system call. UDP input, on the other hand, occurs when IP input receives an IP datagram on its input queue whose protocol field specifies UDP. IP calls the function udp\_input through the pr\_input function in the protocol switch table (Figure 8.15). Since IP input is at the software interrupt level, udp\_input also executes at this level. The goal of udp\_input is to place the UDP datagram onto the appropriate socket's buffer and wake up any process blocked for input on that socket.

We'll divide our discussion of the udp\_input function into three sections:

- 1. the general validation that UDP performs on the received datagram,
- 2. processing UDP datagrams destined for a unicast address: locating the appropriate PCB and placing the datagram onto the socket's buffer, and
- 3. processing UDP datagrams destined for a broadcast or multicast address: the datagram may be delivered to multiple sockets.

This last step is new with the support of multicasting in Net/3, but consumes almost one-third of the code.

#### General Validation of Received UDP Datagram

Figure 23.21 shows the first section of UDP input.

The two arguments to udp\_input are m, a pointer to an mbuf chain containing the IP datagram, and iphlen, the length of the IP header (including possible IP options).

### Discard IP options

- 67-76 If IP options are present they are discarded by ip\_stripoptions. As the comments indicate, UDP should save a copy of the IP options and make them available to the receiving process through the IP\_RECVOPTS socket option, but this isn't implemented yet.
- <sup>77–88</sup> If the length of the first mbuf on the mbuf chain is less than 28 bytes (the size of the IP header plus the UDP header), m\_pullup rearranges the mbuf chain so that at least 28 bytes are stored contiguously in the first mbuf.

770 UDP: User Datagram Protocol

Chapter 23

Se

```
udp_usrreq.c
55 void
56 udp_input(m, iphlen)
57 struct mbuf *m;
58 int
            iphlen;
59 (
60
        struct ip *ip;
61
        struct udphdr *uh;
        struct inpcb *inp;
62
        struct mbuf *opts = 0;
63
64
        int
                len;
65
        struct ip save_ip;
66
        udpstat.udps_ipackets++;
67
        * Strip IP options, if any; should skip this,
68
         * make available to user, and use on returned packets,
69
70
         * but we don't yet have a way to check the checksum
71
         * with options still present.
72
         */
73
        if (iphlen > sizeof(struct ip)) {
74
            ip_stripoptions(m, (struct mbuf *) 0);
75
            iphlen = sizeof(struct ip);
76
        }
77
        /*
78
         *
          Get IP and UDP header together in first mbuf.
         */
79
80
        ip = mtod(m, struct ip *);
        if (m->m_len < iphlen + sizeof(struct udphdr)) {</pre>
81
            if ((m = m_pullup(m, iphlen + sizeof(struct udphdr))) == 0) {
82
83
                udpstat.udps_hdrops++;
84
                return;
85
            }
86
            ip = mtod(m, struct ip *);
87
        }
88
        uh = (struct udphdr *) ((caddr_t) ip + iphlen);
89
        /*
         * Make mbuf data length reflect UDP length.
90
91
         * If not enough data to reflect UDP length, drop.
 92
         */
93
        len = ntohs((u_short) uh->uh_ulen);
94
        if (ip->ip_len != len) {
 95
            if (len > ip->ip_len) {
 96
                udpstat.udps_badlen++;
 97
                goto bad;
 98
 99
            m_adj(m, len - ip->ip_len);
100
            /* ip->ip_len = len; */
101
        }
102
        /*
         * Save a copy of the IP header in case we want to restore
103
104
         * it for sending an ICMP error message in response.
105
         */
106
        save_ip = *ip;
```

107 108 \* Checksum extended UDP header and data. 109 \*/ 110 if (udpcksum && uh->uh\_sum) { 111 ((struct ipovly \*) ip)->ih\_next = 0; 112 ((struct ipovly \*) ip)->ih\_prev = 0; 113 ((struct ipovly \*) ip)->ih\_x1 = 0; 114 ((struct ipovly \*) ip)->ih\_len = uh->uh\_ulen; if (uh->uh\_sum = in\_cksum(m, len + sizeof(struct ip))) { 115 116 udpstat.udps\_badsum++; 117 m\_freem(m); 118 return: 119 } 120 3

- udp\_usrreq.c

# Figure 23.21 udp\_input function: general validation of received UDP datagram.

#### Verify UDP length

89-101

「日本の一方法である」

53

There are two lengths associated with a UDP datagram: the length field in the IP header (ip\_len) and the length field in the UDP header (uh\_ulen). Recall that ipintr subtracted the length of the IP header from ip\_len before calling udp\_input (Figure 10.11). The two lengths are compared and there are three possibilities:

- 1. ip\_len equals uh\_ulen. This is the common case.
- 2. ip\_len is greater than uh\_ulen. The IP datagram is too big, as shown in Figure 23.22.

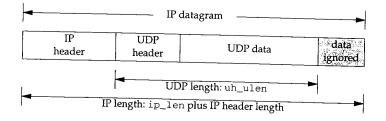


Figure 23.22 UDP length too small.

The code believes the smaller of the two lengths (the UDP header length) and  $m\_adj$  removes the excess bytes of data from the end of the datagram. In the code the second argument to  $m\_adj$  is negative, which we said in Figure 2.20 trims data from the end of the mbuf chain. It is possible in this scenario that the UDP length field has been corrupted. If so, the datagram will probably be discarded shortly, assuming the sender calculated the UDP checksum, that this checksum detects the error, and that the receiver verifies the checksum. The IP length field should be correct since it was verified by IP against the amount of data received from the interface, and the IP length field is covered by the mandatory IP header checksum.

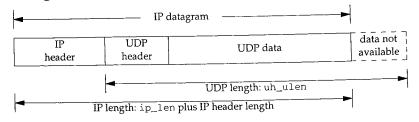
].C

23

and and the first first the

772 UDP: User Datagram Protocol

3. ip\_len is less than uh\_ulen. The IP datagram is smaller than possible, given the length in the UDP header. Figure 23.23 shows this case.





Something is wrong and the datagram is discarded. There is no other choice here: if the UDP length field has been corrupted, it can't be detected with the UDP checksum. The correct UDP length is needed to calculate the checksum.

As we've said, the UDP length is redundant. In Chapter 28 we'll see that TCP does not have a length field in its header—it uses the IP length field, minus the lengths of the IP and TCP headers, to determine the amount of data in the datagram. Why does the UDP length field exist? Possibly to add a small amount of error checking, since UDP checksums are optional.

# Save copy of IP header and verify UDP checksum udp\_input saves a copy of the IP header before verifying the checksum, because

102-106

110

the checksum computation wipes out some of the fields in the original IP header. The checksum is verified only if UDP checksums are enabled for the kernel (udpcksum), and if the sender calculated a UDP checksum (the received checksum is nonzero).

This test is incorrect. If the sender calculated a checksum, it should be verified, regardless of whether outgoing checksums are calculated or not. The variable udpcksum should only specify whether outgoing checksums are calculated. Unfortunately many vendors have copied this incorrect test, although many vendors today finally ship their kernels with UDP checksums enabled by default.

<sup>111–120</sup> Before calculating the checksum, the IP header is referenced as an ipovly structure (Figure 23.18) and the fields are initialized as described in the previous section when the UDP checksum is calculated by udp\_output.

At this point special code is executed if the datagram is destined for a broadcast or multicast IP address. We defer this code until later in the section.

# **Demultiplexing Unicast Datagrams**

Assuming the datagram is destined for a unicast address, Figure 23.24 shows the code that is executed.

206-

210-

Chapter 23

given

pter 23

	그는 것 그는 한 것같이 하는 아님께서 것을 못했다. 나는 것 그는 아님은 것 같은 것을 했다.
206	/*
207	* Locate pcb for unicast datagram.
208	*/
209	inp = udp_last_inpcb;
210	if (inp->inp_lport != uh->uh_dport
211	inp->inp fport != uh->uh_sport
212	inp->inp_faddr.s_addr != ip->ip_src.s_addr
213	<pre>inp-&gt;inp_laddr.s_addr != ip-&gt;ip_dst.s_addr) {</pre>
214	<pre>inp = in_pcblookup(&amp;udb, ip-&gt;ip_src, uh-&gt;uh_sport,</pre>
215	ip->ip_dst, uh->uh_dport, INPLOOKUP_WILDCARD);
216	if (inp)
217	udp_last_inpcb = inp;
218	udpstat.udpps_pcbcachemiss++;
219	}
220	if $(inp == 0)$ {
221	udpstat.udps_noport++;
222	if (m->m_flags & (M_BCAST   M_MCAST)) {
223	udpstat.udps_noportbcast++;
224	goto bad;
225	} .
226	<pre>*ip = save_ip;</pre>
227	ip->ip_len += iphlen;
228	<pre>icmp_error(m, ICMP_UNREACH, ICMP_UNREACH_PORT, 0, 0);</pre>
229	return;
230	}udp_usrrea

Figure 23.24 udp\_input function: demultiplex unicast datagram.

#### Check one-behind cache

UDP maintains a pointer to the last Internet PCB for which it received a datagram, udp\_last\_inpcb. Before calling in\_pcblookup, which might have to search many PCBs on the UDP list, the foreign and local addresses and ports of that last PCB are compared against the received datagram. This is called a *one-behind cache* [Partridge and Pink 1993], and it is based on the assumption that the next datagram received has a high probability of being destined for the same socket as the last received datagram [Mogul 1991]. This cache was introduced with the 4.3BSD Tahoe release.

<sup>210–213</sup> The order of the four comparisons between the cached PCB and the received datagram is intentional. If the PCBs don't match, the comparisons should stop as soon as possible. The highest probability is that the destination port numbers are different—this is therefore the first test. The lowest probability of a mismatch is between the local addresses, especially on a host with just one interface, so this is the last test.

th the 1m. t have a

choice

語を設定した。ため

دھ . ا

×.

nd TCP ;th field onal.

ecause

kernel sum is

dless of dy specbied this bicksums

ucture ien the

cast or

e code

Chapter 23

Unfortunately this one-behind cache, as coded, is practically useless [Partridge and Pink 1993]. The most common type of UDP server binds only its well-known port, leaving its local address, foreign address, and foreign port wildcarded. The most common type of UDP client does not connect its UDP socket; it specifies the destination address for each datagram using sendto. Therefore most of the time the three values in the PCB inp\_laddr, inp\_faddr, and inp\_fport are wildcards. In the cache comparison the four values in the received datagram are never wildcards, meaning the cache entry will compare equal with the received datagram only when the PCB has all four local and foreign values specified to nonwildcard values. This happens only for a connected UDP socket.

On the system bsdi, the counter udpps\_pcbcachemiss was 41,253 and the counter udps\_ipackets was 42,485. This is less than a 3% cache hit rate.

The netstat -s command prints most of the fields in the udpstat structure (Figure 23.5). Unfortunately the Net/3 version, and most vendor's versions, never print udpps\_pcbcachemiss. If you want to see the value, use a debugger to examine the variable in the running kernel.

#### Search all UDP PCBs

Assuming the comparison with the cached PCB fails, in\_pcblookup searches for a match. The INPLOOKUP\_WILDCARD flag is specified, allowing a wildcard match. If a match is found, the pointer to the PCB is saved in udp\_last\_inpcb, which we said is a cache of the last received UDP datagram's PCB.

# Generate ICMP port unreachable error

If a matching PCB is not found, UDP normally generates an ICMP port unreachable error. First the m\_flags for the received mbuf chain is checked to see if the datagram was sent to a link-level broadcast or multicast destination address. It is possible to receive an IP datagram with a unicast IP address that was sent to a broadcast or multicast link-level address, but an ICMP port unreachable error must not be generated. If it is OK to generate the ICMP error, the IP header is restored to its received value (save\_ip) and the IP length is also set back to its original value.

This check for a link-level broadcast or multicast address is redundant. icmp\_error also performs this check. The only advantage in this redundant check is to maintain the counter udps\_noportbcast in addition to the counter udps\_noport.

The addition of iphlen back into ip\_len is a bug. icmp\_error will also do this, causing the IP length field in the IP header returned in the ICMP error to be 20 bytes too large. You can tell if a system has this bug by adding a few lines of code to the Traceroute program (Chapter 8 of Volume 1) to print this field in the ICMP port unreachable that is returned when the destination host is finally reached.

Figure 23.25 is the next section of processing for a unicast datagram, delivering the datagram to the socket corresponding to the destination PCB.

23

d

1-

'n

s

ιe

iıe

ır

۱-

er

i).

at

le

а

а

is

e

n

0

i-

it

e

r-

æ

g

n

8 1-

e

udp\_usrreq.c 231 /\* 232 \* Construct sockaddr format source address. \* Stuff source address and datagram in user buffer. 233 234 \*/ 235 udp\_in.sin\_port = uh->uh\_sport; udp\_in.sin\_addr = ip->ip\_src; 236 237 if (inp->inp\_flags & INP\_CONTROLOPTS) { struct mbuf \*\*mp = &opts; 238 239 if (inp->inp\_flags & INP\_RECVDSTADDR) ( \*mp = udp\_saveopt((caddr\_t) & ip->ip\_dst, 240 241 sizeof(struct in\_addr), IP\_RECVDSTADDR); 242 if (\*mp) 243  $mp = \&(*mp) - >m_next;$ 244 } 245 #ifdef notyet 246 /\* IP options were tossed above \*/ 247 if (inp->inp\_flags & INP\_RECVOPTS) { 248 \*mp = udp\_saveopt((caddr\_t) opts\_deleted\_above, 249 sizeof(struct in\_addr), IP\_RECVOPTS); 250 if (\*mp) 251  $mp = \&(*mp) - m_next;$ 252 } /\* ip\_srcroute doesn't do what we want here, need to fix \*/ 253 if (inp->inp\_flags & INP\_RECVRETOPTS) { 254 \*mp = udp\_saveopt((caddr\_t) ip\_srcroute(), 255 256 sizeof(struct in\_addr), IP\_RECVRETOPTS); 257 if (\*mp) 258  $mp = \& (*mp) - >m_next;$ 259 } 260 #endif 261 } 262 iphlen += sizeof(struct udphdr); 263 m->m\_len -= iphlen; 264 m->m\_pkthdr.len -= iphlen; 265 m->m\_data += iphlen; if (sbappendaddr(&inp->inp\_socket->so\_rcv, (struct sockaddr \*) &udp\_in, 266 267  $m, opts) == 0) \{$ 268 udpstat.udps\_fullsock++; 269 goto bad; 270 } 271 sorwakeup(inp->inp\_socket); 272 return; 273 bad: 274 m\_freem(m); 275 if (opts) 276 m\_freem(opts); 277 } — udp\_usrreq.c

Figure 23.25 udp\_input function: deliver unicast datagram to socket.

#### Return source IP address and source port

231-236

The source IP address and source port number from the received IP datagram are stored in the global sockaddr\_in structure udp\_in. This structure is passed as an argument to sbappendaddr later in the function.

Using a global to hold the IP address and port number is OK because udp\_input is single threaded. When this function is called by ipintr it processes the received data-gram completely before returning. Also, sbappendaddr copies the socket address structure from the global into an mbuf.

#### IP\_RECVDSTADDR socket option

237--244

The constant INP\_CONTROLOPTS is the combination of the three socket options that the process can set to cause control information to be returned through the recvmsg system call for a UDP socket (Figure 22.5). The IP\_RECVDSTADDR socket option returns the destination IP address from the received UDP datagram as control information. The function udp\_saveopt allocates an mbuf of type MT\_CONTROL and stores the 4-byte destination IP address in the mbuf. We show this function in Section 23.8.

This socket option appeared with 4.3BSD Reno and was intended for applications such as TFTP, the Trivial File Transfer Protocol, that should not respond to client requests that are sent to a broadcast address. Unfortunately, even if the receiving application uses this option, it is nontrivial to determine if the destination IP address is a broadcast address or not (Exercise 23.6).

When the multicasting changes were added in 4.4BSD, this code was left in only for datagrams destined for a unicast address. We'll see in Figure 23.26 that this option is not implemented for datagrams sent to a broadcast of multicast address. This defeats the purpose of the option!

### Unimplemented socket options

245-260

This code is commented out because it doesn't work. The intent of the IP\_RECVOPTS socket option is to return the IP options from the received datagram as control information, and the intent of IP\_RECVRETOPTS socket option is to return source route information. The manipulation of the mp variable by all three IP\_RECV socket options is to build a linked list of up to three mbufs that are then placed onto the socket's buffer by sbappendaddr. The code shown in Figure 23.25 only returns one option as control information, so the m\_next pointer of that mbuf is always a null pointer.

# Append data to socket's receive queue

262-272

At this point the received datagram (the mbuf chain pointed to by m), is ready to be placed onto the socket's receive queue along with a socket address structure representing the sender's IP address and port (udp\_in), and optional control information (the destination IP address, the mbuf pointed to by opts). This is done by sbappendaddr. Before calling this function, however, the pointer and lengths of the first mbuf on the chain are adjusted to ignore the IP and UDP headers. Before returning, sorwakeup is called for the receiving socket to wake up any processes asleep on the socket's receive queue.

17

でいたのに対応に

#### Error return

If an error is encountered during UDP input processing, udp\_input jumps to the label bad. The mbuf chain containing the datagram is released, along with the mbuf chain containing any control information (if present).

#### Demultiplexing Multicast and Broadcast Datagrams

We now return to the portion of udp\_input that handles datagrams sent to a broadcast or multicast IP address. The code is shown in Figure 23.26.

As the comments indicate, these datagrams are delivered to *all* sockets that match, not just a single socket. The inadequacy of the UDP interface that is mentioned refers to the inability of a process to receive asynchronous errors on a UDP socket (notably ICMP port unreachables) unless the socket is connected. We described this in Section 22.11.

<sup>139–145</sup> The source IP address and port number are saved in the global sockaddr\_in structure udp\_in, which is passed to sbappendaddr. The mbuf chain's length and data pointer are updated to ignore the IP and UDP headers.

<sup>146–164</sup> The large for loop scans each UDP PCB to find all matching PCBs. in\_pcblookup is not called for this demultiplexing because it returns only one PCB, whereas the broadcast or multicast datagram may be delivered to more than one PCB.

If the local port in the PCB doesn't match the destination port from the received datagram, the entry is ignored. If the local address in the PCB is not the wildcard, it is compared to the destination IP address and the entry is skipped if they're not equal. If the foreign address in the PCB is not a wildcard, it is compared to the source IP address and if they match, the foreign port must also match the source port. This last test assumes that if the socket is connected to a foreign IP address it must also be connected to a foreign port, and vice versa. This is the same logic we saw in in\_pcblookup.

165-177 If this is not the first match found (last is nonnull), a copy of the datagram is placed onto the receive queue for the previous match. Since sbappendaddr releases the mbuf chain when it is done, a copy is first made by m\_copy. Any processes waiting for this data are awakened by sorwakeup. A pointer to this matching socket structure is saved in last.

This use of the variable last avoids calling m\_copy (an expensive operation since an entire mbuf chain is copied) unless there are multiple recipients for a given datagram. In the common case of a single recipient, the for loop just sets last to the single matching PCB, and when the loop terminates, sbappendaddr places the mbuf chain onto the socket's receive queue—a copy is not made.

178-188 If this matching socket doesn't have either the SO\_REUSEPORT or the SO\_REUSEADDR socket option set, then there's no need to check for additional matches and the loop is terminated. The datagram is placed onto the single socket's receive queue in the call to sbappendaddr outside the loop.

189–197 If last is null at the end of the loop, no matches were found. An ICMP error is not generated because the datagram was sent to a broadcast or multicast IP address.

13

Chapter 23

Sec

121	if (IN_MULTICAST(ntohl(ip->ip_dst.s_addr))    udp_usrreq.c
122	in_broadcast(ip->ip_dst, m->m_pkthdr.rcvif)) {
123	struct socket *last;
124	/*
125	* Deliver a multicast or broadcast datagram to *all* sockets
126	* for which the local and remote addresses and ports match
127	* those of the incoming datagram. This allows more than
128	* one process to receive multi/broadcasts on the same port.
129	
130	* (This really ought to be done for unicast datagrams as
130	* well, but that would cause problems with existing
132	* applications that open both address-specific sockets and
133	* a wildcard socket listening to the same port they would
134	* end up receiving duplicates of every unicast datagram.
L34 L35	* Those applications open the multiple sockets to overcome an
136	<pre>* inadequacy of the UDP socket interface, but for backwards * answer in the socket interface.</pre>
L30 L37	* compatibility we avoid the problem here rather than
L37 L38	<pre>* fixing the interface. Maybe 4.5BSD will remedy this?) */</pre>
130	*/
.39	/*
40	* Construct sockaddr format source address.
.41	*/
42	udp_in.sin_port = uh->uh_sport;
43	udp_in.sin_addr = ip->ip_src;
44	<pre>m-&gt;m_len -= sizeof(struct udpiphdr);</pre>
45	m->m_data += sizeof(struct udpiphdr);
46	/*
47	* Locate pcb(s) for datagram.
48	* (Algorithm copied from raw_intr().)
49	*/
.50	last = NULL;
.51	for (inp = udb.inp_next; inp != &udb inp = inp->inp_next) {
.52	if (inp->inp_lport != uh->uh_dport)
.53	continue;
54	if (inp->inp_laddr.s_addr != INADDR_ANY) {
55	if (inp->inp_laddr.s_addr != INADDR_ANY) {
56	ip->ip_dst.s_addr)
57	continue;
58	}
59	; if (inp->inp_faddr.s_addr != INADDR_ANY) {
60	if (inp->inp_faddr.s_addr != INADDR_ANY) {
61	
62	ip->ip_src.s_addr
63	<pre>inp-sinp_fport != uh-suh_sport) continue;</pre>
64	<pre>continue; }</pre>
65	if (last != NULL) {
66	
-	<pre>struct mbuf *n;</pre>
67	if $((n = m_{copy}(m, 0, M_{COPYALL})) != NULL) {$
68	if (sbappendaddr(&last->so_rcv,
69	(struct sockaddr *) &udp_in,
70	n, (struct mbuf *) 0) == 0) {
71	<pre>m_freem(n);</pre>
72	udpstat.udps_fullsock++;

198~

Con

5

```
173
                         } else
                             sorwakeup(last);
174
175
                     }
176
                }
177
                last = inp->inp_socket;
178
                 /*
                 * Don't look for additional matches if this one does
179
180
                 * not have either the SO_REUSEPORT or SO_REUSEADDR
                 * socket options set. This heuristic avoids searching
181
                 * through all pcbs in the common case of a non-shared
182
                  * port. It assumes that an application will never
183
184
                  * clear these options after setting them.
                  * /
185
                if ((last->so_options & (SO_REUSEPORT | SO_REUSEADDR) == 0))
186
                     break;
187
188
            }
            if (last == NULL) {
189
                 /*
190
191
                 * No matching pcb found; discard datagram.
                  * (No need to send an ICMP Port Unreachable
192
                  * for a broadcast or multicast datgram.)
193
194
                  */
195
                 udpstat.udps_noportbcast++;
                 goto bad;
196
197
            }
            if (sbappendaddr(&last->so_rcv, (struct sockaddr *) &udp_in,
198
                             m, (struct mbuf *) 0) == 0) {
199
                 udpstat.udps_fullsock++;
200
201
                 goto bad;
202
             }
             sorwakeup(last);
203
204
             return;
205
         }
                                                                         udp_usrreq.c
```

Figure 23.26 udp\_input function: demultiplexing of broadcast and multicast datagrams.

<sup>198–204</sup> The final matching entry (which could be the only matching entry) has the original datagram (m) placed onto its receive queue. After sorwakeup is called, udp\_input returns, since the processing the broadcast or multicast datagram is complete.

The remainder of the function (shown previously in Figure 23.24) handles unicast datagrams.

#### **Connected UDP Sockets and Multihomed Hosts**

There is a subtle problem when using a connected UDP socket to exchange datagrams with a process on a multihomed host. Datagrams from the peer may arrive with a different source IP address and will not be delivered to the connected socket. Consider the example shown in Figure 23.27.



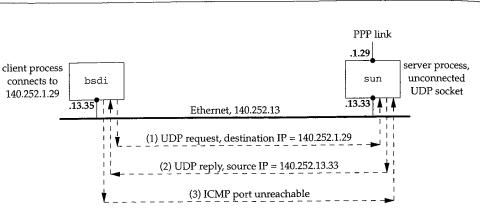


Figure 23.27 Example of connected UDP socket sending datagram to a multihomed host.

Three steps take place.

1. The client on bsdi creates a UDP socket and connects it to 140.252.1.29, the PPP interface on sun, not the Ethernet interface. A datagram is sent on the socket to the server.

The server on sun receives the datagram and accepts it, even though it arrives on an interface that differs from the destination IP address. (sun is acting as a router, so whether it implements the weak end system model or the strong end system model doesn't matter.) The datagram is delivered to the server, which is waiting for client requests on an unconnected UDP socket.

2. The server sends a reply, but since the reply is being sent on an unconnected UDP socket, the source IP address for the reply is chosen by the kernel based on the outgoing interface (140.252.13.33). The destination IP address in the request is not used as the source address for the reply.

When the reply is received by bsdi it is not delivered to the client's connected UDP socket since the IP addresses don't match.

3. bsdi generates an ICMP port unreachable error since the reply can't be demultiplexed. (This assumes that there is not another process on bsdi eligible to receive the datagram.)

The problem in this example is that the server does not use the destination IP address from the request as the source IP address of the reply. If it did, the problem wouldn't exist, but this solution is nontrivial—see Exercise 23.10. We'll see in Figure 28.16 that a TCP server uses the destination IP address from the client as the source IP address from the server, if the server has not explicitly bound a local IP address to its socket. Sect

23

271

29

Chapter 23

- udp\_usrreq.c

- udp\_usrreq.c

# 23.8 udp\_saveopt Function

If a process specifies the IP\_RECVDSTADDR socket option, to receive the destination IP address from the received datagram udp\_saveopt is called by udp\_input:

Figure 23.28 shows this function.

```
278 /*
    * Create a "control" mbuf containing the specified data
279
280 * with the specified type for presentation with a datagram.
281 */
282 struct mbuf *
283 udp_saveopt(p, size, type)
284 caddr_t p;
285 int
            size;
286 int
            type;
287 {
        struct cmsghdr *cp;
288
       struct mbuf *m;
289
        if ((m = m_get(M_DONTWAIT, MT_CONTROL)) == NULL)
290
           return ((struct mbuf *) NULL);
291
        cp = (struct cmsghdr *) mtod(m, struct cmsghdr *);
292
        bcopy(p, CMSG_DATA(cp), size);
293
        size += sizeof(*cp);
294
        m->m_len = size;
295
        cp->cmsg_len = size;
296
        cp->cmsg_level = IPPROTO_IP;
297
        cp->cmsg_type = type;
298
299
        return (m);
300 }
```

Figure 23.28 udp\_saveopt function: create mbuf with control information.

The arguments are p, a pointer to the information to be stored in the mbuf (the destination IP address from the received datagram); size, its size in bytes (4 in this example, the size of an IP address); and type, the type of control information (IP RECVDSTADDR).

An mbuf is allocated, and since the code is executing at the software interrupt layer, M\_DONTWAIT is specified. The pointer cp points to the data portion of the mbuf, and it is cast into a pointer to a cmsghdr structure (Figure 16.14). The IP address is copied from the IP header into the data portion of the cmsghdr structure by bcopy. The length of the mbuf is then set (to 16 in this example), followed by the remainder of the cmsghdr structure. Figure 23.29 shows the final state of the mbuf.

The cmsg\_len field contains the length of the cmsghdr structure (12) plus the size of the cmsg\_data field (4 for this example). If the application calls recvmsg to receive the control information, it must go through the cmsghdr structure to determine the type and length of the cmsg\_data field.

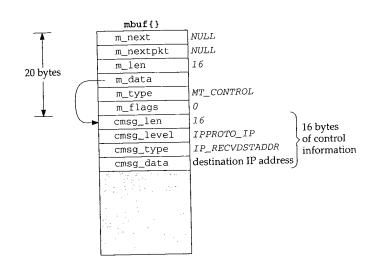


Figure 23.29 Mbuf containing destination address from received datagram as control information.

## 23.9 udp\_ctlinput Function

When icmp\_input receives an ICMP error (destination unreachable, parameter problem, redirect, source quench, and time exceeded) the corresponding protocol's pr\_ctlinput function is called:

For UDP, Figure 22.32 showed that the function udp\_ctlinput is called. We show this function in Figure 23.30.

314-322

The arguments are cmd, one of the PRC\_xxx constants from Figure 11.19; sa, a pointer to a sockaddr\_in structure containing the source IP address from the ICMP message; and ip, a pointer to the IP header that caused the error. For the destination unreachable, parameter problem, source quench, and time exceeded errors, the pointer ip points to the IP header that caused the error. But when udp\_ctlinput is called by pfctlinput for redirects (Figure 22.32), sa points to a sockaddr\_in structure containing the destination address that should be redirected, and ip is a null pointer. There is no loss of information in this final case, since we saw in Section 22.11 that a redirect is applied to all TCP and UDP sockets connected to the destination address. The nonnull third argument is needed, however, for other errors, such as a port unreachable, since the protocol header following the IP header contains the unreachable port.

323-325

If the error is not a redirect, and either the PRC\_xxx value is too large or there is no error code in the global array inetctlerrmap, the ICMP error is ignored. To understand this test we need to review what happens to a received ICMP message.

- 1. icmp\_input converts the ICMP type and code into a PRC\_xxx error code.
- 2. The PRC\_xxx error code is passed to the protocol's control-input function.

udp

301-

326-3

Section

Chapter 23

ter 23

udp\_usrreq.c 314 void 315 udp\_ctlinput(cmd, sa, ip) cmd; 316 int 317 struct sockaddr \*sa; 318 struct ip \*ip; 319 { struct udphdr \*uh; 320 extern struct in\_addr zeroin\_addr; 321 extern u\_char inetctlerrmap[]; 322 if (!PRC\_IS\_REDIRECT(cmd) && 323 ((unsigned) cmd >= PRC\_NCMDS || inetctlerrmap[cmd] == 0)) 324 325 return; 326 if (ip) { uh = (struct udphdr \*) ((caddr\_t) ip + (ip->ip\_hl << 2));</pre> 327 in\_pcbnotify(&udb, sa, uh->uh\_dport, ip->ip\_src, uh->uh\_sport, 328 cmd, udp\_notify); 329 } else 330 in\_pcbnotify(&udb, sa, 0, zeroin\_addr, 0, cmd, udp\_notify); 331 332 } udp\_usrreq.c

Figure 23.30 udp\_ctlinput function: process received ICMP errors.

3. The Internet protocols (TCP and UDP) map the PRC\_xxx error code into one of the Unix errno values using inetctlerrmap, and this value is returned to the process.

Figures 11.1 and 11.2 summarize this processing of ICMP messages.

Returning to Figure 23.30, we can see what happens to an ICMP source quench that arrives in response to a UDP datagram. icmp\_input converts the ICMP message into the error PRC\_QUENCH and udp\_ctlinput is called. But since the error column for this ICMP error is blank in Figure 11.2, the error is ignored.

The function in\_pcbnotify notifies the appropriate PCBs of the ICMP error. If the third argument to udp\_ctlinput is nonnull, the source and destination UDP ports from the datagram that caused the error are passed to in\_pcbnotify along with the source IP address.

### udp\_notify Function

The final argument to in\_pcbnotify is a pointer to a function that in\_pcbnotify calls for each PCB that is to receive the error. The function for UDP is udp\_notify and we show it in Figure 23.31.

<sup>301-313</sup> The errno value, the second argument to this function, is stored in the socket's so\_error variable. By setting this socket variable, the socket becomes readable and writable if the process calls select. Any processes waiting to receive or send on the socket are then awakened to receive the error.

robcol's

this

null .nce

; no

ler-

12 11 1

Chapter 23

udp\_usrreq.c

udp\_usrreq.c

15

. .

State of the second sec

Constant of the second

udp\_usrreq.c

Secl

41

42

43

43

4.

	static void
306	udp_notify(inp, errno)
307	struct inpcb *inp;
308	int errno;
309	{
310	<pre>inp-&gt;inp_socket-&gt;so_error = errno;</pre>
311	<pre>sorwakeup(inp-&gt;inp_socket);</pre>
312	<pre>sowwakeup(inp-&gt;inp_socket);</pre>
313	)

Figure 23.31 udp\_notify function: notify process of an asynchronous error.

# 23.10 udp\_usrreq Function

The protocol's user-request function is called for a variety of operations. We saw in Figure 23.14 that a call to any one of the five write functions on a UDP socket ends up calling UDP's user-request function with a request of PRU\_SEND.

Figure 23.32 shows the beginning and end of udp\_usrreq. The body of the switch is discussed in separate figures following this figure. The function arguments are described in Figure 15.17.

/\* switch cases \*/

```
417 int
418 udp_usrreq(so, req, m, addr, control)
419 struct socket *so;
            req;
420 int
421 struct mbuf *m, *addr, *control;
422 {
        struct inpcb *inp = sotoinpcb(so);
423
                error = 0;
        int
424
425
        int
                s;
        if (req == PRU_CONTROL)
426
            return (in_control(so, (int) m, (caddr_t) addr,
427
                                (struct ifnet *) control));
428
        if (inp == NULL && req != PRU_ATTACH) {
429
            error = EINVAL;
430
             goto release;
431
432
         }
433
         /*
          * Note: need to block udp_input while changing
434
          * the udp pcb queue and/or pcb addresses.
435
          */
436
         switch (reg) {
437
```

udp\_usrreq Function 785

– udp\_usrreq.c

Section 23.10

'3

.c

с

е

S

С

```
default:
522
            panic("udp_usrreq");
523
524
525
      release:
        if (control) {
526
            printf("udp control data unexpectedly retained\n");
527
             m_freem(control);
528
529
         }
         if (m)
530
             m_freem(m);
531
         return (error);
532
533 }
```

Figure 23.32 Body of udp\_usrreq function.

417-428 The PRU\_CONTROL request is from the ioctl system call. The function in\_control processes the request completely.

- <sup>429–432</sup> The socket pointer was converted to the PCB pointer when inp was declared at the beginning of the function. The only time a null PCB pointer is allowed is when a new socket is being created (PRU\_ATTACH).
- <sup>433-436</sup> The comment indicates that whenever entries are being added to or deleted from UDP's PCB list, the code must be protected by splnet. This is done because udp\_usrreq is called as part of a system call, and it doesn't want to be interrupted by UDP input (called by IP input, which is called as a software interrupt) while it is modifying the doubly linked list of PCBs. UDP input is also blocked while modifying the local or foreign addresses or ports in a PCB, to prevent a received UDP datagram from being delivered incorrectly by in\_pcblookup.

We now discuss the individual case statements. The PRU\_ATTACH request, shown in Figure 23.33, is from the socket system call.

<sup>438–447</sup> If the socket structure already points to a PCB, EINVAL is returned. in\_pcballoc allocates a new PCB, adds it to the front of UDP's PCB list, and links the socket structure and the PCB to each other.

448-450 soreserve reserves buffer space for a receive buffer and a send buffer for the socket. As noted in Figure 16.7, soreserve just enforces system limits; the buffer space is not actually allocated. The default values for the send and receive buffer sizes are 9216 bytes (udp\_sendspace) and 41,600 bytes (udp\_recvspace). The former allows for a maximum UDP datagram size of 9200 bytes (to hold 8 Kbytes of data in an NFS packet), plus the 16-byte sockaddr\_in structure for the destination address. The latter allows for 40 1024-byte datagrams to be queued at one time for the socket. The process can change these defaults by calling setsockopt.

451-452

There are two fields in the prototype IP header in the PCB that the process can change by calling setsockopt: the TTL and the TOS. The TTL defaults to 64 (ip\_deftt1) and the TOS defaults to 0 (normal service), since the PCB is initialized to 0 by in\_pcballoc.

456-460

		– udp_usrreq.c
438	case PRU_ATTACH:	uup_uorreg.e
439	if (inp != NULL) {	
440	error = EINVAL;	
441	break;	
442	}	
443	s = splnet();	
444	error = in_pcballoc(so, &udb);	
445	<pre>splx(s);</pre>	
446	if (error)	
447	break;	
448	error = soreserve(so, udp_sendspace, udp_recvspace);	
449	if (error)	
450	break;	
451	((struct inpcb *) so->so_pcb)->inp_ip.ip_ttl = ip_defttl;	
452	break;	
453	case PRU_DETACH:	
454	udp_detach(inp);	
455	break;	– udp_usrreq.c

Figure 23.33 udp\_usrreg function: PRU\_ATTACH and PRU\_DETACH requests.

453-455 The close system call issues the PRU\_DETACH request. The function udp\_detach, shown in Figure 23.34, is called. This function is also called later in this section for the PRU\_ABORT request.

534 s	static void	
535 u	udp_detach(inp)	
536 s	struct inpcb *inp;	
537 {	{	
538	<pre>int s = splnet();</pre>	
539	if (inp == udp_last_inpcb)	
540	udp_last_inpcb = &udb	
541	<pre>in_pcbdetach(inp);</pre>	
542	<pre>splx(s);</pre>	
543 }	}	

Figure 23.34 udp\_detach function: delete a UDP PCB.

If the last-received PCB pointer (the one-behind cache) points to the PCB being detached, the cache pointer is set to the head of the UDP list (udb). The function in\_pcbdetach removes the PCB from UDP's list and releases the PCB.

Returning to udp\_usrreq, a PRU\_BIND request is the result of the bind system call and a PRU\_LISTEN request is the result of the listen system call. Both are shown in Figure 23.35.

All the work for a PRU\_BIND request is done by in\_pcbbind.

<sup>461–463</sup> The PRU\_LISTEN request is invalid for a connectionless protocol—it is used only by connection-oriented protocols.

464-4

Chapter 23

udp\_usrreq.c

udp\_usrreq.c

787 udp\_usrreq Function Section 23.10 udp\_usrreq.c case PRU\_BIND: 456 s = splnet(); 457 error = in\_pcbbind(inp, addr); 458 splx(s); 459 break; 460 case PRU\_LISTEN: 461 error = EOPNOTSUPP; 462 break; udp\_usrreq.c 463

Figure 23.35 udp\_usrreq function: PRU\_BIND and PRU\_LISTEN requests.

We mentioned earlier that a UDP application, either a client or server (normally a client), can call connect. This fixes the foreign IP address and port number that this socket can send to or receive from. Figure 23.36 shows the PRU\_CONNECT, PRU\_CONNECT2, and PRU\_ACCEPT requests.

— udp\_usrreq.c

case PRU_CONNECT:	
error = EISCONN;	
break;	
}	
s = splnet();	
error = in_pcbconnect(inp, addr);	
<pre>splx(s);</pre>	
if (error == 0)	
soisconnected(so);	
break;	
case PRU_CONNECT2:	
error = EOPNOTSUPP;	
break;	
case PRU_ACCEPT:	
error = EOPNOTSUPP;	
break;udp_usrreq.c	
	<pre>if (inp-&gt;inp_faddr.s_addr != INADDR_ANY) (     error = EISCONN;     break;     }     s = splnet();     error = in_pcbconnect(inp, addr);     splx(s);     if (error == 0)         soisconnected(so);     break;  case PRU_CONNECT2:     error = EOPNOTSUPP;     break;  case PRU_ACCEPT:     error = EOPNOTSUPP; </pre>

Figure 23.36 udp\_usrreq function: PRU\_CONNECT, PRU\_CONNECT2, and PRU\_ACCEPT requests.

If the socket is already connected, EISCONN is returned. The socket should never be connected at this point, because a call to connect on an already-connected UDP socket generates a PRU\_DISCONNECT request before this PRU\_CONNECT request. Otherwise in\_pcbconnect does all the work. If no errors are encountered, soisconnected marks the socket structure as being connected.

The socketpair system call issues the PRU\_CONNECT2 request, which is defined only for the Unix domain protocols.

The PRU\_ACCEPT request is from the accept system call, which is defined only for connection-oriented protocols.

r 23

eq.c

eq.c

ion his:

eq.c

eq.c

ing

ion

em

wn

aly

788 UDP: U	Jser Datagram Protocol Chapter 23		Section
			495-49
Т	he PRU_DISCONNECT request can occur in two cases for a UDP socket:	the second s Second second s	495-4
	. When a connected UDP socket is closed, PRU_DISCONNECT is called before DBU_DETACH	1	
2	the second on an already-connected UDP socket, soconnected	t Maria	
	issues the PRU_DISCONNECT request before the PRU_CONNECT request.	575. 2	
Figu	re 23.37 shows the PRU_DISCONNECT request.	 	
	uup_usiteq	н. <b>с</b>	
481	uup_usiteq	τ. 	
481 482	uup_usrieu	. <b>.</b> .	
481 482 483	case PRU_DISCONNECT: if (inp->inp_faddr.s_addr == INADDR_ANY) {		
481 482 483 484	<pre>case PRU_DISCONNECT: if (inp-&gt;inp_faddr.s_addr == INADDR_ANY) { error = ENOTCONN;</pre>		497-5
481 482 483 484 485	<pre>case PRU_DISCONNECT: if (inp-&gt;inp_faddr.s_addr == INADDR_ANY) { error = ENOTCONN; break; } s = splnet();</pre>	A.C.	497-5
481 482 483 484 485 486	<pre>case PRU_DISCONNECT: if (inp-&gt;inp_faddr.s_addr == INADDR_ANY) { error = ENOTCONN; break; } s = splnet(); in ncbdisconnect(inp);</pre>		497-5
481 482 483 484 485	<pre>case PRU_DISCONNECT: if (inp-&gt;inp_faddr.s_addr == INADDR_ANY) { error = ENOTCONN; break; } s = splnet();</pre>		497-5
481 482 483 484 485 486 486 487	<pre>case PRU_DISCONNECT: if (inp-&gt;inp_faddr.s_addr == INADDR_ANY) { error = ENOTCONN; break; } s = splnet(); in_pcbdisconnect(inp); inp-&gt;inp_laddr.s_addr = INADDR_ANY; splx(s);</pre>		497-5
481 482 483 484 485 486 487 488	<pre>case PRU_DISCONNECT: if (inp-&gt;inp_faddr.s_addr == INADDR_ANY) { error = ENOTCONN; break; } s = splnet(); in_pcbdisconnect(inp); inp-&gt;inp_laddr.s_addr = INADDR_ANY; splx(s); so-&gt;so_state &amp;= ~SS_ISCONNECTED; /* XXX */</pre>		497-5
481 482 483 484 485 486 487 488 489	<pre>case PRU_DISCONNECT: if (inp-&gt;inp_faddr.s_addr == INADDR_ANY) { error = ENOTCONN; break; } s = splnet(); in_pcbdisconnect(inp); inp-&gt;inp_laddr.s_addr = INADDR_ANY; splx(s);</pre>		497– <u>5</u>
481 482 483 484 485 486 487 488 489 490	<pre>case PRU_DISCONNECT: if (inp-&gt;inp_faddr.s_addr == INADDR_ANY) { error = ENOTCONN; break; } s = splnet(); in_pcbdisconnect(inp); inp-&gt;inp_laddr.s_addr = INADDR_ANY; splx(s); so-&gt;so_state &amp;= ~SS_ISCONNECTED; /* XXX */</pre>		497-5

12.45 - 12.42 American American American American American

492-494

If the socket is not already connected, ENOTCOMN is returned in\_pcbdisconnect sets the foreign IP address to 0.0.0.0 and the foreign port to 0. The local address is also set to 0.0.0.0, since this PCB variable could have been set by connect.

A call to shutdown specifying that the process has finished sending data generates the PRU\_SHUTDOWN request, although it is rare for a process to issue this system call for a UDP socket. Figure 23.38 shows the PRU\_SHUTDOWN, PRU\_SEND, and PRU\_ABORT requests.

492 493 494	<pre>case PRU_SHUTDOWN:     socantsendmore(so);     break;</pre>
495 496	<pre>case PRU_SEND:     return (udp_output(inp, m, addr, control));</pre>
497 498 499 500	<pre>case PRU_ABORT: soisdisconnected(so); udp_detach(inp); break;</pre>

Figure 23.38 udp\_usrreq function: PRU\_SHUTDOWN, PRU\_SEND, and PRU\_ABORT requests.

socantsendmore sets the socket's flags to prevent any future output.

501-50

udp\_usrreq.c

udp usrreg Function 789

Section 23.10

е

С

<sup>495-496</sup> In Figure 23.14 we showed how the five write functions ended up calling udp\_usrreq with a PRU\_SEND request. udp\_output sends the datagram. udp\_usrreq returns, to avoid falling through to the label release (Figure 23.32), since the mbuf chain containing the data (m) must not be released yet. IP output appends this mbuf chain to the appropriate interface output queue, and the device driver will release the mbuf when the data has been transmitted.

The only buffering of UDP output within the kernel is on the interface's output queue. If there is room in the socket's send buffer for the datagram and destination address, sosend calls udp\_usrreq, which we see calls udp\_output. We saw in Figure 23.20 that ip\_output is then called, which calls ether\_output for an Ethernet, placing the datagram onto the interface's output queue (if there is room). If the process calls sendto faster than the interface can transmit the datagrams, ether\_output can return ENOBUFS, which is returned to the process.

497-500 A PRU\_ABORT request should never be generated for a UDP socket, but if it is, the socket is disconnected and the PCB detached.

The PRU\_SOCKADDR and PRU\_PEERADDR requests are from the getsockname and getpeername system calls, respectively. These two requests, and the PRU\_SENSF request, are shown in Figure 23.39.

udp\_usrreq...

501	case PRU_SOCKADDR:	· - ·
502	in_setsockaddr(inp, addr);	
503	break;	
504	case PRU_PEERADDR:	
505	<pre>in_setpeeraddr(inp, addr);</pre>	
506	break;	
507	case PRU_SENSE:	
508	/*	
509	* fstat: don't bother with a blocksize.	
510	*/	
511	return (0);	

Figure 23.39 udp\_usrreq function: PRU\_SOCKADDR, PRU\_PEERADDR, and PRU\_SENSE requests.

<sup>501-506</sup> The functions in\_setsockaddr and in\_setpeeraddr fetch the information from the PCB, storing the result in the addr argument.

507-511 The fstat system call generates the PRU\_SENSE request. The function returns OM, but doesn't return any other information. We'll see later that TCP returns the size of the send buffer as the st\_blksize element of the stat structure.

The remaining seven PRU\_xxx requests, shown in Figure 23.40, are not supported for a UDP socket.

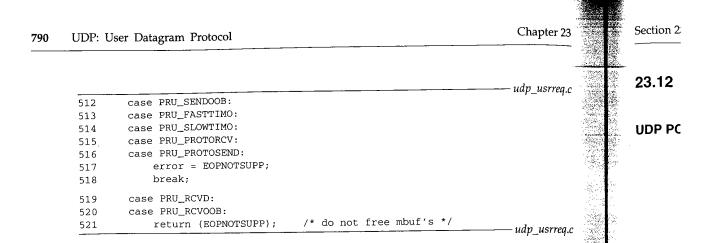


Figure 23.40 udp\_usrreq function: unsupported requests.

There is a slight difference in how the last two are handled because PRU\_RCVD doesn't pass a pointer to an mbuf as an argument (m is a null pointer) and PRU\_RCVOOB passes a pointer to an mbuf for the protocol to fill in. In both cases the error is immediately returned, without breaking out of the switch and releasing the mbuf chain. With PRU\_RCVOOB the caller releases the mbuf that it allocated.

## 23.11 udp\_sysct1 Function

The sysctl function for UDP supports only a single option, the UDP checksum flag. The system administrator can enable or disable UDP checksums using the sysctl(8) program. Figure 23.41 shows the udp\_sysctl function. This function calls sysctl\_int to fetch or set the value of the integer udpcksum.

	udp_usrreq.c
547 udp_sysctl(name, namelen, oldp, oldlenp, newp, newlen)	,_ ,
548 int *name;	
549 u_int namelen;	
550 void *oldp;	
551 size_t *oldlenp;	
552 void *newp;	
553 size_t newlen;	
554 {	
555 /* All sysctl names at this level are terminal. */	
556 if (namelen != 1)	
557 return (ENOTDIR);	
558 switch (name[0]) {	
559 case UDPCTL_CHECKSUM:	
560 return (sysctl_int(oldp, oldlenp, newp, newlen, &udpcksum)	);
561 default:	
562 return (ENOPROTOOPT);	
563 }	
564 /* NOTREACHED */	
565 }	udp_usrreq.c

Figure 23.41 udp\_sysct1 function.

~

## 23.12 Implementation Refinements

#### **UDP PCB Cache**

In Section 22.12 we talked about some general features of PCB searching and how the code we've seen uses a linear search of the protocol's PCB list. We now tie this together with the one-behind cache used by UDP in Figure 23.24.

The problem with the one-behind cache occurs when the cached PCB contains wildcard values (for either the local address, foreign address, or foreign port): the cached value never matches any received datagram. One solution tested in [Partridge and Pink 1993] is to modify the cache to not compare wildcarded values. That is, instead of comparing the foreign address in the PCB with the source address in the datagram, compare these two values only if the foreign address in the PCB is not a wildcard.

There's a subtle problem with this approach [Partridge and Pink 1993]. Assume there are two sockets bound to local port 555. One has the remaining three elements wildcarded, while the other has connected to the foreign address 128.1.2.3 and the foreign port 1600. If we cache the first PCB and a datagram arrives from 128.1.2.3, port 1600, we can't ignore comparing the foreign addresses just because the cached value has a wildcarded foreign address. This is called *cache hiding*. The cached PCB has hidden another PCB that is a better match in this example.

To get around cache hiding requires more work when a new entry is added to or deleted from the cache. Those PCBs that hide other PCBs cannot be cached. This is not a problem, however, because the normal scenario is to have one socket per local port. The example we just gave with two sockets bound to local port 555, while possible (especially on a multihomed host), is rare.

The next enhancement tested in [Partridge and Pink 1993] is to also remember the PCB of the last datagram sent. This is motivated by [Mogul 1991], who shows that half of all datagrams received are replies to the last datagram that was sent. Cache hiding is a problem here also, so PCBs that would hide other PCBs are not cached.

The results of these two caches shown in [Partridge and Pink 1993] on a generalpurpose system measured for around 100,000 received UDP datagrams show a 57% hit rate for the last-received PCB cache and a 30% hit rate for the last-sent PCB cache. The amount of CPU time spent in udp\_input is more than halved, compared to the version with no caching.

These two caches still depend on a certain amount of locality: that with a high probability the UDP datagram that just arrived is either from the same peer as the last UDP datagram received or from the peer to whom the last datagram was sent. The latter is typical for request-response applications that send a datagram and wait for a reply. [McKenney and Dove 1992] show that some applications, such as data entry into an online transaction processing (OLTP) system, don't yield the high cache hit rates that [Partridge and Pink 1993] observed. As we mentioned in Section 22.12, placing the PCBs onto hash chains provided an order of magnitude improvement over the last-received and last-sent caches for a system with thousands of OLTP connections.

Chapter 23

2

#### UDP Checksum

The next area for improving the implementation is to combine the copying of data between the process and the kernel with the calculation of the checksum. In Net/3, each byte of data is processed twice during an output operation: once when copied from the process into an mbuf (the function uiomove, which is called by sosend), and again when the UDP checksum is calculated (by the function in\_cksum, which is called by udp\_output). This happens on input as well as output.

[Partridge and Pink 1993] modified the UDP output processing from what we showed in Figure 23.14 so that a UDP-specific function named udp\_sosend is called instead of sosend. This new function calculates the checksum of the UDP header and the pseudo-header in-line (instead of calling the general-purpose function in\_cksum) and then copies the data from the process into an mbuf chain using a special function named in\_uiomove (instead of the general-purpose uiomove). This new function copies the data *and* updates the checksum. The amount of time spent copying the data and calculating the checksum is reduced with this technique by about 40 to 45%.

On the receive side the scenario is different. UDP calculates the checksum of the UDP header and the pseudo-header, removes the UDP header, and queues the data for the appropriate socket. When the application reads the data, a special version of soreceive (called udp\_soreceive) completes the calculation of the checksum while copying the data into the user's buffer. If the checksum is in error, however, the error is not detected until the entire datagram has been copied into the user's buffer. In the normal case of a blocking socket, udp\_soreceive just waits for the next datagram to arrive. But if the socket is nonblocking, the error EWOULDBLOCK must be returned if another datagram is not ready to be passed to the process. This implies two changes in the socket interface for a nonblocking read from a UDP socket:

- The select function can indicate that a nonblocking UDP socket is readable, yet the error EWOULDBLOCK is unexpectedly returned by one of the read functions if the checksum fails.
- 2. Since a checksum error is detected after the datagram has been copied into the user's buffer, the application's buffer is changed even though no data is returned by the read.

Even with a blocking socket, if the datagram with the checksum error contains 100 bytes of data and the next datagram without an error contains 40 bytes of data, recvfrom returns a length of 40, but the 60 bytes that follow in the user's buffer have also been modified.

[Partridge and Pink 1993] compare the timings for a copy versus a copy-withchecksum for six different computers. They show that the checksum is calculated for free during the copy operation on many architectures. This occurs when memory access speeds and CPU processing speeds are mismatched, as is true for many current RISC processors. Chapter 23

## 23.13 Summary

UDP is a simple, connectionless protocol, which is why we cover it before looking at TCP. UDP output is simple: IP and UDP headers are prepended to the user's data, as much of the header is filled in as possible, and the result is passed to ip\_output. The only complication is calculating the UDP checksum, which involves prepending a pseudo-header just for the checksum computation. We'll encounter a similar pseudo-header for the calculation of the TCP checksum in Chapter 26.

When udp\_input receives a datagram, it first performs a general validation (the length and checksum); the processing then differs depending on whether the destination IP address is a unicast address or a broadcast or multicast address. A unicast datagram is delivered to at most one process, but a broadcast or multicast datagram can be delivered to multiple processes. A one-behind cache is maintained for unicast datagrams, which maintains a pointer to the last Internet PCB for which a UDP datagram was received. We saw, however, that because of the prevalence of wildcard addressing with UDP applications, this cache is practically useless.

The udp\_ctlinput function is called to handle received ICMP messages, and the udp\_usrreq function handles the PRU\_xxx requests from the socket layer.

#### Exercises

- 23.1 List the five types of mbuf chains that udp\_output passes to ip\_output. (*Hint*: look at sosend.)
- **23.2** What happens to the answer for the previous exercise when the process specifies IP options for the outgoing datagram?
- 23.3 Does a UDP client need to call bind? Why or why not?
- **23.4** What happens to the processor priority level in udp\_output if the socket is unconnected and the call to M\_PREPEND in Figure 23.15 fails?
- **23.5** udp\_output does not check for a destination port of 0. Is it possible to send a UDP datagram with a destination port of 0?
- **23.6** Assuming the IP\_RECVDSTADDR socket option worked when a datagram was sent to a broadcast address, how can you then determine if this address is a broadcast address?
- 23.7 Who releases the mbuf that udp\_saveopt (Figure 23.28) allocates?
- 23.8 How can a process disconnect a connected UDP socket? That is, the process calls connect and exchanges datagrams with that peer, and then the process wants to disconnect the socket, allowing it to call sendto and send a datagram to some other host.
- 23.9 In our discussion of Figure 22.25 we noted that a UDP application that calls connect with a foreign IP address of 255.255.255.255 actually sends datagrams out the primary interface with a destination IP address corresponding to the broadcast address of that interface. What happens if a UDP application uses an unconnected socket instead, calling sendto with a destination address of 255.255.255.255.255?

23

ta

3,

n

n

·y

e d

£

ı)

1

٦

9

2

r

1.00

2

- **23.10** After discussing the problem with Figure 23.27, we mentioned that this problem would not exist if the server used the destination IP address from the request as the source IP address of the reply. Explain how the server could do this.
- **23.11** Implement changes to allow a process to perform path MTU discovery using UDP: the process must be able to set the "don't fragment" bit in the resulting IP datagram and be told if the corresponding ICMP destination unreachable error is received.
- 23.12 Does the variable udp\_in need to be global?
- 23.13 Modify udp\_input to save the IP options and make them available to the receiver with the IP\_RECVOPTS socket option.
- **23.14** Fix the one-behind cache in Figure 23.24.
- **23.15** Fix udp\_input to implement the IP\_RECVOPTS and IP\_RETOPTS socket options.
- **23.16** Fix udp\_input so that the IP\_RECVDSTADDR socket option works for datagrams sent to a broadcast or multicast address.

24.