

```

231  /* udp_usrreq.c
232  * Construct sockaddr format source address.
233  * Stuff source address and datagram in user buffer.
234  */
235  udp_in.sin_port = uh->uh_sport;
236  udp_in.sin_addr = ip->ip_src;
237  if (inp->inp_flags & INP_CONTROLOPTS) {
238      struct mbuf **mp = &opts;
239      if (inp->inp_flags & INP_RECVDSTADDR) {
240          *mp = udp_saveopt((caddr_t) & ip->ip_dst,
241                          sizeof(struct in_addr), IP_RECVDSTADDR);
242          if (*mp)
243              mp = &(*mp)->m_next;
244      }
245  #ifndef 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      }
253      /* ip_srcroute doesn't do what we want here, need to fix */
254      if (inp->inp_flags & INP_RECVRETOPTS) {
255          *mp = udp_saveopt((caddr_t) ip_srcroute(),
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;
266  if (sbappendaddr(&inp->inp_socket->so_rcv, (struct sockaddr *) &udp_in,
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

231-238 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 datagram 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

267-272 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, `sockwakep` is called for the receiving socket to wake up any processes asleep on the socket's receive queue.

Error return

273-276 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.

127-138 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 unreachable) unless the socket is connected. We described this in Section 22.11.

139-145 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.

146-164 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 `sbwakeup`. 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.

```

121 if (IN_MULTICAST(ntohl(ip->ip_dst.s_addr)) ||
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      * (This really ought to be done for unicast datagrams as
130      * well, but that would cause problems with existing
131      * applications that open both address-specific sockets and
132      * a wildcard socket listening to the same port -- they would
133      * end up receiving duplicates of every unicast datagram,
134      * Those applications open the multiple sockets to overcome an
135      * inadequacy of the UDP socket interface, but for backwards
136      * compatibility we avoid the problem here rather than
137      * fixing the interface. Maybe 4.5BSD will remedy this?)
138      */
139     /*
140      * Construct sockaddr format source address.
141      */
142     udp_in.sin_port = uh->uh_sport;
143     udp_in.sin_addr = ip->ip_src;
144     m->m_len -= sizeof(struct udpiphdr);
145     m->m_data += sizeof(struct udpiphdr);
146     /*
147      * Locate pcb(s) for datagram.
148      * (Algorithm copied from raw_intr().)
149      */
150     last = NULL;
151     for (inp = udb.inp_next; inp != &udb; inp = inp->inp_next) {
152         if (inp->inp_lport != uh->uh_dport)
153             continue;
154         if (inp->inp_laddr.s_addr != INADDR_ANY) {
155             if (inp->inp_laddr.s_addr !=
156                 ip->ip_dst.s_addr)
157                 continue;
158         }
159         if (inp->inp_faddr.s_addr != INADDR_ANY) {
160             if (inp->inp_faddr.s_addr !=
161                 ip->ip_src.s_addr ||
162                 inp->inp_fport != uh->uh_sport)
163                 continue;
164         }
165         if (last != NULL) {
166             struct mbuf *n;
167             if ((n = m_copy(m, 0, M_COPYALL)) != NULL) {
168                 if (sbappendaddr(&last->so_rcv,
169                                 (struct sockaddr *) &udp_in,
170                                 n, (struct mbuf *) 0) == 0) {
171                     m_freem(n);
172                     udpstat.udps_fullsock++;

```



```

173         } else
174             sorwakeup(last);
175     }
176 }
177 last = inp->inp_socket;
178 /*
179  * Don't look for additional matches if this one does
180  * not have either the SO_REUSEPORT or SO_REUSEADDR
181  * socket options set. This heuristic avoids searching
182  * through all pcb's in the common case of a non-shared
183  * port. It assumes that an application will never
184  * clear these options after setting them.
185  */
186 if ((last->so_options & (SO_REUSEPORT | SO_REUSEADDR) == 0))
187     break;
188 }
189 if (last == NULL) {
190     /*
191     * No matching pcb found; discard datagram.
192     * (No need to send an ICMP Port Unreachable
193     * for a broadcast or multicast datagram.)
194     */
195     udpstat.udps_noportbcst++;
196     goto bad;
197 }
198 if (sbappendaddr(&last->so_rcv, (struct sockaddr *) &udp_in,
199                 m, (struct mbuf *) 0) == 0) {
200     udpstat.udps_fullsock++;
201     goto bad;
202 }
203 sorwakeup(last);
204 return;
205 }

```

udp_usrreq.c

Figure 23.26 udp_input function: demultiplexing of broadcast and multicast datagrams.

198-204 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.

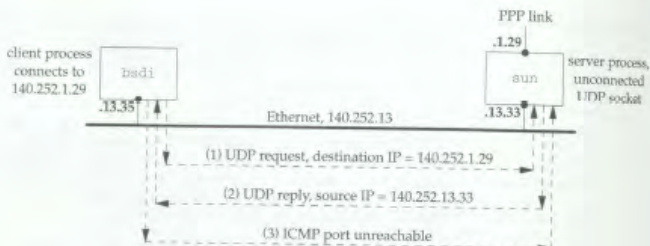


Figure 23.27 Example of connected UDP socket sending datagram to a multihomed host.

Three steps take place.

1. The client on `bsd1` 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 destination 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 `bsd1` it is not delivered to the client's connected UDP socket since the IP addresses don't match.

3. `bsd1` generates an ICMP port unreachable error since the reply can't be demultiplexed. (This assumes that there is not another process on `bsd1` 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.

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`:

```
*mp = udp_saveopt((caddr_t) &ip->ip_dst, sizeof(struct in_addr),
                 IP_RECVDSTADDR);
```

Figure 23.28 shows this function.

```
278 /* udp_usrreq.c
279  * Create a "control" mbuf containing the specified data
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 {
288     struct cmsghdr *cp;
289     struct mbuf *m;
290     if ((m = m_get(M_DONTWAIT, MT_CONTROL)) == NULL)
291         return ((struct mbuf *) NULL);
292     cp = (struct cmsghdr *) mtod(m, struct cmsghdr *);
293     bcopy(p, CMSG_DATA(cp), size);
294     size += sizeof(*cp);
295     m->m_len = size;
296     cp->cmsg_len = size;
297     cp->cmsg_level = IPPROTO_IP;
298     cp->cmsg_type = type;
299     return (m);
300 }
```

Figure 23.28 `udp_saveopt` function: create mbuf with control information.

278-289 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`).

290-299 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.

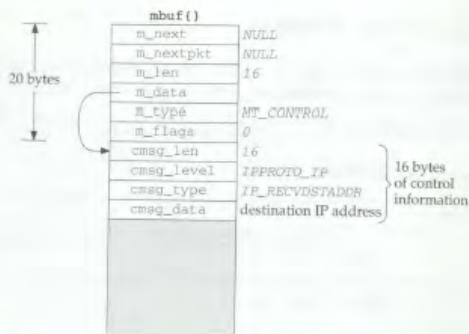


Figure 23.29 Mbuf containing destination address from received datagram as control information.

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:

```
if (ctlfunc = inetaw[ ip_protox[icp->icmp_ip.ip_p] ].pr_ctlinput)
    (*ctlfunc)(code, (struct sockaddr *) &icmpsrc, &icp->icmp_ip);
```

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.


```

314 void                                     udp_usrreq.c
315 udp_ctlinput(cmd, sa, ip)
316 int cmd;
317 struct sockaddr *sa;
318 struct ip *ip;
319 {
320     struct udphdr *uh;
321     extern struct in_addr zeroin_addr;
322     extern u_char inetctlerrmap[];
323     if (!PRC_IS_REDIRECT(cmd) &&
324         (unsigned) cmd >= PRC_NCMDS || inetctlerrmap[cmd] == 0)
325         return;
326     if (ip) {
327         uh = (struct udphdr *) ((caddr_t) ip + (ip->ip_hl << 2));
328         in_pcbnotify(&udb, sa, uh->uh_dport, ip->ip_src, uh->uh_sport,
329                     cmd, udp_notify);
330     } else
331         in_pcbnotify(&udb, sa, 0, zeroin_addr, 0, cmd, udp_notify);
332 }

```

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 `errno` column for this ICMP error is blank in Figure 11.2, the error is ignored.

326-332 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.

301-313 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.

```

305 static void
306 udp_notify(inp, errno)
307 struct inpcb *inp;
308 int     errno;
309 {
310     inp->inp_socket->so_error = errno;
311     sowakeup(inp->inp_socket);
312     scwakeup(inp->inp_socket);
313 }

```

udp_usrreq.c

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.

```

417 int
418 udp_usrreq(so, req, m, addr, control)
419 struct socket *so;
420 int     req;
421 struct mbuf *m, *addr, *control;
422 {
423     struct inpcb *inp = sotoinpcb(so);
424     int     error = 0;
425     int     s;
426
427     if (req == PRU_CONTROL)
428         return (in_control(so, (int) m, (caddr_t) addr,
429                               (struct ifnet *) control));
430     if (inp == NULL && req != PRU_ATTACH) {
431         error = EINVAL;
432         goto release;
433     }
434     /*
435      * Note: need to block udp_input while changing
436      * the udp pcb queue and/or pcb addresses.
437      */
438     switch (req) {

```

/* switch cases */

```

522     default:
523         panic("udp_usrreq");
524     }
525     release:
526     if (control)
527         printf("udp control data unexpectedly retained\n");
528         m_freem(control);
529     }
530     if (m)
531         m_freem(m);
532     return (error);
533 }

```

udp_usrreq.c

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.

429-432 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`).

433-438 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.

439-447 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_defttl`) and the TOS defaults to 0 (normal service), since the PCB is initialized to 0 by `in_pcballoc`.

```

438     case PRU_ATTACH:
439         if (inp != NULL) {
440             error = EINVAL;
441             break;
442         }
443         s = splnet();
444         error = in_pcballot(so, &udb);
445         splx(s);
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_usrreq* 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 static void
535 udp_detach(inp)
536 struct inpcb *inp;
537 {
538     int    s = splnet();
539     if (inp == udp_last_inpcb)
540         udp_last_inpcb = &udb;
541     in_pcbdetach(inp);
542     splx(s);
543 }

```

udp_usrreq.c

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.

456-460 All the work for a PRU_BIND request is done by `in_pcbbind`.

461-463 The PRU_LISTEN request is invalid for a connectionless protocol—it is used only by connection-oriented protocols.


```

456 case PRU_BIND:
457     s = splnet();
458     error = in_pcbbind(inp, addr);
459     splx(s);
460     break;

461 case PRU_LISTEN:
462     error = EOPNOTSUPP;
463     break;

```

udp_usrreq.c

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.

```

464 case PRU_CONNECT:
465     if (inp->inp_faddr.s_addr != INADDR_ANY) {
466         error = EISCONN;
467         break;
468     }
469     s = splnet();
470     error = in_pcbconnect(inp, addr);
471     splx(s);
472     if (error == 0)
473         soisconnected(so);
474     break;

475 case PRU_CONNECT2:
476     error = EOPNOTSUPP;
477     break;

478 case PRU_ACCEPT:
479     error = EOPNOTSUPP;
480     break;

```

udp_usrreq.c

Figure 23.36 udp_usrreq function: PRU_CONNECT, PRU_CONNECT2, and PRU_ACCEPT requests.

464-474 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.

475-477 The `socketpair` system call issues the `PRU_CONNECT2` request, which is defined only for the Unix domain protocols.

478-480 The `PRU_ACCEPT` request is from the `accept` system call, which is defined only for connection-oriented protocols.

The PRU_DISCONNECT request can occur in two cases for a UDP socket:

1. When a connected UDP socket is closed, PRU_DISCONNECT is called before PRU_DETACH.
2. When a connect is issued on an already-connected UDP socket, soconnect issues the PRU_DISCONNECT request before the PRU_CONNECT request.

Figure 23.37 shows the PRU_DISCONNECT request.

```

481     case PRU_DISCONNECT:
482         if (inp->inp_faddr.s_addr == INADDR_ANY) {
483             error = ENOTCONN;
484             break;
485         }
486         s = splnet();
487         in_pcbdisconnect(inp);
488         inp->inp_laddr.s_addr = INADDR_ANY;
489         splx(s);
490         so->so_state &= ~SS_ISCONNECTED; /* XXX */
491         break;

```

— udp_usrreq.c

— udp_usrreq.c

Figure 23.37 udp_usrreq function: PRU_DISCONNECT request.

If the socket is not already connected, ENOTCONN is returned. Otherwise 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     case PRU_SHUTDOWN:
493         socantsendmore(so);
494         break;
495     case PRU_SEND:
496         return (udp_output(inp, m, addr, control));
497     case PRU_ABORT:
498         soisdisconnected(so);
499         udp_detach(inp);
500         break;

```

— udp_usrreq.c

Figure 23.38 udp_usrreq function: PRU_SHUTDOWN, PRU_SEND, and PRU_ABORT requests.

492-494 socantsendmore sets the socket's flags to prevent any future output.

495-496 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, `send` 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 `send` 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_SENSE` request, are shown in Figure 23.39.

```

501     case PRU_SOCKADDR:
502         in_setsockaddr(inp, addr);
503         break;
504     case PRU_PEERADDR:
505         in_setpeeraddr(inp, addr);
506         break;
507     case PRU_SENSE:
508         /*
509          * fstat: don't bother with a blocksize.
510          */
511         return (0);

```

udp_usrreq.c

Figure 23.39 `udp_usrreq` function: `PRU_SOCKADDR`, `PRU_PEERADDR`, and `PRU_SENSE` requests.

501-506 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 OK, 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.

```

512     case PRU_SENDOOB:
513     case PRU_FASTTIMO:
514     case PRU_SLOWTIMO:
515     case PRU_PROTORCV:
516     case PRU_PROTOSEND:
517         error = EOPNOTSUPP;
518         break;
519     case PRU_RCVD:
520     case PRU_RCVOOB:
521         return (EOPNOTSUPP); /* do not free mbuf's */

```

udp_usrreq.c

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_sysctl 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 `udpcsum`.

```

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, &udpcsum));
561     default:
562         return (ENOPROTOOPT);
563     }
564     /* NOTREACHED */
565 }

```

udp_usrreq.c

Figure 23.41 udp_sysctl 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 general-purpose 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 on-line 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.

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:

1. 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-with-checksum 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.

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_RECVSTADDE` 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?

- 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

TCP: Transmission Control Protocol

24.1 Introduction

The Transmission Control Protocol, or TCP, provides a connection-oriented, reliable, byte-stream service between the two end points of an application. This is completely different from UDP's connectionless, unreliable, datagram service.

The implementation of UDP presented in Chapter 23 comprised 9 functions and about 800 lines of C code. The TCP implementation we're about to describe comprises 28 functions and almost 4,500 lines of C code. Therefore we divide the presentation of TCP into multiple chapters.

These chapters are not an introduction to TCP. We assume the reader is familiar with the operation of TCP from Chapters 17–24 of Volume 1.

24.2 Code Introduction

The TCP functions appear in six C files and numerous TCP definitions are in seven headers, as shown in Figure 24.1.

Figure 24.2 shows the relationship of the various TCP functions to other kernel functions. The shaded ellipses are the nine main TCP functions that we cover. Eight of these functions appear in the TCP `protosw` structure (Figure 24.8) and the ninth is `tcp_output`.

File	Description
netinet/tcp.h	tcphdr structure definition
netinet/tcp_debug.h	tcp_debug structure definition
netinet/tcp_fsm.h	definitions for TCP's finite state machine
netinet/tcp_seq.h	macros for comparing TCP sequence numbers
netinet/tcp_timer.h	definitions for TCP timers
netinet/tcp_var.h	tcpcb (control block) and tcpstat (statistics) structure definitions
netinet/tcpip.h	TCP plus IP header definition
netinet/tcp_debug.c	support for SO_DEBUG socket debugging (Section 27.10)
netinet/tcp_input.c	tcp_input and ancillary functions (Chapters 28 and 29)
netinet/tcp_output.c	tcp_output and ancillary functions (Chapter 26)
netinet/tcp_subr.c	miscellaneous TCP subroutines (Chapter 27)
netinet/tcp_timer.c	TCP timer handling (Chapter 25)
netinet/tcp_usrreq.c	PRU_xxx request handling (Chapter 30)

Figure 24.1 Files discussed in the TCP chapters.

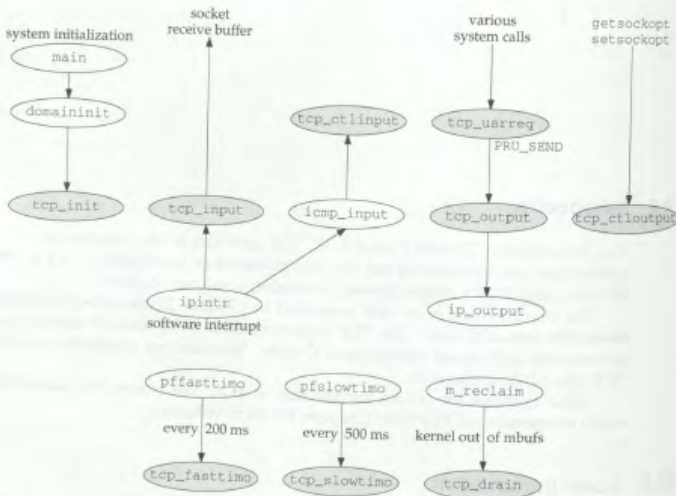


Figure 24.2 Relationship of TCP functions to rest of the kernel.

Global Variables

Figure 24.3 shows the global variables we encounter throughout the TCP functions.

Variable	Datatype	Description
<code>tcp</code>	<code>struct inpcb</code>	head of the TCP Internet PCB list
<code>tcp_last_inpcb</code>	<code>struct inpcb *</code>	pointer to PCB for last received segment; one-behind cache
<code>tcpstat</code>	<code>struct tcpstat</code>	TCP statistics (Figure 24.4)
<code>tcp_outflags</code>	<code>u_char</code>	array of output flags, indexed by connection state (Figure 24.16)
<code>tcp_recvspace</code>	<code>u_long</code>	default size of socket receive buffer (8192 bytes)
<code>tcp_sendspace</code>	<code>u_long</code>	default size of socket send buffer (8192 bytes)
<code>tcp_iss</code>	<code>tcp_seq</code>	initial send sequence number (ISS)
<code>tcp_rmxmtthresh</code>	<code>int</code>	number of duplicate ACKs to trigger fast retransmit (3)
<code>tcp_msdfilt</code>	<code>int</code>	default MSS (512 bytes)
<code>tcp_rttdfilt</code>	<code>int</code>	default RTT if no data (3 seconds)
<code>tcp_so_rfc1323</code>	<code>int</code>	if true (default), request window scale and timestamp options
<code>tcp_now</code>	<code>u_long</code>	500 ms counter for RFC 1323 timestamps
<code>tcp_keepidle</code>	<code>int</code>	keepalive: idle time before first probe (2 hours)
<code>tcp_keepintvl</code>	<code>int</code>	keepalive: interval between probes when no response (75 sec) (also used as timeout for <code>connect</code>)
<code>tcp_maxidle</code>	<code>int</code>	keepalive: time after probing before giving up (10 min)

Figure 24.3 Global variables introduced in the following chapters.

Statistics

Various TCP statistics are maintained in the global structure `tcpstat`, described in Figure 24.4. We'll see where these counters are incremented as we proceed through the code.

Figure 24.5 shows some sample output of these statistics, from the `netstat -s` command. These statistics were collected after the host had been up for 30 days. Since some counters come in pairs—one counts the number of packets and the other the number of bytes—we abbreviate these in the figure. For example, the two counters for the second line of the table are `tcps_sndpack` and `tcps_sndbyte`.

The counter for `tcps_sndbyte` should be 3,722,884,824, not -22,194,928 bytes. This is an average of about 405 bytes per segment, which makes sense. Similarly, the counter for `tcps_rcvackbyte` should be 3,738,811,552, not -21,264,360 bytes (for an average of about 565 bytes per segment). These numbers are incorrectly printed as negative numbers because the `printf` calls in the `netstat` program use `%d` (signed decimal) instead of `%lu` (long integer, unsigned decimal). All the counters are unsigned long integers, and these two counters are near the maximum value of an unsigned 32-bit long integer ($2^{32} - 1 = 4,294,967,295$).

tcpstat member	Description	Used by SNMP
tcpstat_accepts	#connections accepted passively	•
tcpstat_closed	#connections closed (includes drops)	•
tcpstat_connatempt	#connections initiated (calls to connect)	•
tcpstat_conndrops	#embryonic connections dropped (before SYN received)	•
tcpstat_connects	#connections established actively (by connect)	•
tcpstat_delayack	#delayed ACKs sent	
tcpstat_drops	#connections dropped (after SYN received)	•
tcpstat_keepdrops	#connections dropped in keepalive (established or awaiting SYN)	
tcpstat_keepprobe	#keepalive probes sent	
tcpstat_keeptimeo	#times keepalive timer or connection-establishment timer expire	
tcpstat_pawdrop	#segments dropped due to PAWS	
tcpstat_pcbcachemiss	#times PCB cache comparison fails	
tcpstat_persisttimeo	#times persist timer expires	
tcpstat_predack	#times header prediction correct for ACKs	
tcpstat_predat	#times header prediction correct for data packets	
tcpstat_rcvackbyte	#bytes ACKed by received ACKs	
tcpstat_rcvackpack	#received ACK packets	
tcpstat_rcvacktoomuch	#received ACKs for unsent data	
tcpstat_rcvafterclose	#packets received after connection closed	
tcpstat_rcvbadoff	#packets received with invalid header length	•
tcpstat_rcvbadsum	#packets received with checksum errors	•
tcpstat_rcvbyte	#bytes received in sequence	
tcpstat_rcvbyteafterwin	#bytes received beyond advertised window	
tcpstat_rcvdupack	#duplicate ACKs received	
tcpstat_rcvdupbyte	#bytes received in completely duplicate packets	
tcpstat_rcvduppack	#packets received with completely duplicate bytes	
tcpstat_rcvoobbyte	#out-of-order bytes received	
tcpstat_rcvoopack	#out-of-order packets received	
tcpstat_rcvpack	#packets received in sequence	
tcpstat_rcvpackafterwin	#packets with some data beyond advertised window	
tcpstat_rcvpartdupbyte	#duplicate bytes in part-duplicate packets	
tcpstat_rcvpartduppack	#packets with some duplicate data	
tcpstat_rcvshort	#packets received too short	•
tcpstat_rcvtotal	total #packets received	•
tcpstat_rcwinprobe	#window probe packets received	
tcpstat_rcwinupd	#received window update packets	
tcpstat_rexmttimeo	#retransmit timeouts	
tcpstat_rttupdated	#times RTT estimators updated	
tcpstat_segatimed	#segments for which TCP tried to measure RTT	
tcpstat_sndacks	#ACK-only packets sent (data length = 0)	
tcpstat_sndbyte	#data bytes sent	
tcpstat_sndctl	#control (SYN, FIN, RST) packets sent (data length = 0)	
tcpstat_sndpack	#data packets sent (data length > 0)	
tcpstat_sndprobe	#window probes sent (1 byte of data forced by persist timer)	
tcpstat_sndrexmitbyte	#data bytes retransmitted	•
tcpstat_sndrexmitpack	#data packets retransmitted	•
tcpstat_sndtotal	total #packets sent	•
tcpstat_sndurg	#packets sent with URG-only (data length = 0)	
tcpstat_sndwinupd	#window update-only packets sent (data length = 0)	
tcpstat_timeoutdrop	#connections dropped in retransmission timeout	

Figure 24.4 TCP statistics maintained in the tcpstat structure.

netstat -s output	tcpstat members
10,655,999 packets sent	tcps_sndtotal
9,177,823 data packets (-22,194,928 bytes)	tcps_snd(pack,byte)
257,295 data packets (81,075,086 bytes) retransmitted	tcps_sndrexit(pack,byte)
862,900 ack-only packets (531,285 delayed)	tcps_sndacks,tcps_slack
229 URG-only packets	tcps_sndurg
3,453 window probe packets	tcps_sndprobe
74,825 window update packets	tcps_sndwinup
279,387 control packets	tcps_sndctrl
8,801,953 packets received	tcps_rcvtotal
6,617,079 acks (for -21,264,360 bytes)	tcps_rcvack(pack,byte)
235,311 duplicate acks	tcps_rcvdupack
0 acks for unsent data	tcps_rcvacktoomuch
4,670,615 packets (324,965,351 bytes) rcvd in-sequence	tcps_rcv(pack,byte)
46,953 completely duplicate packets (1,549,785 bytes)	tcps_rcvdup(pack,byte)
22 old duplicate packets	tcps_rcwdrop
3,442 packets with some dup. data (54,483 bytes duped)	tcps_rcvpartdup(pack,byte)
77,114 out-of-order packets (13,938,456 bytes)	tcps_rcvoo(pack,byte)
1,892 packets (1,755 bytes) of data after window	tcps_rcv(pack,byte)afterwin
1,755 window probes	tcps_rcvwinprobe
175,476 window update packets	tcps_rcvwinup
1,017 packets received after close	tcps_rcvafterclose
60,370 discarded for bad checksums	tcps_rcvbadsum
279 discarded for bad header offset fields	tcps_rcvbadoff
0 discarded because packet too short	tcps_rcvshort
144,020 connection requests	tcps_connattemp
92,595 connection accepts	tcps_accepts
126,820 connections established (including accepts)	tcps_connects
237,743 connections closed (including 1,061 drops)	tcps_closed,tcps_drops
110,016 embryonic connections dropped	tcps_conndrops
6,363,546 segments updated rtt (of 6,444,667 attempts)	tcps_(rttupdated,segstimed)
114,797 retransmit timeouts	tcps_rexmttimeo
86 connection dropped by rexmit timeout	tcps_timeoutdrop
1,173 persist timeouts	tcps_persisttimeo
16,419 keepalive timeouts	tcps_keeptimeo
6,899 keepalive probes sent	tcps_keepprobe
3,219 connections dropped by keepalive	tcps_keepprops
733,130 correct ACK header predictions	tcps_predack
1,266,889 correct data packet header predictions	tcps_preddat
1,851,557 cache misses	tcps_pcbcachemiss

Figure 24.5 Sample TCP statistics.

SNMP Variables

Figure 24.6 shows the 14 simple SNMP variables in the TCP group and the counters from the `tcpstat` structure implementing that variable. The constant values shown for the first four entries are fixed by the Net/3 implementation. The counter `tcpCurrEstab` is computed as the number of Internet PCBs on the TCP PCB list.

Figure 24.7 shows `tcpTable`, the TCP listener table.

SNMP variable	tcpstat members or constant	Description
tcpRtoAlgorithm	4	algorithm used to calculate retransmission timeout value: 1 = none of the following, 2 = a constant RTO, 3 = MIL-STD-1778 Appendix B, 4 = Van Jacobson's algorithm.
tcpRtoMin	1000	minimum retransmission timeout value, in milliseconds
tcpRtoMax	64000	maximum retransmission timeout value, in milliseconds
tcpMaxConn	-1	maximum #TCP connections (-1 if dynamic)
tcpActiveOpens	tcps_connatempt	#transitions from CLOSED to SYN_SENT states
tcpPassiveOpens	tcps_accepts	#transitions from LISTEN to SYN_RCVD states
tcpAttemptFails	tcps_conndrops	#transitions from SYN_SENT or SYN_RCVD to CLOSED, plus #transitions from SYN_RCVD to LISTEN
tcpEstabResets	tcps_drops	#transitions from ESTABLISHED or CLOSE_WAIT states to CLOSED
tcpCurrEstab	(see text)	#connections currently in ESTABLISHED or CLOSE_WAIT states
tcpInSegs	tcps_rcvtotal	total #segments received
tcpOutSegs	tcps_sndtotal - tcps_sndrexmitpack	total #segments sent, excluding those containing only retransmitted bytes
tcpRetransSegs	tcps_sndrexmitpack	total #retransmitted segments
tcpInErrs	tcps_rcvbadsum + tcps_rcvbadoff + tcps_rcvshort	total #segments received with an error
tcpOutRsts	(not implemented)	total #segments sent with RST flag set

Figure 24.6 Simple SNMP variables in tcp group.

index = <tcpConnLocalAddress><tcpConnLocalPort><tcpConnRemAddress><tcpConnRemPort>		
SNMP variable	PCB variable	Description
tcpConnState	t_state	state of connection: 1 = CLOSED, 2 = LISTEN, 3 = SYN_SENT, 4 = SYN_RCVD, 5 = ESTABLISHED, 6 = FIN_WAIT_1, 7 = FIN_WAIT_2, 8 = CLOSE_WAIT, 9 = LAST_ACK, 10 = CLOSING, 11 = TIME_WAIT, 12 = delete TCP control block.
tcpConnLocalAddress	inp_laddr	local IP address
tcpConnLocalPort	inp_lport	local port number
tcpConnRemAddress	inp_faddr	foreign IP address
tcpConnRemPort	inp_fport	foreign port number

Figure 24.7 Variables in TCP listener table: tcpTable.

The first PCB variable (`t_state`) is from the TCP control block (Figure 24.13) and the remaining four are from the Internet PCB (Figure 22.4).

24.3 TCP protosw Structure

Figure 24.8 lists the TCP protosw structure, the protocol switch entry for TCP.

Member	inetow[2]	Description
pr_type	SOCK_STREAM	TCP provides a byte-stream service
pr_domain	&inetdomain	TCP is part of the Internet domain
pr_protocol	IPPROTO_TCP (6)	appears in the ip_p field of the IP header
pr_flags	PR_CONMREQUIRED/PR_WANTRCVD	socket layer flags, not used by protocol processing
pr_input	tcp_input	receives messages from IP layer
pr_output	0	not used by TCP
pr_ctlinput	tcp_ctlinput	control input function for ICMP errors
pr_ctloutput	tcp_ctloutput	respond to administrative requests from a process
pr_usrreq	tcp_usrreq	respond to communication requests from a process
pr_init	tcp_init	initialization for TCP
pr_fasttimo	tcp_fasttimo	fast timeout function, called every 200 ms
pr_slowtimo	tcp_slowtimo	slow timeout function, called every 500 ms
pr_drain	tcp_drain	called when kernel runs out of mbufs
pr_sysctl	0	not used by TCP

Figure 24.8 The TCP protosw structure.

24.4 TCP Header

The TCP header is defined as a `tcphdr` structure. Figure 24.9 shows the C structure and Figure 24.10 shows a picture of the TCP header.

```

40 struct tcphdr {
41     u_short th_sport;          /* source port */
42     u_short th_dport;         /* destination port */
43     tcp_seq th_seq;           /* sequence number */
44     tcp_seq th_ack;           /* acknowledgement number */
45 #if BYTE_ORDER == LITTLE_ENDIAN
46     u_char  th_x2:4;           /* (unused) */
47     th_off:4;                 /* data offset */
48 #endif
49 #if BYTE_ORDER == BIG_ENDIAN
50     u_char  th_off:4;         /* data offset */
51     th_x2:4;                 /* (unused) */
52 #endif
53     u_char  th_flags;         /* ACK, FIN, PUSH, RST, SYN, URG */
54     u_short th_win;           /* advertised window */
55     u_short th_sum;           /* checksum */
56     u_short th_urp;           /* urgent offset */
57 };

```

Figure 24.9 `tcphdr` structure.

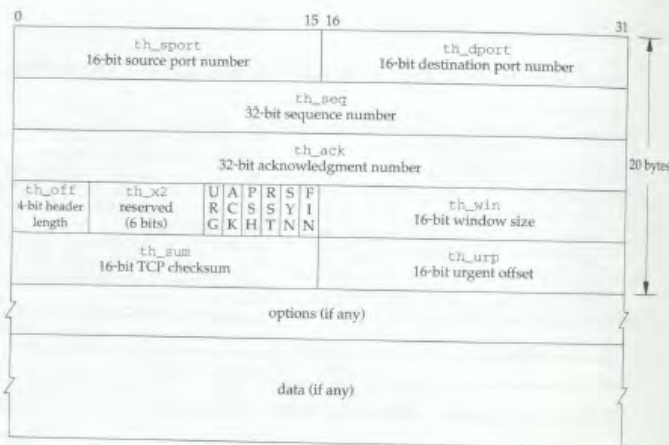


Figure 24.10 TCP header and optional data.

Most RFCs, most books (including Volume 1), and the code we'll examine call `th_urp` the *urgent pointer*. A better term is the *urgent offset*, since this field is a 16-bit unsigned offset that must be added to the sequence number field (`th_seq`) to give the 32-bit sequence number of the last byte of urgent data. (There is a continuing debate over whether this sequence number points to the last byte of urgent data or to the byte that follows. This is immaterial for the present discussion.) We'll see in Figure 24.13 that TCP correctly calls the 32-bit sequence number of the last byte of urgent data *send_urgent_pointer*. But using the term *pointer* for the 16-bit offset in the TCP header is misleading. In Exercise 26.6 we'll reiterate the distinction between the urgent pointer and the urgent offset.

The 4-bit header length, the 6 reserved bits that follow, followed by 8 bits of flags. To handle the difference in the order of these 4-bit fields within an 8-bit byte, the code contains an `#ifdef` based on the byte order of the system.

Also notice that we call the 4-bit `th_off` the *header length*, while the C code calls it the *data offset*. Both are correct since it is the length of the TCP header, including options, in 32-bit words, which is the offset of the first byte of data.

The `th_flags` member contains 6 flag bits, accessed using the names in Figure 24.11.

In Net/3 the TCP header is normally referenced as an IP header immediately followed by a TCP header. This is how `tcp_input` processes received IP datagrams and how `tcp_output` builds outgoing IP datagrams. This combined IP/TCP header is a `tciphdr` structure, shown in Figure 24.12.

th_flags	Description
TH_ACK	the acknowledgment number (th_ack) is valid
TH_FIN	the sender is finished sending data
TH_PUSH	receiver should pass the data to application without delay
TH_RST	reset the connection
TH_SYN	synchronize sequence numbers (establish connection)
TH_URG	the urgent offset (th_urg) is valid

Figure 24.11 th_flags values.

```

38 struct tcphdr {
39     struct ipovly ti_i;      /* overlaid ip structure */
40     struct tcphdr ti_t;    /* tcp header */
41 };

42 #define ti_next      ti_i.ih_next
43 #define ti_prev      ti_i.ih_prev
44 #define ti_x1        ti_i.ih_x1
45 #define ti_pr        ti_i.ih_pr
46 #define ti_len        ti_i.ih_len
47 #define ti_src        ti_i.ih_src
48 #define ti_dst        ti_i.ih_dst
49 #define ti_sport      ti_t.th_sport
50 #define ti_dport      ti_t.th_dport
51 #define ti_seq        ti_t.th_seq
52 #define ti_ack        ti_t.th_ack
53 #define ti_x2        ti_t.th_x2
54 #define ti_off        ti_t.th_off
55 #define ti_flags      ti_t.th_flags
56 #define ti_win        ti_t.th_win
57 #define ti_sum        ti_t.th_sum
58 #define ti_urg        ti_t.th_urg

```

Figure 24.12 tcphdr structure: combined IP/TCP header.

38-58 The 20-byte IP header is defined as an ipovly structure, which we showed earlier in Figure 23.12. As we discussed with Figure 23.19, this structure is not a real IP header, although the lengths are the same (20 bytes).

24.5 TCP Control Block

In Figure 22.1 we showed that TCP maintains its own control block, a `tcpcb` structure, in addition to the standard Internet PCB. In contrast, UDP has everything it needs in the Internet PCB—it doesn't need its own control block.

The TCP control block is a large structure, occupying 140 bytes. As shown in Figure 22.1 there is a one-to-one relationship between the Internet PCB and the TCP control block, and each points to the other. Figure 24.13 shows the definition of the TCP control block.

```

41 struct tcpcb {
42     struct tcpiphdr *seg_next; /* reassembly queue of received segments */
43     struct tcpiphdr *seg_prev; /* reassembly queue of received segments */
44     short t_state; /* connection state (Figure 24.16) */
45     short t_timer[TCPT_NTIMERS]; /* tcp timers (Chapter 25) */
46     short t_rxtshift; /* log(2) of rextm exp. backoff */
47     short t_rxtcur; /* current retransmission timeout (#ticks) */
48     short t_dupacks; /* #consecutive duplicate ACKs received */
49     u_short t_maxseg; /* maximum segment size to send */
50     char t_force; /* 1 if forcing out a byte (persist/OOB) */
51     u_short t_flags; /* (Figure 24.14) */
52     struct tcpiphdr *t_template; /* skeletal packet for transmit */
53     struct inpcb *t_inpcb; /* back pointer to internet PCB */
54 } /*
55  * The following fields are used as in the protocol specification.
56  * See RFC783, Dec. 1981, page 21.
57  */
58 /* send sequence variables */
59 tcp_seq snd_una; /* send unacknowledged */
60 tcp_seq snd_nxt; /* send next */
61 tcp_seq snd_up; /* send urgent pointer */
62 tcp_seq snd_wl1; /* window update seg seq number */
63 tcp_seq snd_wl2; /* window update seg ack number */
64 tcp_seq iss; /* initial send sequence number */
65 u_long snd_wnd; /* send window */
66 /* receive sequence variables */
67 u_long rcv_wnd; /* receive window */
68 tcp_seq rcv_nxt; /* receive next */
69 tcp_seq rcv_up; /* receive urgent pointer */
70 tcp_seq irs; /* initial receive sequence number */
71 /*
72  * Additional variables for this implementation.
73  */
74 /* receive variables */
75 tcp_seq rcv_adv; /* advertised window by other end */
76 /* retransmit variables */
77 tcp_seq snd_max; /* highest sequence number sent;
78  * used to recognize retransmits */
79 /* congestion control (slow start, source quench, retransmit after loss) */
80 u_long snd_cwnd; /* congestion-controlled window */
81 u_long snd_ssthresh; /* snd_cwnd size threshold for slow start
82  * exponential to linear switch */
83 /*
84  * transmit timing stuff. See below for scale of srtt and rttvar.
85  * "Variance" is actually smoothed difference.
86  */
87 short t_idle; /* inactivity time */
88 short t_rtt; /* round-trip time */
89 tcp_seq t_rtseq; /* sequence number being timed */
90 short t_srtt; /* smoothed round-trip time */
91 short t_rttvar; /* variance in round-trip time */
92 u_short t_rttmin; /* minimum rtt allowed */
93 u_long max_sndwnd; /* largest window peer has offered */

```

```

94 /* out-of-band data */
95 char   t_oobflags;           /* TCPOOB_HAVEDATA, TCPOOB_HADDATA */
96 char   t_ioobc;             /* input character, if not SO_OOBLINK */
97 short  t_softerror;         /* possible error not yet reported */
98 /* RFC 1323 variables */
99 u_char  snd_scale;          /* scaling for send window [0-14] */
100 u_char  rcv_scale;          /* scaling for receive window [0-14] */
101 u_char  request_r_scale;    /* our pending window scale */
102 u_char  requested_s_scale;  /* peer's pending window scale */
103 u_long  ts_recent;          /* timestamp echo data */
104 u_long  ts_recent_age;      /* when last updated */
105 tcp_seq last_ack_sent;      /* sequence number of last ack field */
106 };
107 #define intotcpb(ip) ((struct tcpb *) (ip)->inp_ppcb)
108 #define sototcpb(so) (intotcpb(sotoinpcb(so)))

```

tcp_var.h

Figure 24.13 tcpb structure: TCP control block.

We'll save the discussion of these variables until we encounter them in the code.

Figure 24.14 shows the values for the `t_flags` member.

<code>t_flags</code>	Description
<code>TF_ACKNOW</code>	send ACK immediately
<code>TF_DELACK</code>	send ACK, but try to delay it
<code>TF_NODELAY</code>	don't delay packets to coalesce (disable Nagle algorithm)
<code>TF_NOOPT</code>	don't use TCP options (never set)
<code>TF_SENTFIN</code>	have sent FIN
<code>TF_RCVD_SCALE</code>	set when other side sends window scale option in SYN
<code>TF_RCVD_TSTMP</code>	set when other side sends timestamp option in SYN
<code>TF_REQ_SCALE</code>	have/will request window scale option in SYN
<code>TF_REQ_TSTMP</code>	have/will request timestamp option in SYN

Figure 24.14 `t_flags` values.

24.6 TCP State Transition Diagram

Many of TCP's actions, in response to different types of segments arriving on a connection, can be summarized in a state transition diagram, shown in Figure 24.15. We also duplicate this diagram on one of the front end papers, for easy reference while reading the TCP chapters.

These state transitions define the TCP finite state machine. Although the transition from LISTEN to SYN_SENT is allowed by TCP, there is no way to do this using the sockets API (i.e., a connect is not allowed after a listen).

The `t_state` member of the control block holds the current state of a connection, with the values shown in Figure 24.16.

This figure also shows the `tcp_outflags` array, which contains the outgoing flags for `tcp_output` to use when the connection is in that state.

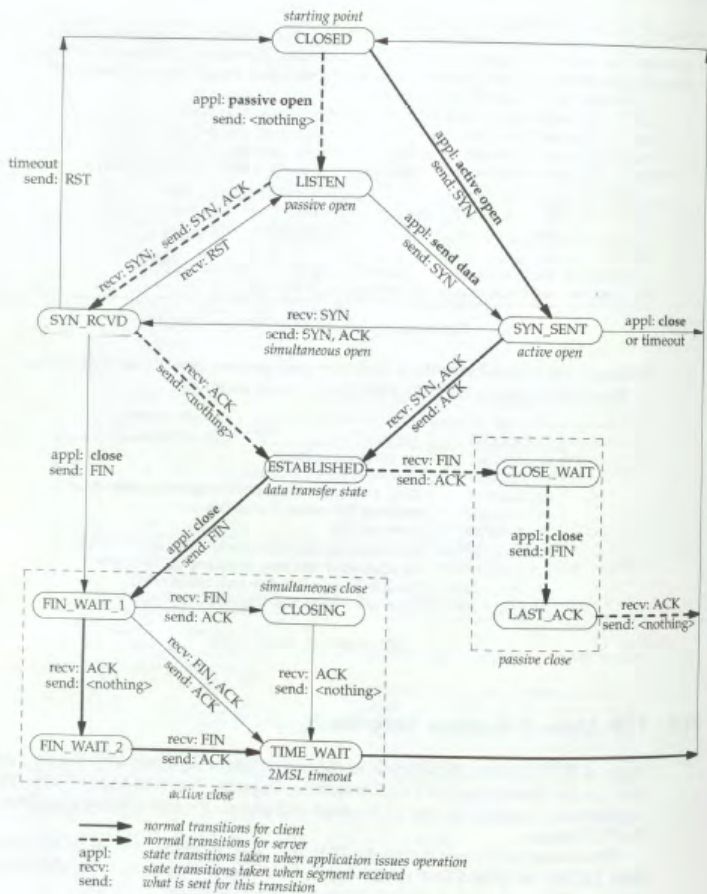


Figure 24.15 TCP state transition diagram.

<code>t_state</code>	value	Description	<code>tcp_outflags()</code>
<code>TCP_CLOSED</code>	0	closed	<code>TH_RST TH_ACK</code>
<code>TCP_LISTEN</code>	1	listening for connection (passive open)	0
<code>TCP_SYN_SENT</code>	2	have sent SYN (active open)	<code>TH_SYN</code>
<code>TCP_SYN_RECEIVED</code>	3	have sent and received SYN; awaiting ACK	<code>TH_SYN TH_ACK</code>
<code>TCP_ESTABLISHED</code>	4	established (data transfer)	<code>TH_ACK</code>
<code>TCP_CLOSE_WAIT</code>	5	received FIN, waiting for application close	<code>TH_ACK</code>
<code>TCP_FIN_WAIT_1</code>	6	have closed, sent FIN; awaiting ACK and FIN	<code>TH_FIN TH_ACK</code>
<code>TCP_CLOSING</code>	7	simultaneous close; awaiting ACK	<code>TH_FIN TH_ACK</code>
<code>TCP_LAST_ACK</code>	8	received FIN have closed; awaiting ACK	<code>TH_FIN TH_ACK</code>
<code>TCP_FIN_WAIT_2</code>	9	have closed; awaiting FIN	<code>TH_ACK</code>
<code>TCP_TIME_WAIT</code>	10	2MSL wait state after active close	<code>TH_ACK</code>

Figure 24.16 `t_state` values.

Figure 24.16 also shows the numerical values of these constants since the code uses their numerical relationships. For example, the following two macros are defined:

```
#define TCPS_HAVERCVDSYN(s) ((s) >= TCPS_SYN_RECEIVED)
#define TCPS_HAVERCVDFIN(s) ((s) >= TCPS_TIME_WAIT)
```

Similarly, we'll see that `tcp_notify` handles ICMP errors differently when the connection is not yet established, that is, when `t_state` is less than `TCPS_ESTABLISHED`.

The name `TCPS_HAVERCVDSYN` is correct, but the name `TCPS_HAVERCVDFIN` is misleading. A FIN has also been received in the `CLOSE_WAIT`, `CLOSING`, and `LAST_ACK` states. We encounter this macro in Chapter 29.

Half-Close

When a process calls `shutdown` with a second argument of 1, it is called a *half-close*. TCP sends a FIN but allows the process to continue receiving on the socket. (Section 18.5 of Volume 1 contains examples of TCP's half-close.)

For example, even though we label the `ESTABLISHED` state "data transfer," if the process does a half-close, moving the connection to the `FIN_WAIT_1` and then the `FIN_WAIT_2` states, data can continue to be received by the process in these two states.

24.7 TCP Sequence Numbers

Every byte of data exchanged across a TCP connection, along with the SYN and FIN flags, is assigned a 32-bit *sequence number*. The sequence number field in the TCP header (Figure 24.10) contains the sequence number of the first byte of data in the segment. The *acknowledgment number* field in the TCP header contains the next sequence number that the sender of the ACK expects to receive, which acknowledges all data bytes through the acknowledgment number minus 1. In other words, the acknowledgment number is the *next* sequence number expected by the sender of the ACK. The acknowledgment number is valid only if the ACK flag is set in the header. We'll see

that TCP always sets the ACK flag except for the first SYN sent by an active open (the SYN_SENT state; see `tcp_outflags[2]` in Figure 24.16) and in some RST segments.

Since a TCP connection is *full-duplex*, each end must maintain a set of sequence numbers for both directions of data flow. In the TCP control block (Figure 24.13) there are 13 sequence numbers: eight for the send direction (the *send sequence space*) and five for the receive direction (the *receive sequence space*).

Figure 24.17 shows the relationship of four of the variables in the send sequence space: `snd_wnd`, `snd_una`, `snd_nxt`, and `snd_max`. In this example we number the bytes 1 through 11.

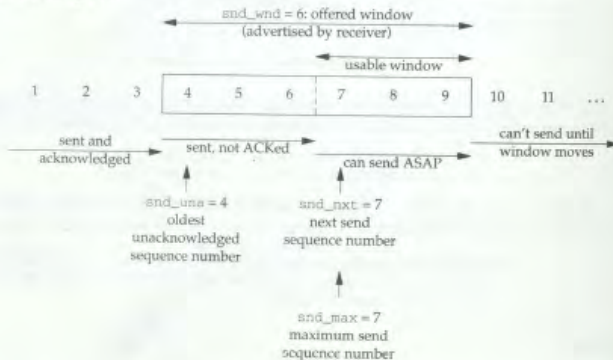


Figure 24.17 Example of send sequence space.

An *acceptable ACK* is one for which the following inequality holds:

$$\text{snd_una} < \text{acknowledgment field} \leq \text{snd_max}$$

In Figure 24.17 an acceptable ACK has an acknowledgment field of 5, 6, or 7. An acknowledgment field less than or equal to `snd_una` is a duplicate ACK—it acknowledges data that has already been ACKed, or else `snd_una` would not have incremented past those bytes.

We encounter the following test a few times in `tcp_output`, which is true if a segment is being retransmitted:

$$\text{snd_nxt} < \text{snd_max}$$

Figure 24.18 shows the other end of the connection in Figure 24.17: the receive sequence space, assuming the segment containing sequence numbers 4, 5, and 6 has not been received yet. We show the three variables `rcv_nxt`, `rcv_wnd`, and `rcv_adv`.

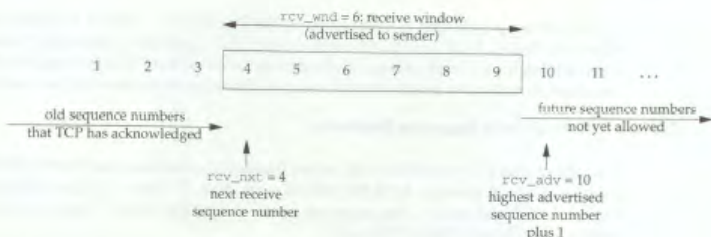


Figure 24.18 Example of receive sequence space.

The receiver considers a received segment valid if it contains data within the window, that is, if either of the following two inequalities is true:

$$rcv_nxt \leq \text{beginning sequence number of segment} < rcv_nxt + rcv_wnd$$

$$rcv_nxt \leq \text{ending sequence number of segment} < rcv_nxt + rcv_wnd$$

The beginning sequence number of a segment is just the sequence number field in the TCP header, ti_seq . The ending sequence number is the sequence number field plus the number of bytes of TCP data, minus 1.

For example, Figure 24.19 could represent the TCP segment containing the 3 bytes with sequence numbers 4, 5, and 6 in Figure 24.17.

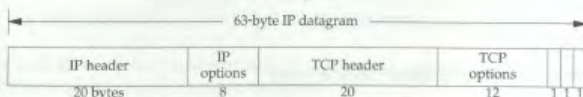


Figure 24.19 TCP segment transmitted as an IP datagram.

We assume that there are 8 bytes of IP options and 12 bytes of TCP options. Figure 24.20 shows the values of the relevant variables.

Variable	Value	Description
ip_hl	7	length of IP header + options in 32-bit words (= 28 bytes)
ip_len	63	length of IP datagram in bytes (20 + 8 + 20 + 12 + 3)
ti_off	8	length of TCP header + options in 32-bit words (= 32 bytes)
ti_seq	4	sequence number of first byte of data
ti_len	3	#bytes of TCP data: $ip_len - (ip_hl \times 4) - (ti_off \times 4)$
	6	sequence number of last byte of data: $ti_seq + ti_len - 1$

Figure 24.20 Values of variables corresponding to Figure 24.19.

`seq_len` is not a field that is transmitted in the TCP header. Instead, it is computed as shown in Figure 24.20 and stored in the overlaid IP structure (Figure 24.12) once the received header fields have been checksummed and verified. The last value in this figure is not stored in the header, but is computed from the other values when needed.

Modular Arithmetic with Sequence Numbers

A problem that TCP must deal with is that the sequence numbers are from a finite 32-bit number space: 0 through 4,294,967,295. If more than 2^{32} bytes of data are exchanged across a TCP connection, the sequence numbers will be reused. Sequence numbers wrap around from 4,294,967,295 to 0.

Even if less than 2^{32} bytes of data are exchanged, wrap around is still a problem because the sequence numbers for a connection don't necessarily start at 0. The initial sequence number for each direction of data flow across a connection can start anywhere between 0 and 4,294,967,295. This complicates the comparison of sequence numbers. For example, sequence number 1 is "greater than" 4,294,967,295, as we discuss below.

TCP sequence numbers are defined as unsigned longs in `tcp.h`:

```
typedef unsigned long tcp_seq;
```

The four macros shown in Figure 24.21 compare sequence numbers.

```

40 #define SEQ_LT(a,b)      ((int)((a)-(b)) < 0)
41 #define SEQ_LEQ(a,b)    ((int)((a)-(b)) <= 0)
42 #define SEQ_GT(a,b)    ((int)((a)-(b)) > 0)
43 #define SEQ_GEQ(a,b)    ((int)((a)-(b)) >= 0)

```

Figure 24.21 Macros for TCP sequence number comparison.

Example—Sequence Number Comparisons

Let's look at an example to see how TCP's sequence numbers operate. Assume 3-bit sequence numbers, 0 through 7. Figure 24.22 shows these eight sequence numbers, their 3-bit binary representation, and their two's complement representation. (To form the two's complement take the binary number, convert each 0 to a 1 and vice versa, then add 1.) We show the two's complement because to form $a - b$ we just add a to the two's complement of b .

The final three columns of this table are 0 minus x , 1 minus x , and 2 minus x . In these final three columns, if the value is considered to be a signed integer (notice the cast to `int` in all four macros in Figure 24.21), the value is less than 0 (the `SEQ_LT` macro) if the high-order bit is 1, and the value is greater than 0 (the `SEQ_GT` macro) if the high-order bit is 0 and the value is not 0. We show horizontal lines in these final three columns to distinguish between the four negative and the four nonnegative values.

If we look at the fourth column of Figure 24.22, (labeled " $0 - x$ "), we see that 0 (i.e., x), is less than 1, 2, 3, and 4 (the high-order bit of the result is 1), and 0 is greater than 5, 6, and 7 (the high-order bit is 0 and the result is not 0). We show this relationship pictorially in Figure 24.23.

x	binary	two's complement	0-x	1-x	2-x
0	000	000	000	001	010
1	001	111	111	000	001
2	010	110	110	111	000
3	011	101	101	110	111
4	100	100	100	101	110
5	101	011	011	100	101
6	110	010	010	011	100
7	111	001	001	010	011

Figure 24.22 Example using 3-bit sequence numbers.

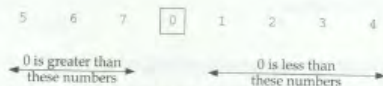


Figure 24.23 TCP sequence number comparisons for 3-bit sequence numbers.

Figure 24.24 shows a similar figure using the fifth row of the table (1-x).

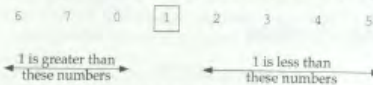


Figure 24.24 TCP sequence number comparisons for 3-bit sequence numbers.

Figure 24.25 is another representation of the two previous figures, using circles to reiterate the wrap around of sequence numbers.



Figure 24.25 Another way to visualize Figures 24.23 and 24.24.

With regard to TCP, these sequence number comparisons determine whether a given sequence number is in the future or in the past (a retransmission). For example, using Figure 24.24, if TCP is expecting sequence number 1 and sequence number 6 arrives, since 6 is less than 1 using the sequence number arithmetic we showed, the data byte is considered a retransmission of a previously received data byte and is discarded. But if sequence number 5 is received, since it is greater than 1 it is considered a future

data byte and is saved by TCP, awaiting the arrival of the missing bytes 2, 3, and 4 (assuming byte 5 is within the receive window).

Figure 24.26 is an expansion of the left circle in Figure 24.25, using TCP's 32-bit sequence numbers instead of 3-bit sequence numbers.



Figure 24.26 Comparisons against 0, using 32-bit sequence numbers.

The right circle in Figure 24.26 is to reiterate that one-half of the 32-bit sequence space uses 2^{31} numbers.

24.8 tcp_init Function

The `domaininit` function calls TCP's initialization function, `tcp_init` (Figure 24.27), at system initialization time.

```

43 void
44 tcp_init()
45 {
46     tcp_iss = 1;
47     tcb.inp_next = tcb.inp_prev = &tcb;
48     if (max_protohdr < sizeof(struct tcpiphdr))
49         max_protohdr = sizeof(struct tcpiphdr);
50     if (max_linkhdr + sizeof(struct tcpiphdr) > MHLLEN)
51         panic("tcp_init");
52 }

```

tcp_subr.c

tcp_subr.c

Figure 24.27 `tcp_init` function.

Set initial send sequence number (ISS)

46 The initial send sequence number (ISS), `tcp_iss`, is initialized to 1. As the comment indicates, this is wrong. We discuss the implications behind this choice shortly, when we describe TCP's *quiet time*. Compare this to the initialization of the IP identifier in Figure 7.23, which used the time-of-day clock.

Initialize linked list of TCP Internet PCBs

47. The next and previous pointers in the head PCB (`tcbb`) point to itself. This is an empty doubly linked list. The remainder of the `tcbb` PCB is initialized to 0 (all uninitialized globals are set to 0), although the only other field used in this head PCB is `inp_lport`, the next TCP ephemeral port number to allocate. The first ephemeral port used by TCP will be 1024, for the reasons described in the solution for Exercise 22.4.

Calculate maximum protocol header length

- 48-51. If the maximum protocol header encountered so far is less than 40 bytes, `max_protohdr` is set to 40 (the size of the combined IP and TCP headers, without any options). This variable is described in Figure 7.17. If the sum of `max_linkhdr` (normally 16) and 40 is greater than the amount of data that fits into an mbuf with a packet header (100 bytes, `MHLEN` from Figure 2.7), the kernel panics (Exercise 24.2).

MSL and Quiet Time Concept

TCP requires any host that crashes without retaining any knowledge of the last sequence numbers used on active connections to refrain from sending any TCP segments for one MSL (2 minutes, the quiet time) on reboot. Few TCPs, if any, retain this knowledge over a crash or operator shutdown.

MSL is the *maximum segment lifetime*. Each implementation chooses a value for the MSL. It is the maximum amount of time any segment can exist in the network before being discarded. A connection that is actively closed remains in the `CLOSE_WAIT` state (Figure 24.15) for twice the MSL.

RFC 793 [Postel 1981c] recommends an MSL of 2 minutes, but Net/3 uses an MSL of 30 seconds (the constant `TCPTV_MSL` in Figure 25.3).

The problem occurs if packets are delayed somewhere in the network (RFC 793 calls these *wandering duplicates*). Assume a Net/3 system starts up, initializes `tcp_iss` to 1 (as in Figure 24.27) and then crashes just after the sequence numbers wrap. We'll see in Section 25.5 that TCP increments `tcp_iss` by 128,000 every second, causing the wrap around of the ISS to occur about 9.3 hours after rebooting. Also, `tcp_iss` is incremented by 64,000 each time a `connect` is issued, which can cause the wrap around to occur earlier than 9.3 hours. The following scenario is one example of how an old segment can incorrectly be delivered to a connection:

1. A client and server have an established connection. The client's port number is 1024. The client sends a data segment with a starting sequence number of 2. This data segment gets trapped in a routing loop somewhere between the two end points and is not delivered to the server. This data segment becomes a wandering duplicate.
2. The client retransmits the data segment starting with sequence number 2, which is delivered to the server.
3. The client closes the connection.

4. The client host crashes.
5. The client host reboots about 40 seconds after crashing, causing TCP to initialize `tcp_iss` to 1 again.
6. Another connection is immediately established by the same client to the same server, using the same socket pair: the client uses 1024 again, and the server uses its well-known port. The client's SYN uses sequence number 1. This new connection using the same socket pair is called a new *incarnation* of the old connection.
7. The wandering duplicate from step 1 is delivered to the server, and it thinks this datagram belongs to the new connection, when it is really from the old connection.

Figure 24.28 is a time line of this sequence of steps.

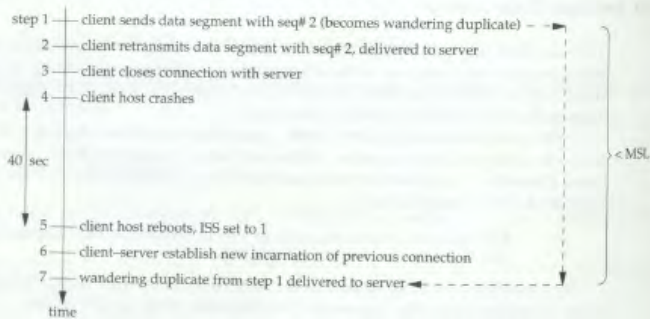


Figure 24.28 Example of old segment delivered to new incarnation of a connection.

This problem exists even if the rebooting TCP were to use an algorithm based on its time-of-day clock to choose the ISS on rebooting: regardless of the ISS for the previous incarnation of a connection, because of sequence number wrap it is possible for the ISS after rebooting to nearly equal the sequence number in use before the reboot.

Besides saving the sequence number of all established connections, the only other way around this problem is for the rebooting TCP to be quiet (i.e., not send any TCP segments) for MSL seconds after crashing. Few TCPs do this, however, since it takes most hosts longer than MSL seconds just to reboot.

24.9 Summary

This chapter is an introduction to the TCP source code in the six chapters that follow. TCP maintains its own control block for each connection, containing all the variable and state information for the connection.

A state transition diagram is defined for TCP that shows under what conditions TCP moves from one state to another and what segments get sent by TCP for each transition. This diagram shows how connections are established and terminated. We'll refer to this state transition diagram frequently in our description of TCP.

Every byte exchanged across a TCP connection has an associated sequence number, and TCP maintains numerous sequence numbers in the connection control block: some for sending and some for receiving (since TCP is full-duplex). Since these sequence numbers are from a finite 32-bit sequence space, they wrap around from the maximum value back to 0. We explained how the sequence numbers are compared to each other using less-than and greater-than tests, which we'll encounter repeatedly in the TCP code.

Finally, we looked at one of the simplest of the TCP functions, `tcp_init`, which initializes TCP's linked list of Internet PCBs. We also discussed TCP's choice of an initial send sequence number, which is used when actively opening a connection.

Exercises

- 24.1 What is the average number of bytes transmitted and received per connection from the statistics in Figure 24.5?
- 24.2 Is the kernel panic in `tcp_init` reasonable?
- 24.3 Execute `netstat -a` to see how many TCP end points your system currently has active.

Case Summary

The following is a summary of the case. The case involves a dispute between Cavium, Inc. and Alacritech, Inc. regarding the use of certain technology. Cavium, Inc. claims that Alacritech, Inc. has infringed on its patents. Alacritech, Inc. denies the allegations and claims that its technology is prior art. The court has granted summary judgment in favor of Cavium, Inc., finding that Alacritech, Inc. has infringed on Cavium, Inc.'s patents. The court also awarded damages to Cavium, Inc. and granted its request for attorneys' fees and costs.

Conclusion

In conclusion, the court has ruled in favor of Cavium, Inc. and awarded damages and attorneys' fees. This decision is a significant victory for Cavium, Inc. and a setback for Alacritech, Inc. The court's ruling is based on the fact that Alacritech, Inc. cannot prove that its technology is prior art. The court also found that Alacritech, Inc. has infringed on Cavium, Inc.'s patents. This decision is a clear statement of the law and will serve as a warning to other companies in the industry. Cavium, Inc. is pleased with the court's decision and will continue to protect its intellectual property.

TCP Timers

25.1 Introduction

We start our detailed description of the TCP source code by looking at the various TCP timers. We encounter these timers throughout most of the TCP functions.

TCP maintains seven timers for *each* connection. They are briefly described here, in the approximate order of their occurrence during the lifetime of a connection.

1. A *connection-establishment* timer starts when a SYN is sent to establish a new connection. If a response is not received within 75 seconds, the connection establishment is aborted.
2. A *retransmission* timer is set when TCP sends data. If the data is not acknowledged by the other end when this timer expires, TCP retransmits the data. The value of this timer (i.e., the amount of time TCP waits for an acknowledgment) is calculated dynamically, based on the round-trip time measured by TCP for this connection, and based on the number of times this data segment has been retransmitted. The retransmission timer is bounded by TCP to be between 1 and 64 seconds.
3. A *delayed ACK* timer is set when TCP receives data that must be acknowledged, but need not be acknowledged immediately. Instead, TCP waits up to 200 ms before sending the ACK. If, during this 200-ms time period, TCP has data to send on this connection, the pending acknowledgment is sent along with the data (called *piggybacking*).

4. A *persist* timer is set when the other end of a connection advertises a window of 0, stopping TCP from sending data. Since window advertisements from the other end are not sent reliably (that is, ACKs are not acknowledged, only data is acknowledged), there's a chance that a future window update, allowing TCP to send some data, can be lost. Therefore, if TCP has data to send and the other end advertises a window of 0, the persist timer is set and when it expires, 1 byte of data is sent to see if the window has opened. Like the retransmission timer, the persist timer value is calculated dynamically, based on the round-trip time. The value of this is bounded by TCP to be between 5 and 60 seconds.
5. A *keepalive* timer can be set by the process using the `SO_KEEPAIVE` socket option. If the connection is idle for 2 hours, the keepalive timer expires and a special segment is sent to the other end, forcing it to respond. If the expected response is received, TCP knows that the other host is still up, and TCP won't probe it again until the connection is idle for another 2 hours. Other responses to the keepalive probe tell TCP that the other host has crashed and rebooted. If no response is received to a fixed number of keepalive probes, TCP assumes that the other end has crashed, although it can't distinguish between the other end being down (i.e., it crashed and has not yet rebooted) and a temporary lack of connectivity to the other end (i.e., an intermediate router or phone line is down).
6. A `FIN_WAIT_2` timer. When a connection moves from the `FIN_WAIT_1` state to the `FIN_WAIT_2` state (Figure 24.15) and the connection cannot receive any more data (implying the process called `close`, instead of taking advantage of TCP's half-close with `shutdown`), this timer is set to 10 minutes. When this timer expires it is reset to 75 seconds, and when it expires the second time the connection is dropped. The purpose of this timer is to avoid leaving a connection in the `FIN_WAIT_2` state forever, if the other end never sends a FIN. (We don't show this timeout in Figure 24.15.)
7. A `TIME_WAIT` timer, often called the *2MSL* timer. The term *2MSL* means twice the MSL, the maximum segment lifetime defined in Section 24.8. It is set when a connection enters the `TIME_WAIT` state (Figure 24.15), that is, when the connection is actively closed. Section 18.6 of Volume 1 describes the reasoning for the 2MSL wait state in detail. The timer is set to 1 minute (Net/3 uses an MSL of 30 seconds) when the connection enters the `TIME_WAIT` state and when it expires, the TCP control block and Internet PCB are deleted, allowing that socket pair to be reused.

TCP has two timer functions: one is called every 200 ms (the fast timer) and the other every 500 ms (the slow timer). The delayed ACK timer is different from the other six: when the delayed ACK timer is set for a connection it means that a delayed ACK must be sent the next time the 200-ms timer expires (i.e., the elapsed time is between 0 and 200 ms). The other six timers are decremented every 500 ms, and only when the counter reaches 0 does the corresponding action take place.

25.2 Code Introduction

The delayed ACK timer is enabled for a connection when the `TF_DELACK` flag (Figure 24.14) is set in the TCP control block. The array `t_timer` in the TCP control block contains four (`TCPT_NTIMERS`) counters used to implement the other six timers. The indexes into this array are shown in Figure 25.1. We describe briefly how the six timers (other than the delayed ACK timer) are implemented by these four counters.

Constant	Value	Description
<code>TCPT_REXMT</code>	0	retransmission timer
<code>TCPT_PERSIST</code>	1	persist timer
<code>TCPT_KEEP</code>	2	keepalive timer <i>or</i> connection-establishment timer
<code>TCPT_2MSL</code>	3	2MSL timer <i>or</i> <code>FIN_WAIT_2</code> timer

Figure 25.1 Indexes into the `t_timer` array.

Each entry in the `t_timer` array contains the number of 500-ms clock ticks until the timer expires, with 0 meaning that the timer is not set. Since each timer is a `short`, if 16 bits hold a `short`, the maximum timer value is 16,383.5 seconds, or about 4.5 hours.

Notice in Figure 25.1 that four “timer counters” implement six TCP “timers,” because some of the timers are mutually exclusive. We’ll distinguish between the counters and the timers. The `TCPT_KEEP` counter implements both the keepalive timer and the connection-establishment timer, since the two timers are never used at the same time for a connection. Similarly, the 2MSL timer and the `FIN_WAIT_2` timer are implemented using the `TCPT_2MSL` counter, since a connection is only in one state at a time. The first section of Figure 25.2 summarizes the implementation of the seven TCP timers. The second and third sections of the table show how four of the seven timers are initialized using three global variables from Figure 24.3 and two constants from Figure 25.3. Notice that two of the three globals are used with multiple timers. We’ve already said that the delayed ACK timer is tied to TCP’s 200-ms timer, and we describe how the other two timers are set later in this chapter.

	conn. estab.	rexmit	delayed ACK	persist	keep-alive	FIN_WAIT_2	2MSL
<code>t_timer[TCPT_REXMT]</code>		*					
<code>t_timer[TCPT_PERSIST]</code>				*			
<code>t_timer[TCPT_KEEP]</code>	*				*		
<code>t_timer[TCPT_2MSL]</code>						*	*
<code>t_flags & TF_DELACK</code>			*				
<code>tcp_keepidle (2 hr)</code>					*		
<code>tcp_keepintvl (75 sec)</code>					*	*	
<code>tcp_maxidle (10 min)</code>					*	*	
<code>2 * TCPTV_MSL (60 sec)</code>							*
<code>TCPTV_KEEP_INIT (75 sec)</code>	*						

Figure 25.2 Implementation of the seven TCP timers.

Figure 25.3 shows the fundamental timer values for the Net/3 implementation.

Constant	#500-ms clock ticks	#sec	Description
<i>TCPTV_MSL</i>	60	30	MSL, maximum segment lifetime
<i>TCPTV_MIN</i>	2	1	minimum value of retransmission timer
<i>TCPTV_REXMTMAX</i>	128	64	maximum value of retransmission timer
<i>TCPTV_PERSMIN</i>	10	5	minimum value of persist timer
<i>TCPTV_PERSMAX</i>	120	60	maximum value of persist timer
<i>TCPTV_KEEP_INIT</i>	150	75	connection-establishment timer value
<i>TCPTV_KEEP_IDLE</i>	14400	7200	idle time for connection before first probe (2 hours)
<i>TCPTV_KEEPIPTVL</i>	150	75	time between probes when no response
<i>TCPTV_SRTTBASE</i>	0		special value to denote no measurements yet for connection
<i>TCPTV_SRTTDFLT</i>	6	3	default RTT when no measurements yet for connection

Figure 25.3 Fundamental timer values for the implementation.

Figure 25.4 shows other timer constants that we'll encounter.

Constant	Value	Description
<i>TCP_LINGERTIME</i>	120	maximum #seconds for <i>SO_LINGER</i> socket option
<i>TCP_MAXRTSHIFT</i>	12	maximum #retransmissions waiting for an ACK
<i>TCPTV_KEEPCNT</i>	8	maximum #keepalive probes when no response received

Figure 25.4 Timer constants.

The *TCPT_RANGESET* macro, shown in Figure 25.5, sets a timer to a given value, making certain the value is between the specified minimum and maximum.

```

102 #define TCPT_RANGESET(tv, value, tvmin, tvmax) ( \
103     (tv) = (value); \
104     if ((tv) < (tvmin)) \
105         (tv) = (tvmin); \
106     else if ((tv) > (tvmax)) \
107         (tv) = (tvmax); \
108 )

```

tcp_timer.h

tcp_timer.h

Figure 25.5 *TCPT_RANGESET* macro.

We see in Figure 25.3 that the retransmission timer and the persist timer have upper and lower bounds, since their values are calculated dynamically, based on the measured round-trip time. The other timers are set to constant values.

There is one additional timer that we allude to in Figure 25.4 but don't discuss in this chapter: the linger timer for a socket, set by the *SO_LINGER* socket option. This is a socket-level timer used by the *close* system call (Section 15.15). We will see in Figure 30.12 that when a socket is closed, TCP checks whether this socket option is set and whether the linger time is 0. If so, the connection is aborted with an RST instead of TCP's normal close.

25.3 tcp_canceltimers Function

The function `tcp_canceltimers`, shown in Figure 25.6, is called by `tcp_input` when the `TIME_WAIT` state is entered. All four timer counters are set to 0, which turns off the retransmission, persist, keepalive, and `FIN_WAIT_2` timers, before `tcp_input` sets the 2MSL timer.

```

107 void                                     tcp_timer.c
108 tcp_canceltimers(tp)
109 struct tcpcb *tp;
110 {
111     int    i;
112     for (i = 0; i < TCPT_NTIMERS; i++)
113         tp->t_timer[i] = 0;
114 }

```

Figure 25.6 `tcp_canceltimers` function.

25.4 tcp_fasttimo Function

The function `tcp_fasttimo`, shown in Figure 25.7, is called by `pr_fasttimo` every 200 ms. It handles only the delayed ACK timer.

```

41 void                                     tcp_timer.c
42 tcp_fasttimo()
43 {
44     struct inpcb *inp;
45     struct tcpcb *tp;
46     int    s = splnet();
47     inp = tcb.inp_next;
48     if (inp)
49         for (; inp != &tcb; inp = inp->inp_next)
50             if ((tp = (struct tcpcb *) inp->inp_ppcb) &&
51                 (tp->t_flags & TF_DELACK)) {
52                 tp->t_flags &= ~TF_DELACK;
53                 tp->t_flags |= TF_ACKNOW;
54                 tcbstat.tcpa_delack++;
55                 (void) tcp_output(tp);
56             }
57     splx(s);
58 }

```

Figure 25.7 `tcp_fasttimo` function, which is called every 200 ms.

Each Internet PCB on the TCP list that has a corresponding TCP control block is checked. If the `TF_DELACK` flag is set, it is cleared and the `TF_ACKNOW` flag is set instead. `tcp_output` is called, and since the `TF_ACKNOW` flag is set, an ACK is sent.

How can TCP have an Internet PCB on its PCB list that doesn't have a TCP control block (the test at line 50)? When a socket is created (the `PRU_ATTACH` request, in response to the `socket` system call) we'll see in Figure 30.11 that the creation of the Internet PCB is done first, followed by the creation of the TCP control block. Between these two operations a high-priority clock interrupt can occur (Figure 1.13), which calls `tcp_fasttimo`.

25.5 `tcp_slowtimo` Function

The function `tcp_slowtimo`, shown in Figure 25.8, is called by `pr_slowtimo` every 500 ms. It handles the other six TCP timers: connection establishment, retransmission, persist, keepalive, `FIN_WAIT_2`, and 2MSL.

71 `tcp_maxidle` is initialized to 10 minutes. This is the maximum amount of time TCP will send keepalive probes to another host, waiting for a response from that host. This variable is also used with the `FIN_WAIT_2` timer, as we describe in Section 25.6. This initialization statement could be moved to `tcp_init`, since it only needs to be evaluated when the system is initialized (see Exercise 25.2).

Check each timer counter in all TCP control blocks

72-85 Each Internet PCB on the TCP list that has a corresponding TCP control block is checked. Each of the four timer counters for each connection is tested, and if nonzero, the counter is decremented. When the timer reaches 0, a `PRU_SLOWTIMO` request is issued. We'll see that this request calls the function `tcp_timers`, which we describe later in this chapter.

The fourth argument to `tcp_usrreq` is a pointer to an mbuf. But this argument is actually used for different purposes when the mbuf pointer is not required. Here we see the index `i` is passed, telling the request which timer has expired. The funny-looking cast of `i` to an mbuf pointer is to avoid a compile-time error.

Check if TCP control block has been deleted

90-93 Before examining the timers for a control block, a pointer to the next Internet PCB is saved in `ipnxt`. Each time the `PRU_SLOWTIMO` request returns, `tcp_slowtimo` checks whether the next PCB in the TCP list still points to the PCB that's being processed. If not, it means the control block has been deleted—perhaps the 2MSL timer expired or the retransmission timer expired and TCP is giving up on this connection—causing a jump to `tpgone`, skipping the remaining timers for this control block, and moving on to the next PCB.

Count idle time

94 `t_idle` is incremented for the control block. This counts the number of 500-ms clock ticks since the last segment was received on this connection. It is set to 0 by `tcp_input` when a segment is received on the connection and used for three purposes: (1) by the keepalive algorithm to send a probe after the connection is idle for 2 hours, (2) to drop a connection in the `FIN_WAIT_2` state that is idle for 10 minutes and 75 seconds, and (3) by `tcp_output` to return to the slow start algorithm after the connection has been idle for a while.


```

64 void                                     tcp_timer.c
65 tcp_slowtimo()
66 {
67     struct inpcb *ip, *ipnxt;
68     struct tcpcb *tp;
69     int     s = splnet();
70     int     i;

71     tcp_maxidle = TCPTV_KEEPCNT * tcp_keepintvl;
72     /*
73      * Search through tcb's and update active timers.
74      */
75     ip = tcb.inp_next;
76     if (ip == 0) {
77         splx(s);
78         return;
79     }
80     for (; ip != &tcb; ip = ipnxt) {
81         ipnxt = ip->inp_next;
82         tp = intotcpcb(ip);
83         if (tp == 0)
84             continue;
85         for (i = 0; i < TCPT_TIMERS; i++) {
86             if (tp->t_timer[i] && --tp->t_timer[i] == 0) {
87                 (void) tcp_usrreq(tp->t_inpcb->inp_socket,
88                                     PRU_SLOWTIMO, (struct mbuf *) 0,
89                                     (struct mbuf *) i, (struct mbuf *) 0);
90                 if (ipnxt->inp_prev != ip)
91                     goto tpgone;
92             }
93         }
94         tp->t_idle++;
95         if (tp->t_rtt)
96             tp->t_rtt++;
97         tpgone:
98         ;
99     }
100     tcp_iss += TCP_ISSINCR / PR_SLOWHZ; /* increment iss */
101     tcp_now++; /* for timestamps */
102     splx(s);
103 }

```

Figure 25.8 tcp_slowtimo function, which is called every 500 ms.

Increment RTT counter

95-96

If this connection is timing an outstanding segment, `t_rtt` is nonzero and counts the number of 500-ms clock ticks until that segment is acknowledged. It is initialized to 1 by `tcp_output` when a segment is transmitted whose RTT should be timed. `tcp_slowtimo` increments this counter.

Increment initial send sequence number

100 `tcp_iss` was initialized to 1 by `tcp_init`. Every 500 ms it is incremented by 64,000: 128,000 (`TCP_ISSINCR`) divided by 2 (`PR_SLOWHZ`). This is a rate of about once every 8 microseconds, although `tcp_iss` is incremented only twice a second. We'll see that `tcp_iss` is also incremented by 64,000 each time a connection is established, either actively or passively.

RFC 793 specifies that the initial sequence number should increment roughly every 4 microseconds, or 250,000 times a second. The Net/3 value increments at about one-half this rate.

Increment RFC 1323 timestamp value

101 `tcp_now` is initialized to 0 on bootstrap and incremented every 500 ms. It is used by the timestamp option defined in RFC 1323 [Jacobson, Braden, and Borman 1992], which we describe in Section 26.6.

75-79 Notice that if there are no TCP connections active on the host (`tcb.inp_next` is null), neither `tcp_iss` nor `tcp_now` is incremented. This would occur only when the system is being initialized, since it would be rare to find a Unix system attached to a network without a few TCP servers active.

25.6 tcp_timers Function

The function `tcp_timers` is called by TCP's `PRU_SLOWTIMO` request (Figure 30.10):

```
case PRU_SLOWTIMO:
    tp = tcp_timers(tp, (int)nam);
```

when any one of the four TCP timer counters reaches 0 (Figure 25.8).

The structure of the function is a switch statement with one case per timer, as outlined in Figure 25.9.

```
120 struct tcpcb *
121 tcp_timers(tp, timer)
122 struct tcpcb *tp;
123 int timer;
124 {
125     int rexmt;
126     switch (timer) {
127
128         /* switch cases */
129
130     }
131     return (tp);
132 }
```

tcp_timer.c

Figure 25.9 `tcp_timers` function: general organization.

We now discuss three of the four timer counters (five of TCP's timers), saving the retransmission timer for Section 25.11.

FIN_WAIT_2 and 2MSL Timers

TCP's TCPT_2MSL counter implements two of TCP's timers.

1. **FIN_WAIT_2 timer.** When `tcp_input` moves from the `FIN_WAIT_1` state to the `FIN_WAIT_2` state and the socket cannot receive any more data (implying the process called `close`, instead of taking advantage of TCP's half-close with `shutdown`), the `FIN_WAIT_2` timer is set to 10 minutes (`tcp_maxidle`). We'll see that this prevents the connection from staying in the `FIN_WAIT_2` state forever.
2. **2MSL timer.** When TCP enters the `TIME_WAIT` state, the 2MSL timer is set to 60 seconds (TCPTV_MSL times 2).

Figure 25.10 shows the case for the 2MSL timer—executed when the timer reaches 0.

```

127      /*                                     tcp_timer.c
128      * 2 MSL timeout in shutdown went off.  If we're closed but
129      * still waiting for peer to close and connection has been idle
130      * too long, or if 2MSL time is up from TIME_WAIT, delete connection
131      * control block.  Otherwise, check again in a bit.
132      */
133      case TCPT_2MSL:
134          if (tp->t_state != TCPS_TIME_WAIT &&
135              tp->t_idle <= tcp_maxidle)
136              tp->t_timer[TCPT_2MSL] = tcp_keepintvl;
137          else
138              tp = tcp_close(tp);
139          break;

```

Figure 25.10 `tcp_slowtimo` function: expiration of 2MSL timer counter.

2MSL timer

¹²⁷⁻¹³⁹ The puzzling logic in the conditional is because the two different uses of the `TCPT_2MSL` counter are intermixed (Exercise 25.4). Let's first look at the `TIME_WAIT` state. When the timer expires after 60 seconds, `tcp_close` is called and the control blocks are released. We have the scenario shown in Figure 25.11. This figure shows the series of function calls that occurs when the 2MSL timer expires. We also see that setting one of the timers for N seconds in the future ($2 \times N$ ticks), causes the timer to expire somewhere between $2 \times N - 1$ and $2 \times N$ ticks in the future, since the time until the first decrement of the counter is between 0 and 500 ms in the future.

FIN_WAIT_2 timer

¹²⁷⁻¹³⁹ If the connection state is not `TIME_WAIT`, the `TCPT_2MSL` counter is the `FIN_WAIT_2` timer. As soon as the connection has been idle for more than 10 minutes (`tcp_maxidle`) the connection is closed. But if the connection has been idle for less than or equal to 10 minutes, the `FIN_WAIT_2` timer is reset for 75 seconds in the future. Figure 25.12 shows the typical scenario.

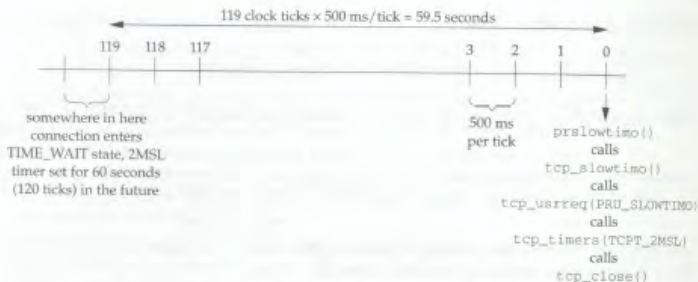


Figure 25.11 Setting and expiration of 2MSL timer in `TIME_WAIT` state.

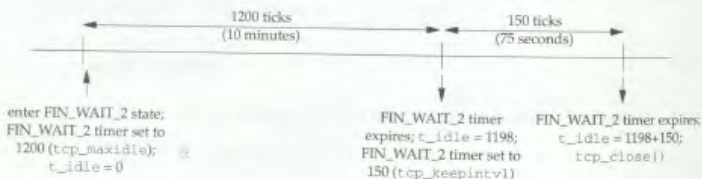


Figure 25.12 `FIN_WAIT_2` timer to avoid infinite wait in `FIN_WAIT_2` state.

The connection moves from the `FIN_WAIT_1` state to the `FIN_WAIT_2` state on the receipt of an ACK (Figure 24.15). Receiving this ACK sets `t_idle` to 0 and the `FIN_WAIT_2` timer is set to 1200 (`tcp_maxidle`). In Figure 25.12 we show the up arrow just to the right of the tick mark starting the 10-minute period, to reiterate that the first decrement of the counter occurs between 0 and 500 ms after the counter is set. After 1199 ticks the timer expires, but since `t_idle` is incremented *after* the test and decrement of the four counters in Figure 25.8, `t_idle` is 1198. (We assume the connection is idle for this 10-minute period.) The comparison of 1198 as less than or equal to 1200 is true, so the `FIN_WAIT_2` timer is set to 150 (`tcp_keepintvl`). When the timer expires again in 75 seconds, assuming the connection is still idle, `t_idle` is now 1348, the test is false, and `tcp_close` is called.

The reason for the 75-second timeout after the first 10-minute timeout is as follows: a connection in the `FIN_WAIT_2` state is not dropped until the connection has been idle for *more than* 10 minutes. There's no reason to test `t_idle` until at least 10 minutes have expired, but once this time has passed, the value of `t_idle` is checked every 75 seconds. Since a duplicate segment could be received, say a duplicate of the ACK that

moved the connection from the `FIN_WAIT_1` state to the `FIN_WAIT_2` state, the 10-minute wait is restarted when the segment is received (since `t_idle` will be set to 0).

Terminating an idle connection after more than 10 minutes in the `FIN_WAIT_2` state violates the protocol specification, but this is practical. In the `FIN_WAIT_2` state the process has called `close`, all outstanding data on the connection has been sent and acknowledged, the other end has acknowledged the FIN, and TCP is waiting for the process at the other end of the connection to issue its `close`. If the other process never closes its end of the connection, our end can remain in the `FIN_WAIT_2` forever. A counter should be maintained for the number of connections terminated for this reason, to see how often this occurs.

Persist Timer

Figure 25.13 shows the case for when the persist timer expires.

```

210      /*
211      * Persistence timer into zero window.
212      * Force a byte to be output, if possible.
213      */
214      case TCPT_PERSIST:
215          tcpstat.tcps_persisttime++;
216          tcp_setpersist(tp);
217          tp->t_force = 1;
218          (void) tcp_output(tp);
219          tp->t_force = 0;
220          break;

```

tcp_timer.c

tcp_timer.c

Figure 25.13 `tcp_slowtimo` function: expiration of persist timer.

Force window probe segment

210-220 When the persist timer expires, there is data to send on the connection but TCP has been stopped by the other end's advertisement of a zero-sized window. `tcp_setpersist` calculates the next value for the persist timer and stores it in the `TCPT_PERSIST` counter. The flag `t_force` is set to 1, forcing `tcp_output` to send 1 byte, even though the window advertised by the other end is 0.

Figure 25.14 shows typical values of the persist timer for a LAN, assuming the retransmission timeout for the connection is 1.5 seconds (see Figure 22.1 of Volume 1).



Figure 25.14 Time line of persist timer when probing a zero window.

Once the value of the persist timer reaches 60 seconds, TCP continues sending window probes every 60 seconds. The reason the first two values are both 5, and not 1.5 and 3, is that the persist timer is lower bounded at 5 seconds. It is also upper bounded at 60 seconds. The multiplication of each value by 2 to give the next value is called an *exponential backoff*, and we describe how it is calculated in Section 25.9.

Connection Establishment and Keepalive Timers

TCP's `TCPT_KEEPC` counter implements two timers:

1. When a SYN is sent, the connection-establishment timer is set to 75 seconds (`TCPTV_KEEPC_INIT`). This happens when `connect` is called, putting a connection into the `SYN_SENT` state (active open), or when a connection moves from the `LISTEN` to the `SYN_RCVD` state (passive open). If the connection doesn't enter the `ESTABLISHED` state within 75 seconds, the connection is dropped.
2. When a segment is received on a connection, `tcp_input` resets the keepalive timer for that connection to 2 hours (`tcp_keepidle`), and the `idle` counter for the connection is reset to 0. This happens for every TCP connection on the system, whether the keepalive option is enabled for the socket or not. If the keepalive timer expires (2 hours after the last segment was received on the connection), and if the socket option is set, a keepalive probe is sent to the other end. If the timer expires and the socket option is not set, the keepalive timer is just reset for 2 hours in the future.

Figure 25.15 shows the case for TCP's `TCPT_KEEPC` counter.

Connection-establishment timer expires after 75 seconds

221-228

If the state is less than `ESTABLISHED` (Figure 24.16), the `TCPT_KEEPC` counter is the connection-establishment timer. At the label `dropit`, `tcp_drop` is called to terminate the connection attempt with an error of `ETIMEDOUT`. We'll see that this error is the default error—if, for example, a soft error such as an ICMP host unreachable was received on the connection, the error returned to the process will be changed to `EHOSTUNREACH` instead of the default.

In Figure 30.4 we'll see that when TCP sends a SYN, two timers are initialized: the connection-establishment timer as we just described, with a value of 75 seconds, and the retransmission timer, to cause the SYN to be retransmitted if no response is received. Figure 25.16 shows these two timers.

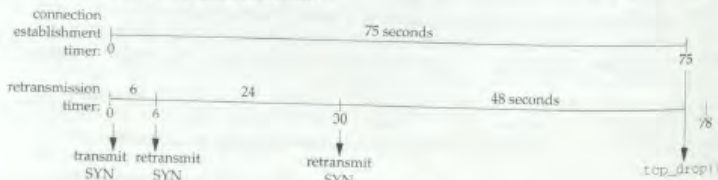


Figure 25.16 Connection-establishment timer and retransmission timer after SYN is sent.

The retransmission timer is initialized to 6 seconds for a new connection (Figure 25.19), and successive values are calculated to be 24 and 48 seconds. We describe how these values are calculated in Section 25.7. The retransmission timer causes the SYN to be

```

221      /*
222      * Keep-alive timer went off; send something
223      * or drop connection if idle for too long.
224      */
225      case TCPT_KEEF:
226          tcpstat.tcps_keeptimeo++;
227          if (tp->t_state < TCPS_ESTABLISHED)
228              goto dropit; /* connection establishment timer */
229
230          if (tp->t_inpcb->inp_socket->so_options & SO_KEEPAIVE &&
231              tp->t_state <= TCPS_CLOSE_WAIT) {
232              if (tp->t_idle >= tcp_keepidle + tcp_maxidle)
233                  goto dropit;
234
235              /*
236               * Send a packet designed to force a response
237               * if the peer is up and reachable:
238               * either an ACK if the connection is still alive
239               * or an RST if the peer has closed the connection
240               * due to timeout or reboot.
241               * Using sequence number tp->and_una-1
242               * causes the transmitted zero-length segment
243               * to lie outside the receive window;
244               * by the protocol spec, this requires the
245               * correspondent TCP to respond.
246               */
247              tcpstat.tcps_keepprobe++;
248              tcp_respond(tp, tp->t_template, (struct mbuf *) NULL,
249                          tp->rcv_nxt, tp->snd_una - 1, 0);
250              tp->t_timer[TCPT_KEEF] = tcp_keepintvl;
251          } else
252              tp->t_timer[TCPT_KEEF] = tcp_keepidle;
253          break;
254      dropit:
255          tcpstat.tcps_keepdrops++;
256          tp = tcp_drop(tp, ETIMEDOUT);
257          break;

```

Figure 25.15 tcp_slowtimo function: expiration of keepalive timer.

transmitted a total of three times, at times 0, 6, and 30. At time 75, 3 seconds before the retransmission timer would expire again, the connection-establishment timer expires, and `tcp_drop` terminates the connection attempt.

Keepalive timer expires after 2 hours of idle time

229-230

This timer expires after 2 hours of idle time on every connection, not just ones with the `SO_KEEPAIVE` socket option enabled. If the socket option is set, probes are sent only if the connection is in the `ESTABLISHED` or `CLOSE_WAIT` states (Figure 24.15). Once the process calls `close` (the states greater than `CLOSE_WAIT`), keepalive probes are not sent, even if the connection is idle for 2 hours.

Drop connection when no response

231-232 If the total idle time for the connection is greater than or equal to 2 hours (`tcp_keepidle`) plus 10 minutes (`tcp_maxidle`), the connection is dropped. This means that TCP has sent its limit of nine keepalive probes, 75 seconds apart (`tcp_keepintvl`), with no response. One reason TCP must send multiple keepalive probes before considering the connection dead is that the ACKs sent in response do not contain data and therefore are not reliably transmitted by TCP. An ACK that is a response to a keepalive probe can get lost.

Send a keepalive probe

233-248 If TCP hasn't reached the keepalive limit, `tcp_respond` sends a keepalive packet. The acknowledgment field of the keepalive packet (the fourth argument to `tcp_respond`) contains `rcv_nxt`, the next sequence number expected on the connection. The sequence number field of the keepalive packet (the fifth argument) deliberately contains `snd_una` minus 1, which is the sequence number of a byte of data that the other end has already acknowledged (figure 24.17). Since this sequence number is outside the window, the other end must respond with an ACK, specifying the next sequence number it expects.

Figure 25.17 summarizes this use of the keepalive timer.

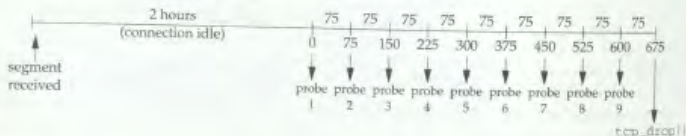


Figure 25.17 Summary of keepalive timer to detect unreachability of other end.

The nine keepalive probes are sent every 75 seconds, starting at time 0, through time 600. At time 675 (11.25 minutes after the 2-hour timer expired) the connection is dropped. Notice that nine keepalive probes are sent, even though the constant `TCPTV_KEEPCNT` (Figure 25.4) is 8. This is because the variable `t_idle` is incremented in Figure 25.8 *after* the timer is decremented, compared to 0, and possibly handled. When `tcp_input` receives a segment on a connection, it sets the keepalive timer to 14400 (`tcp_keepidle`) and `t_idle` to 0. The next time `tcp_slowtimo` is called, the keepalive timer is decremented to 14339 and `t_idle` is incremented to 1. About 2 hours later, when the keepalive timer is decremented from 1 to 0 and `tcp_timers` is called, the value of `t_idle` will be 14399. We can build the table in Figure 25.18 to see the value of `t_idle` each time `tcp_timers` is called.

The code in Figure 25.15 is waiting for `t_idle` to be greater than or equal to 15600 (`tcp_keepidle` + `tcp_maxidle`) and that only happens at time 675 in Figure 25.17, after nine keepalive probes have been sent.

probe#	time in Figure 25.17	<code>t_idle</code>
1	0	14399
2	75	14549
3	150	14699
4	225	14849
5	300	14999
6	375	15149
7	450	15299
8	525	15449
9	600	15599
	675	15749

Figure 25.18 The value of `t_idle` when `tcp_timers` is called for keepalive processing.

Reset keepalive timer

249-250 If the socket option is not set or the connection state is greater than `CLOSE_WAIT`, the keepalive timer for this connection is reset to 2 hours (`tcp_keepidle`).

Unfortunately the counter `tcps_keepdrops` (line 253) counts both uses of the `TCPT_KEEPC` counter: the connection-establishment timer and the keepalive timer.

25.7 Retransmission Timer Calculations

The timers that we've described so far in this chapter have fixed times associated with them: 200 ms for the delayed ACK timer, 75 seconds for the connection-establishment timer, 2 hours for the keepalive timer, and so on. The final two timers that we describe, the retransmission timer and the persist timer, have values that depend on the measured RTT for the connection. Before going through the source code that calculates and sets these timers we need to understand how TCP measures the RTT for a connection.

Fundamental to the operation of TCP is setting a retransmission timer when a segment is transmitted and an ACK is required from the other end. If the ACK is not received when the retransmission timer expires, the segment is retransmitted. TCP requires an ACK for data segments but does not require an ACK for a segment without data (i.e., a pure ACK segment). If the calculated retransmission timeout is too small, it can expire prematurely, causing needless retransmissions. If the calculated value is too large, after a segment is lost, additional time is lost before the segment is retransmitted, degrading performance. Complicating this is that the round-trip times between two hosts can vary widely and dynamically over the course of a connection.

TCP in Net/3 calculates the retransmission timeout (*RTO*) by measuring the round-trip time (*rticks*) of data segments and keeping track of the smoothed RTT estimator (*srtt*) and a smoothed mean deviation estimator (*rttvar*). The mean deviation is a good approximation of the standard deviation, but easier to compute since, unlike the standard deviation, the mean deviation does not require square root calculations. [Jacobson 1988b] provides additional details on these RTT measurements, which lead to the following equations:

$$\begin{aligned} \mathit{delta} &= \mathit{nticks} - \mathit{srtt} \\ \mathit{srtt} &\leftarrow \mathit{srtt} + g \times \mathit{delta} \\ \mathit{rttvar} &\leftarrow \mathit{rttvar} + h(|\mathit{delta}| - \mathit{rttvar}) \\ \mathit{RTO} &= \mathit{srtt} + 4 \times \mathit{rttvar} \end{aligned}$$

delta is the difference between the measured round trip just obtained (*nticks*) and the current smoothed RTT estimator (*srtt*). *g* is the gain applied to the RTT estimator and equals $\frac{1}{8}$. *h* is the gain applied to the mean deviation estimator and equals $\frac{1}{4}$. The two gains and the multiplier 4 in the RTO calculation are purposely powers of 2, so they can be calculated using shift operations instead of multiplying or dividing.

[Jacobson 1988b] specified $2 \times \mathit{rttvar}$ in the calculation of RTO, but after further research, [Jacobson 1990d] changed the value to $4 \times \mathit{rttvar}$, which is what appeared in the Net/1 implementation.

We now describe the variables and calculations used to calculate TCP's retransmission timer, as we'll encounter them throughout the TCP code. Figure 25.19 lists the variables in the control block related to the retransmission timer.

tcpcb member	Units	tcp_newtcpcb initial value	#sec	Description
<i>t_srtt</i>	ticks \times 8	0		smoothed RTT estimator: <i>srtt</i> \times 8
<i>t_rttvar</i>	ticks \times 4	24	3	smoothed mean deviation estimator: <i>rttvar</i> \times 4
<i>t_rxtcur</i>	ticks	12	6	current retransmission timeout: RTO
<i>t_rttmin</i>	ticks	2	1	minimum value for retransmission timeout
<i>t_rxtshift</i>	n.a.	0		index into <i>tcp_backoff[]</i> array (exponential backoff)

Figure 25.19 Control block variables for calculation of retransmission timer.

We show the *tcp_backoff* array at the end of Section 25.9. The *tcp_newtcpcb* function sets the initial values for these variables, and we cover it in the next section. The term *shift* in the variable *t_rxtshift* and its limit *TCP_MAXRXTSHIFT* is not entirely accurate. The former is not used for bit shifting, but as Figure 25.19 indicates, it is an index into an array.

The confusing part of TCP's timeout calculations is that the two smoothed estimators maintained in the C code (*t_srtt* and *t_rttvar*) are fixed-point integers, instead of floating-point values. This is done to avoid floating-point calculations within the kernel, but it complicates the code.

To keep the scaled and unscaled variables distinct, we'll use the italic variables *srtt* and *rttvar* to refer to the unscaled variables in the earlier equations, and *t_srtt* and *t_rttvar* to refer to the scaled variables in the TCP control block.

Figure 25.20 shows four constants we encounter, which define the scale factors of 8 for *t_srtt* and 4 for *t_rttvar*.

Constant	Value	Description
TCP_RTT_SCALE	8	multiplier: $t_srtt = srtt \times 8$
TCP_RTT_SHIFT	3	shift: $t_srtt = srtt \ll 3$
TCP_RTTVAR_SCALE	4	multiplier: $t_rttvar = rttvar \times 4$
TCP_RTTVAR_SHIFT	2	shift: $t_rttvar = rttvar \ll 2$

Figure 25.20 Multipliers and shifts for RTT estimators.

25.8 tcp_newtcpcb Function

A new TCP control block is allocated and initialized by `tcp_newtcpcb`, shown in Figure 25.21. This function is called by TCP's `PRU_ATTACH` request when a new socket is created (Figure 30.2). The caller has previously allocated an Internet PCB for this connection, pointed to by the argument `inp`. We present this function now because it initializes the TCP timer variables.

```

167 struct tcpcb *
168 tcp_newtcpcb(inp)
169 struct inpcb *inp;
170 {
171     struct tcpcb *tp;
172     tp = malloc(sizeof(*tp), M_PCB, M_NOWAIT);
173     if (tp == NULL)
174         return ((struct tcpcb *) 0);
175     bzero((char *) tp, sizeof(struct tcpcb));
176     tp->seg_next = tp->seg_prev = (struct tcpiphdr *) tp;
177     tp->t_maxseg = tcp_maxdfilt;
178     tp->t_flags = tcp_do_rfc1323 ? (TF_RBQ_SCALE | TF_RBQ_TSTMP) : 0;
179     tp->t_inpcb = inp;
180     /*
181      * Init srtt to TCPTV_SRTTBASE (0), so we can tell that we have no
182      * rtt estimate. Get rttvar so that srtt + 2 * rttvar gives
183      * reasonable initial retransmit time.
184      */
185     tp->t_srtt = TCPTV_SRTTBASE;
186     tp->t_rttvar = tcp_rtt0flt * PR_SLOWHZ << 2;
187     tp->t_rttmin = TCPTV_MIN;
188     TCPTV_RANGESET(tp->t_rxtcur,
189                    ((TCPTV_SRTTBASE >> 2) + (TCPTV_SRTTDFLT << 2)) >> 1,
190                    TCPTV_MIN, TCPTV_REXMTMAX);
191     tp->snd_cwnd = TCP_MAXWIN << TCP_MAX_WINSHIFT;
192     tp->snd_ssthresh = TCP_MAXWIN << TCP_MAX_WINSHIFT;
193     inp->inp_ip.ip_ttl = ip_defttl;
194     inp->inp_ppcb = (caddr_t) tp;
195     return (tp);
196 }

```

Figure 25.21 `tcp_newtcpcb` function: create and initialize a new TCP control block.

167-175 The kernel's `malloc` function allocates memory for the control block, and `bzero` sets it to 0.

176 The two variables `seg_next` and `seg_prev` point to the reassembly queue for out-of-order segments received for this connection. We discuss this queue in detail in Section 27.9.

177-178 The maximum segment size to send, `t_maxseg`, defaults to 512 (`tcp_mssdflt`). This value can be changed by the `tcp_mss` function after an MSS option is received from the other end. (TCP also sends an MSS option to the other end when a new connection is established.) The two flags `TF_REQ_SCALE` and `TF_REQ_TSTMP` are set if the system is configured to request window scaling and timestamps as defined in RFC 1323 (the global `tcp_do_rfc1323` from Figure 24.3, which defaults to 1). The `t_inpcb` pointer in the TCP control block is set to point to the Internet PCB passed in by the caller.

189-195 The four variables `t_srtt`, `t_rttvar`, `t_rttmin`, and `t_rxtcur`, described in Figure 25.19, are initialized. First, the smoothed RTT estimator `t_srtt` is set to 0 (`TCPTV_SRTTBASE`), which is a special value that means no RTT measurements have been made yet for this connection. `tcp_xmit_timer` recognizes this special value when the first RTT measurement is made.

186-187 The smoothed mean deviation estimator `t_rttvar` is set to 24: 3 (`tcp_rttdeflt`, from Figure 24.3) times 2 (`PR_SLOWHZ`) multiplied by 4 (the left shift of 2 bits). Since this scaled estimator is 4 times the variable `rttvar`, this value equals 6 clock ticks, or 3 seconds. The minimum RTO, stored in `t_rttmin`, is 2 ticks (`TCPTV_MIN`).

188-190 The current RTO in clock ticks is calculated and stored in `t_rxtcur`. It is bounded by a minimum value of 2 ticks (`TCPTV_MIN`) and a maximum value of 128 ticks (`TCPTV_REXMTMAX`). The value calculated as the second argument to `TCPT_RANGESET` is 12 ticks, or 6 seconds. This is the first RTO for the connection.

Understanding these C expressions involving the scaled RTT estimators can be a challenge. It helps to start with the unscaled equation and substitute the scaled variables. The unscaled equation we're solving is

$$RTO = srtt + 2 \times rttvar$$

where we use the multiplier of 2 instead of 4 to calculate the first RTO.

The use of the multiplier 2 instead of 4 appears to be a leftover from the original 4.3BSD Tahoe code [Paxson 1994].

Substituting the two scaling relationships

$$t_srtt = 8 \times srtt$$

$$t_rttvar = 4 \times rttvar$$

we get

$$\begin{aligned} RTO &= \frac{t_srtt}{8} + 2 \times \frac{t_rttvar}{4} \\ &= \frac{t_srtt}{4} + \frac{t_rttvar}{2} \end{aligned}$$

which is the C code for the second argument to `TCPT_RANGESET`. In this code the variable `t_rttvar` is not used—the constant `TCPTV_SRTTDFLT`, whose value is 6 ticks, is used instead, and it must be multiplied by 4 to have the same scale as `t_rttvar`.

191-192 The congestion window (`snd_cwnd`) and slow start threshold (`snd_ssthresh`) are set to 1,073,725,440 (approximately one gigabyte), which is the largest possible TCP window if the window scale option is in effect. (Slow start and congestion avoidance are described in Section 21.6 of Volume I.) It is calculated as the maximum value for the window size field in the TCP header (65535, `TCP_MAXWIN`) times 2^{14} , where 14 is the maximum value for the window scale factor (`TCP_MAX_WINSHIFT`). We'll see that when a SYN is sent or received on the connection, `tcp_mss` resets `snd_cwnd` to a single segment.

193-194 The default IP TTL in the Internet PCB is set to 64 (`ip_defttl`) and the PCB is set to point to the new TCP control block.

Not shown in this code is that numerous variables, such as the shift variable `t_rxtshift`, are implicitly initialized to 0 since the control block is initialized by `bzero`.

25.9 tcp_setpersist Function

The next function we look at that uses TCP's retransmission timeout calculations is `tcp_setpersist`. In Figure 25.13 we saw this function called when the persist timer expired. This timer is set when TCP has data to send on a connection, but the other end is advertising a window of 0. This function, shown in Figure 25.22, calculates and stores the next value for the timer.

```

493 void
494 tcp_setpersist(tp)
495 struct tcpcb *tp;
496 {
497     t = ((tp->t_srtt >> 2) + tp->t_rttvar) >> 1;
498     if (tp->t_timer[TCPT_REXMT])
499         panic("tcp_output REXMT");
500     /*
501      * Start/restart persistence timer.
502      */
503     TCPT_RANGESET(tp->t_timer[TCPT_PERSIST],
504                  t * tcp_backoff[tp->t_rxtshift],
505                  TCPTV_PERSMIN, TCPTV_PERSMAX);
506     if (tp->t_rxtshift < TCP_MAXRXTSHIFT)
507         tp->t_rxtshift++;
508 }

```

tcp_output.c

Figure 25.22 `tcp_setpersist` function: calculate and store a new value for the persist timer.

Check retransmission timer not enabled

493-499 A check is made that the retransmission timer is not enabled when the persist timer is about to be set, since the two timers are mutually exclusive: if data is being sent, the

other side must be advertising a nonzero window, but the persist timer is being set only if the advertised window is 0.

Calculate RTO

500-505 The variable t is set to the RTO value that was calculated at the beginning of the function. The equation being solved is

$$RTO = srtt + 2 \times rttvar$$

which is identical to the formula used at the end of the previous section. With substitution we get

$$RTO = \frac{\frac{t_srtt}{4} + t_rttvar}{2}$$

which is the value computed for the variable t .

Apply exponential backoff

506-507 An exponential backoff is also applied to the RTO. This is done by multiplying the RTO by a value from the `tcp_backoff` array:

```
int tcp_backoff[TCP_MAXRXTSHIFT + 1] =
    { 1, 2, 4, 8, 16, 32, 64, 64, 64, 64, 64, 64, 64 };
```

When `tcp_output` initially sets the persist timer for a connection, the code is

```
tp->t_rxtshift = 0;
tcp_setpersist(tp);
```

so the first time `tcp_setpersist` is called, `t_rxtshift` is 0. Since the value of `tcp_backoff[0]` is 1, t is used as the persist timeout. The `TCPT_RANGESET` macro bounds this value between 5 and 60 seconds. `t_rxtshift` is incremented by 1 until it reaches a maximum of 12 (`TCP_MAXRXTSHIFT`), since `tcp_backoff[12]` is the final entry in the array.

25.10 tcp_xmit_timer Function

The next function we look at, `tcp_xmit_timer`, is called each time an RTT measurement is collected, to update the smoothed RTT estimator (`srtt`) and the smoothed mean deviation estimator (`rttvar`).

The argument `rtt` is the RTT measurement to be applied. It is the value `nicks + 1`, using the notation from Section 25.7. It can be from one of two sources:

1. If the timestamp option is present in a received segment, the measured RTT is the current time (`tcp_now`) minus the timestamp value. We'll examine the timestamp option in Section 26.6, but for now all we need to know is that `tcp_now` is incremented every 500 ms (Figure 25.8). When a data segment is sent, `tcp_now` is sent as the timestamp, and the other end echoes this timestamp in the acknowledgment it sends back.

2. If timestamps are not in use and a data segment is being timed, we saw in Figure 25.8 that the counter `t_rtt` is incremented every 500 ms for the connection. We also mentioned in Section 25.5 that this counter is initialized to 1, so when the acknowledgment is received the counter is the measured RTT (in ticks) plus 1.

Typical code in `tcp_input` that calls `tcp_xmit_timer` is

```
if (ts_present)
    tcp_xmit_timer(tp, tcp_now - ts_eckr + 1);

else if (tp->t_rtt && SEQ_GT(ti->ti_ack, tp->t_rtsseq))
    tcp_xmit_timer(tp, tp->t_rtt);
```

If a timestamp was present in the segment (`ts_present`), the RTT estimators are updated using the current time (`tcp_now`) minus the echoed timestamp (`ts_eckr`) plus 1. (We describe the reason for adding 1 below.)

If a timestamp is not present, the RTT estimators are updated only if the received segment acknowledges a data segment that was being timed. There is only one RTT counter per TCP control block (`t_rtt`), so only one outstanding data segment can be timed per connection. The starting sequence number of that segment is stored in `t_rtsseq` when the segment is transmitted, to tell when an acknowledgment is received that covers that sequence number. If the received acknowledgment number (`ti_ack`) is greater than the starting sequence number of the segment being timed (`t_rtsseq`), the RTT estimators are updated using `t_rtt` as the measured RTT.

Before RFC 1323 timestamps were supported, TCP measured the RTT only by counting clock ticks in `t_rtt`. But this variable is also used as a flag that specifies whether a segment is being timed (Figure 25.8): if `t_rtt` is greater than 0, then `tcp_slowtimo` adds 1 to it every 500 ms. Hence when `t_rtt` is nonzero, it is the number of ticks plus 1. We'll see shortly that `tcp_xmit_timer` always decrements its second argument by 1 to account for this offset. Therefore when timestamps are being used, 1 is added to the second argument to account for the decrement by 1 in `tcp_xmit_timer`.

The greater-than test of the sequence numbers is because ACKs are cumulative: if TCP sends and times a segment with sequence numbers 1–1024 (`t_rtsseq` equals 1), then immediately sends (but can't time) a segment with sequence numbers 1025–2048, and then receives an ACK with `ti_ack` equal to 2049, this is an ACK for sequence numbers 1–2048 and the ACK acknowledges the first segment being timed as well as the second (untimed) segment. Notice that when RFC 1323 timestamps are in use there is no comparison of sequence numbers. If the other end sends a timestamp option, it chooses the echo reply value (`ts_eckr`) to allow TCP to calculate the RTT.

Figure 25.23 shows the first part of the function that updates the estimators.

Update smoothed estimators

1010-1125 Recall that `tcp_newtcpcb` initialized the smoothed RTT estimator (`t_srtt`) to 0, indicating that no measurements have been made for this connection. `delta` is the difference between the measured RTT and the current value of the smoothed RTT estimator, in unscaled ticks. `t_srtt` is divided by 8 to convert from scaled to unscaled ticks.

```

1310 void
1311 tcp_xmit_timer(tp, rtt)
1312 struct tcpcb *tp;
1313 short rtt;
1314 {
1315     short delta;

1316     tcpstat.tcps_rttupdated++;
1317     if (tp->t_srtt != 0) {
1318         /*
1319          * srtt is stored as fixed point with 3 bits after the
1320          * binary point (i.e., scaled by 8). The following magic
1321          * is equivalent to the smoothing algorithm in rfc793 with
1322          * an alpha of .875 (srtt = rtt/8 + srtt*7/8 in fixed
1323          * point). Adjust rtt to origin 0.
1324          */
1325         delta = rtt - 1 - (tp->t_srtt >> TCP_RTT_SHIFT);
1326         if ((tp->t_srtt += delta) <= 0)
1327             tp->t_srtt = 1;
1328         /*
1329          * We accumulate a smoothed rtt variance (actually, a
1330          * smoothed mean difference), then set the retransmit
1331          * timer to smoothed rtt * 4 times the smoothed variance.
1332          * rttvar is stored as fixed point with 2 bits after the
1333          * binary point (scaled by 4). The following is
1334          * equivalent to rfc793 smoothing with an alpha of .75
1335          * (rttvar = rttvar*3/4 + |delta| / 4). This replaces
1336          * rfc793's wired-in beta.
1337          */
1338         if (delta < 0)
1339             delta = -delta;
1340         delta -= (tp->t_rttvar >> TCP_RTTVAR_SHIFT);
1341         if ((tp->t_rttvar += delta) <= 0)
1342             tp->t_rttvar = 1;
1343     } else {
1344         /*
1345          * No rtt measurement yet - use the unsmoothed rtt.
1346          * Set the variance to half the rtt (so our first
1347          * retransmit happens at 3*rtt).
1348          */
1349         tp->t_srtt = rtt << TCP_RTT_SHIFT;
1350         tp->t_rttvar = rtt << (TCP_RTTVAR_SHIFT - 1);
1351     }
}

```

Figure 25.23 tcp_xmit_timer function: apply new RTT measurement to smoothed estimators.

1326-1327

The smoothed RTT estimator is updated using the equation

$$srtt \leftarrow srtt + g \times delta$$

Since the gain g is $\frac{1}{8}$, this equation is

$$8 \times srtt \leftarrow 8 \times srtt + \mathit{delta}$$

which is

$$t_srtt \leftarrow t_srtt + \mathit{delta}$$

1326-1342 The mean deviation estimator is updated using the equation

$$rttvar \leftarrow rttvar + h(|\mathit{delta}| - rttvar)$$

Substituting $\frac{1}{4}$ for h and the scaled variable t_rttvar for $4 \times rttvar$, we get

$$\frac{t_rttvar}{4} \leftarrow \frac{t_rttvar}{4} + \frac{|\mathit{delta}| - \frac{t_rttvar}{4}}{4}$$

which is

$$t_rttvar \leftarrow t_rttvar + |\mathit{delta}| - \frac{t_rttvar}{4}$$

This final equation corresponds to the C code.

Initialize smoothed estimators on first RTT measurement

1343-1350 If this is the first RTT measured for this connection, the smoothed RTT estimator is initialized to the measured RTT. These calculations use the value of the argument rtt , which we said is the measured RTT plus 1 ($nticks + 1$), whereas the earlier calculation of delta subtracted 1 from rtt .

$$srtt = nticks + 1$$

or

$$\frac{t_srtt}{8} = nticks + 1$$

which is

$$t_srtt = (nticks + 1) \times 8$$

The smoothed mean deviation is set to one-half of the measured RTT:

$$rttvar = \frac{srtt}{2}$$

which is

$$\frac{t_rttvar}{4} = \frac{nticks + 1}{2}$$

or

$$t_rttvar = (nticks + 1) \times 2$$

The comment in the code states that this initial setting for the smoothed mean deviation yields an initial RTO of $3 \times srtt$. Since the RTO is calculated as

$$RTO = srtt + 4 \times rttvar$$

substituting for $rttvar$ gives us

$$RTO = srtt + 4 \times \frac{srtt}{2}$$

which is indeed

$$RTO = 3 \times srtt$$

Figure 25.24 shows the final part of the `tcp_xmit_timer` function.

```

1352     tp->t_rtt = 0;
1353     tp->t_rxtshift = 0;
1354     /*
1355     * the retransmit should happen at rtt + 4 * rttvar.
1356     * Because of the way we do the smoothing, srtt and rttvar
1357     * will each average +1/2 tick of bias. When we compute
1358     * the retransmit timer, we want 1/2 tick of rounding and
1359     * 1 extra tick because of +-1/2 tick uncertainty in the
1360     * firing of the timer. The bias will give us exactly the
1361     * 1.5 tick we need. But, because the bias is
1362     * statistical, we have to test that we don't drop below
1363     * the minimum feasible timer (which is 2 ticks).
1364     */
1365     TCPT_RANGESET(tp->t_rxtcur, TCP_REXMTVAL(tp),
1366                 tp->t_rttmin, TCP_TV_REXMTMAX);
1367     /*
1368     * We received an ack for a packet that wasn't retransmitted;
1369     * it is probably safe to discard any error indications we've
1370     * received recently. This isn't quite right, but close enough
1371     * for now (a route might have failed after we sent a segment,
1372     * and the return path might not be symmetrical).
1373     */
1374     tp->t_softerror = 0;
1375 }

```

tcp_input.c

Figure 25.24 `tcp_xmit_timer` function: final part.

1352-1353 The RIT counter (`t_rtt`) and the retransmission shift count (`t_rxtshift`) are both reset to 0 in preparation for timing and transmission of the next segment.

1354-1366 The next RTO to use for the connection (`t_rxtcur`) is calculated using the macro

```

#define TCP_REXMTVAL(tp) \
    (((tp)->t_srtt >> TCP_RTT_SHIFT) + (tp)->t_rttvar)

```

This is the now-familiar equation

$$RTO = srtt + 4 \times rttvar$$

using the scaled variables updated by `tcp_xmit_timer`. Substituting these scaled variables for $srtt$ and $rttvar$, we have

$$RTO = \frac{t_srtt}{8} + 4 \times \frac{t_rttvar}{4}$$

$$= \frac{t_srtt}{8} + t_rttvar$$

which corresponds to the macro. The calculated value for the *RTO* is bounded by the minimum *RTO* for this connection (`t_rttmin`, which `t_newtcpcb` set to 2 ticks), and 128 ticks (`TCPTV_REXMTMAX`).

Clear soft error variable

1867-1374 Since `tcp_xmit_timer` is called only when an acknowledgment is received for a data segment that was sent, if a soft error was recorded for this connection (`t_softerror`), that error is discarded. We describe soft errors in more detail in the next section.

25.11 Retransmission Timeout: `tcp_timers` Function

We now return to the `tcp_timers` function and cover the final case that we didn't present in Section 25.6: the one that handles the expiration of the retransmission timer. This code is executed when a data segment that was transmitted has not been acknowledged by the other end within the *RTO*.

Figure 25.25 summarizes the actions caused by the retransmission timer. We assume that the first timeout calculated by `tcp_output` is 1.5 seconds, which is typical for a LAN (see Figure 21.1 of Volume 1).

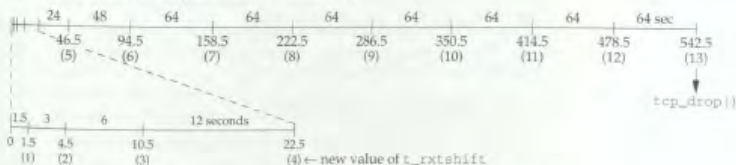


Figure 25.25 Summary of retransmission timer when sending data.

The x-axis is labeled with the time in seconds: 0, 1.5, 4.5, and so on. Below each of these numbers we show the value of `t_rxtshift` that is used in the code we're about to examine. Only after 12 retransmissions and a total of 542.5 seconds (just over 9 minutes) does TCP give up and drop the connection.

RFC 793 recommended that an open of a new connection, active or passive, allow a parameter specifying the total timeout period for data sent by TCP. This is the total amount of time TCP will try to send a given segment before giving up and terminating the connection. The recommended default was 5 minutes.

RFC 1122 requires that an application must be able to specify a parameter for a connection giving either the total number of retransmissions or the total timeout value for data sent by TCP. This parameter can be specified as "infinity," meaning TCP never gives up, allowing, perhaps, an interactive user the choice of when to give up.

We'll see in the code described shortly that Net/3 does not give the application any of this control: a fixed number of retransmissions (12) always occurs before TCP gives up, and the total timeout before giving up depends on the RTT.

The first half of the `tcp_timers` function is shown in Figure 25.26.

```

140      /*
141      * Retransmission timer went off. Message has not
142      * been acked within retransmit interval. Back off
143      * to a longer retransmit interval and retransmit one segment.
144      */
145      case TCPT_REXMT:
146      if (++tp->t_rxtshift > TCP_MAXRXTSHIFT) {
147          tp->t_rxtshift = TCP_MAXRXTSHIFT;
148          tcpstat.tcps_timeoutdrop++;
149          tp = tcp_drop(tp, tp->t_softerror ?
150                      tp->t_softerror : ETIMEDOUT);
151          break;
152      }
153      tcpstat.tcps_rexmttimeo++;
154      rexmt = TCP_REXMTVAL(tp) * tcp_backoff(tp->t_rxtshift);
155      TCP_RANGESET(tp->t_rxtcur, rexmt,
156                 tp->t_rttmin, TCPTV_REXMTMAX);
157      tp->t_timer[TCPT_REXMT] = tp->t_rxtcur;
158      /*
159      * If losing, let the lower level know and try for
160      * a better route. Also, if we backed off this far,
161      * our srtt estimate is probably bogus. Clobber it
162      * so we'll take the next rtt measurement as our srtt;
163      * move the current srtt into rttvar to keep the current
164      * retransmit times until then.
165      */
166      if (tp->t_rxtshift > TCP_MAXRXTSHIFT / 4) {
167          in_losing(tp->t_inpcb);
168          tp->t_rttvar += (tp->t_srtt >> TCP_RTT_SHIFT);
169          tp->t_srtt = 0;
170      }
171      tp->snd_nxt = tp->snd_una;
172      /*
173      * If timing a segment in this window, stop the timer.
174      */
175      tp->t_rtt = 0;

```

tcp_timer.c

Figure 25.26 `tcp_slowtimo` function: expiration of retransmission timer, first half.

Increment shift count

186

The retransmission shift count (`t_rxtshift`) is incremented, and if the value exceeds 12 (`TCP_MAXRXTSHIFT`) it is time to drop the connection. This new value of `t_rxtshift` is what we show in Figure 25.25. Notice the difference between this dropping of a connection because an acknowledgment is not received from the other end in response to data sent by TCP, and the keepalive timer, which drops a connection after a

long period of inactivity and no response from the other end. Both report the error `ETIMEDOUT` to the process, unless a soft error is received for the connection.

Drop connection

141-152 A *soft error* is one that doesn't cause TCP to terminate an established connection or an attempt to establish a connection, but the soft error is recorded in case TCP gives up later. For example, if TCP retransmits a SYN segment to establish a connection, receiving nothing in response, the error returned to the process will be `ETIMEDOUT`. But if during the retransmissions an ICMP host unreachable is received for the connection, that is considered a soft error and stored in `t_softerror` by `tcp_notify`. If TCP finally gives up the retransmissions, the error returned to the process will be `EHOSTUNREACH` instead of `ETIMEDOUT`, providing more information to the process. If TCP receives an RST on the connection in response to the SYN, that's considered a *hard error* and the connection is terminated immediately with an error of `ECONNREFUSED` (Figure 28.18).

Calculate new RTO

153-159 The next *RTO* is calculated using the `TCP_REXMTVAL` macro, applying an exponential backoff. In this code, `t_rxtshift` will be 1 the first time a given segment is retransmitted, so the *RTO* will be twice the value calculated by `TCP_REXMTVAL`. This value is stored in `t_rxtcur` and as the retransmission timer for the connection, `t_timer[TCP_REXMT]`. The value stored in `t_rxtcur` is used in `tcp_input` when the retransmission timer is restarted (Figures 28.12 and 29.6).

Ask IP to find a new route

159-167 If this segment has been retransmitted four or more times, `in_losing` releases the cached route (if there is one), so when the segment is retransmitted by `tcp_output` (at the end of this `case` statement in Figure 25.27) a new, and hopefully better, route will be chosen. In Figure 25.25 `in_losing` is called each time the retransmission timer expires, starting with the retransmission at time 22.5.

Clear estimators

169-170 The smoothed RTT estimator (`t_srtt`) is set to 0, which is what `t_newtcpcb` did. This forces `tcp_xmit_timer` to use the next measured RTT as the smoothed RTT estimator. This is done because the retransmitted segment has been sent four or more times, implying that TCP's smoothed RTT estimator is probably way off. But if the retransmission timer expires again, at the beginning of this `case` statement the *RTO* is calculated by `TCP_REXMTVAL`. That calculation should generate the same value as it did for this retransmission (which will then be exponentially backed off), even though `t_srtt` is set to 0. (The retransmission at time 42.464 in Figure 25.28 is an example of what's happening here.)

To accomplish this the value of `t_rttvar` is changed as follows. The next time the *RTO* is calculated, the equation

$$RTO = \frac{t_srtt}{8} + t_rttvar$$

is evaluated. Since `t_srtt` will be 0, if `t_rttvar` is increased by `t_srtt` divided by

8, *RTO* will have the same value. If the retransmission timer expires again for this segment (e.g., times 84,064 through 217,184 in Figure 25.28), when this code is executed again `t_rtt` will be 0, so `t_rttvar` won't change.

Force retransmission of oldest unacknowledged data

171 The next send sequence number (`snd_nxt`) is set to the oldest unacknowledged sequence number (`snd_una`). Recall from Figure 24.17 that `snd_nxt` can be greater than `snd_una`. By moving `snd_nxt` back, the retransmission will be the oldest segment that hasn't been acknowledged.

Karn's algorithm

172-175 The RTT counter, `t_rtt`, is set to 0, in case the last segment transmitted was being timed. Karn's algorithm says that even if an ACK of that segment is received, since the segment is about to be retransmitted, any timing of the segment is worthless since the ACK could be for the first transmission or for the retransmission. The algorithm is described in [Karn and Partridge 1987] and in Section 21.3 of Volume 1. Therefore the only segments that are timed using the `t_rtt` counter and used to update the RTT estimators are those that are not retransmitted. We'll see in Figure 29.6 that the use of RFC 1323 timestamps overrides Karn's algorithm.

Slow Start and Congestion Avoidance

The second half of this case is shown in Figure 25.27. It performs slow start and congestion avoidance and retransmits the oldest unacknowledged segment.

Since a retransmission timeout has occurred, this is a strong indication of congestion in the network. TCP's *congestion avoidance algorithm* comes into play, and when a segment is eventually acknowledged by the other end, TCP's *slow start algorithm* will continue the data transmission on the connection at a slower rate. Sections 20.6 and 21.6 of Volume 1 describe the two algorithms in detail.

176-205 `win` is set to one-half of the current window size (the minimum of the receiver's advertised window, `snd_wnd`, and the sender's congestion window, `snd_cwnd`) in segments, not bytes (hence the division by `t_maxseg`). Its minimum value is two segments. This records one-half of the window size when the congestion occurred, assuming one cause of the congestion is our sending segments too rapidly into the network. This becomes the slow start threshold, `t_ssthresh` (which is stored in bytes, hence the multiplication by `t_maxseg`). The congestion window, `snd_cwnd`, is set to one segment, which forces slow start.

This code is enclosed in braces because it was added between the 4.3BSD and Net/1 releases and required its own local variable (`win`).

206 The counter of consecutive duplicate ACKs, `t_dupacks` (which is used by the fast retransmit algorithm in Section 29.4), is set to 0. We'll see how this counter is used with TCP's fast retransmit and fast recovery algorithms in Chapter 29.

208 `tcp_output` resends a segment containing the oldest unacknowledged sequence number. This is the retransmission caused by the retransmission timer expiring.

```

176      /*
177      * Close the congestion window down to one segment
178      * (we'll open it by one segment for each ack we get).
179      * Since we probably have a window's worth of unacked
180      * data accumulated, this "slow start" keeps us from
181      * dumping all that data as back-to-back packets (which
182      * might overwhelm an intermediate gateway).
183      *
184      * There are two phases to the opening: Initially we
185      * open by one mss on each ack. This makes the window
186      * size increase exponentially with time. If the
187      * window is larger than the path can handle, this
188      * exponential growth results in dropped packet(s)
189      * almost immediately. To get more time between
190      * drops but still "push" the network to take advantage
191      * of improving conditions, we switch from exponential
192      * to linear window opening at some threshold size.
193      * For a threshold, we use half the current window
194      * size, truncated to a multiple of the mss.
195      *
196      * (the minimum cwnd that will give us exponential
197      * growth is 2 mss. We don't allow the threshold
198      * to go below this.)
199      */
200      {
201          u_int win = min(tp->snd_wnd, tp->snd_cwnd) / 2 / tp->t_maxseg;
202          if (win < 2)
203              win = 2;
204          tp->snd_cwnd = tp->t_maxseg;
205          tp->snd_ssthresh = win * tp->t_maxseg;
206          tp->t_dupacks = 0;
207      }
208      (void) tcp_output(tp);
209      break;

```

tcp_timer.c

tcp_timer.c

Figure 25.27 tcp_slowtimo function: expiration of retransmission timer, second half.

Accuracy

How accurate are these estimators that TCP maintains? At first they appear too coarse, since the RTTs are measured in multiples of 500 ms. The mean and mean deviation are maintained with additional accuracy (factors of 8 and 4 respectively), but LANs have RTTs on the order of milliseconds, and a transcontinental RTT is around 60 ms. What these estimators provide is a solid upper bound on the RTT so that the retransmission timeout can be set without worrying that the timeout is too small, causing unnecessary and wasteful retransmissions.

[Brakmo, O'Malley, and Peterson 1994] describe a TCP implementation that provides higher-resolution RTT measurements. This is done by recording the system clock (which has a much higher resolution than 500 ms) when a segment is transmitted and reading the system clock when the ACK is received, calculating a higher-resolution RTT.

The timestamp option provided by Net/3 (Section 26.6) can provide higher-resolution RTTs, but Net/3 sets the resolution of these timestamps to 500 ms.

25.12 An RTT Example

We now go through an actual example to see how the calculations are performed. We transfer 12288 bytes from the host `bsd1` to `vangogh.cs.berkeley.edu`. During the transfer we purposely bring down the PPP link being used and then bring it back up, to see how timeouts and retransmissions are handled. To transfer the data we use our `sock` program (described in Appendix C of Volume 1) with the `-D` option, to enable the `SO_DEBUG` socket option (Section 27.10). After the transfer is complete we examine the debug records left in the kernel's circular buffer using the `trpt(8)` program and print the desired timer variables from the TCP control block.

Figure 25.28 shows the calculations that occur at the various times. We use the notation $M:N$ to mean that sequence numbers M through and including $N-1$ are sent. Each segment in this example contains 512 bytes. The notation "ack M " means that the acknowledgment field of the ACK is M . The column labeled "actual delta (ms)" shows the time difference between the RTT timer going on and going off. The column labeled "rtt (arg.)" shows the second argument to the `tcp_xmit_timer` function: the number of clock ticks plus 1 between the RTT timer going on and going off.

The function `tcp_newtcpcb` initializes `t_srtt`, `t_rttvar`, and `t_rxtcur` to the values shown at time 0.0.

The first segment timed is the initial SYN. When its ACK is received 365 ms later, `tcp_xmit_timer` is called with an `rtt` argument of 2. Since this is the first RTT measurement (`t_srtt` is 0), the `else` clause in Figure 25.23 calculates the first values of the smoothed estimators.

The data segment containing bytes 1 through 512 is the next segment timed, and the RTT variables are updated at time 1.259 when its ACK is received.

The next three segments show how ACKs are cumulative. The timer is started at time 1.260 when bytes 513 through 1024 are sent. Another segment is sent with bytes 1025 through 1536, and the ACK received at time 2.206 acknowledges both data segments. The RTT estimators are then updated, since the ACK covers the starting sequence number being timed (513).

The segment with bytes 1537 through 2048 is transmitted at time 2.206 and the timer is started. Just that segment is acknowledged at time 3.132, and the estimators updated.

The data segment at time 3.132 is timed and the retransmission timer is set to 5 ticks (the current value of `t_rxtcur`). Somewhere around this time the PPP link between the routers `sun` and `netb` is taken down and then brought back up, a procedure that takes a few minutes. When the retransmission timer expires at time 6.064, the code in Figure 25.26 is executed to update the RTT variables. `t_rxtshift` is incremented from 0 to 1 and `t_rxtcur` is set to 10 ticks (the exponential backoff). A segment starting with the oldest unacknowledged sequence number (`sequna`, which is 3073) is retransmitted. After 5 seconds the timer expires again, `t_rxtshift` is incremented to 2, and the retransmission timer is set to 20 ticks.

xmit time	send	rcv	RTT timer	actual delta (ms)	rtt arg.	t_srtt (ticks × 8)	t_rttvar (ticks × 4)	t_rxtcur (ticks)	t_rxtshift
0.0	SYN		on			0	24	12	
0.365		SYN,ACK	off	365	2	16	4	6	
0.365	ACK								
0.415	1:513		on						
1.259		ack 513	off	844	2	15	4	5	
1.260	513:1025		on						
1.261	1025:1537								
2.206		ack 1537	off	946	3	16	4	6	
2.206	1537:2049		on						
2.207	2049:2561								
2.209	2561:3073								
3.132		ack 2049	off	926	3	16	3	5	
3.132	3073:3585		on						
3.133	3585:4097								
3.736		ack 2561							
3.736	4097:4609								
3.737	4609:5121								
3.739		ack 3073							
3.739	5121:5633								
3.740	5633:6145								
6.064	3073:3585		off			16	3	10	1
11.264	3073:3585		off			16	3	20	2
21.664	3073:3585		off			16	3	40	3
42.464	3073:3585		off			0	5	80	4
84.064	3073:3585		off			0	5	128	5
180.624	3073:3585		off			0	5	128	6
217.184	3073:3585		off			0	5	128	7
217.944		ack 6145							
217.944	6145:6657		on						
217.945	6657:7169								
218.834		ack 6657	off	890	3	24	6	9	
218.834	7169:7681		on						
218.836	7681:8193								
219.209		ack 7169							
219.209	8193:8705								
219.760		ack 7681	off	926	2	22	7	9	
219.760	8705:9217		on						
220.103		ack 8705							
220.103	9217:9729								
220.105	9729:10241								
220.106	10241:10753								
220.821		ack 9217	off	1061	3	22	6	8	
220.821	10753:11265		on						
221.310		ack 9729							
221.310	11265:11777								
221.312		ack 10241							
221.312	11777:12289								
221.674		ack 10753							
221.955		ack 11265	off	1134	3	22	5	7	

Figure 25.28 Values of RTT variables and estimators during example.

When the retransmission timer expires at time 42.464, `t_srtt` is set to 0 and `t_rttvar` is set to 5. As we mentioned in our discussion of Figure 25.26, this leaves the calculation of `t_rxtcur` the same (so the next calculation yields 160), but by setting `t_srtt` to 0, the next time the RTT estimators are updated (at time 218.834), the measured RTT becomes the smoothed RTT, as if the connection were starting fresh.

The rest of the data transfer continues, and the estimators are updated a few more times.

25.13 Summary

The two functions `tcp_fasttimo` and `tcp_slowtimo` are called by the kernel every 200 ms and every 500 ms, respectively. These two functions drive TCP's per-connection timer maintenance.

TCP maintains the following seven timers for each connection:

- a connection-establishment timer,
- a retransmission timer,
- a delayed ACK timer,
- a persist timer,
- a keepalive timer,
- a FIN_WAIT_2 timer, and
- a 2MSL timer.

The delayed ACK timer is different from the other six, since when it is set it means a delayed ACK must be sent the next time TCP's 200-ms timer expires. The other six timers are counters that are decremented by 1 every time TCP's 500-ms timer expires. When any one of the counters reaches 0, the appropriate action is taken: drop the connection, retransmit a segment, send a keepalive probe, and so on, as described in this chapter. Since some of the timers are mutually exclusive, the six timers are really implemented using four counters, which complicates the code.

This chapter also introduced the recommended way to calculate values for the retransmission timer. TCP maintains two smoothed estimators for a connection: the round-trip time and the mean deviation of the RTT. Although the algorithms are simple and elegant, these estimators are maintained as scaled fixed-point numbers (to provide adequate precision without using floating-point code within the kernel), which complicates the code.

Exercises

- 25.1 How efficient is TCP's fast timeout function? (*Hint:* Look at the number of delayed ACKs in Figure 24.5.) Suggest alternative implementations.
- 25.2 Why do you think the initialization of `tcp_maxidle` is in the `tcp_slowtimo` function instead of the `tcp_init` function?
- 25.3 `tcp_slowtimo` increments `t_idle`, which we said counts the clock ticks since a segment was last received on the connection. Should TCP also count the idle time since a segment was last sent on a connection?
- 25.4 Rewrite the code in Figure 25.10 to separate the logic for the two different uses of the `TCPT_2MSL` counter.
- 25.5 75 seconds after the connection in Figure 25.12 enters the `FIN_WAIT_2` state a duplicate ACK is received on the connection. What happens?
- 25.6 A connection has been idle for 1 hour when the application sets the `SO_KEEPAALIVE` option. Will the first keepalive probe be sent 1 or 2 hours in the future?
- 25.7 Why is `tcp_rttdeflt` a global variable and not a constant?
- 25.8 Rewrite the code related to Exercise 25.6 to implement the alternate behavior.

[The page contains extremely faint, illegible text, likely bleed-through from the reverse side of the document. The text is too light to transcribe accurately.]

TCP Output

26.1 Introduction

The function `tcp_output` is called whenever a segment needs to be sent on a connection. There are numerous calls to this function from other TCP functions:

- `tcp_usrreq` calls it for various requests: `PRU_CONNECT` to send the initial SYN, `PRU_SHUTDOWN` to send a FIN, `PRU_RCVD` in case a window update can be sent after the process has read some data from the socket receive buffer, `PRU_SEND` to send data, and `PRU_SENDOOB` to send out-of-band data.
- `tcp_fasttimo` calls it to send a delayed ACK.
- `tcp_timers` calls it to retransmit a segment when the retransmission timer expires.
- `tcp_timers` calls it to send a persist probe when the persist timer expires.
- `tcp_drop` calls it to send an RST.
- `tcp_disconnect` calls it to send a FIN.
- `tcp_input` calls it when output is required or when an immediate ACK should be sent.
- `tcp_input` calls it when a pure ACK is processed by the header prediction code and there is more data to send. (A *pure ACK* is a segment without data that just acknowledges data.)
- `tcp_input` calls it when the third consecutive duplicate ACK is received, to send a single segment (the fast retransmit algorithm).

`tcp_output` first determines whether a segment should be sent or not. TCP output is controlled by numerous factors other than data being ready to send to the other end of the connection. For example, the other end might be advertising a window of size 0 that stops TCP from sending anything, the Nagle algorithm prevents TCP from sending lots of small segments, and slow start and congestion avoidance limit the amount of data TCP can send on a connection. Conversely, some functions set flags just to force `tcp_output` to send a segment, such as the `TF_ACKNOW` flag that means an ACK should be sent immediately and not delayed. If `tcp_output` decides not to send a segment, the data (if any) is left in the socket's send buffer for a later call to this function.

26.2 `tcp_output` Overview

`tcp_output` is a large function, so we'll discuss it in 14 parts. Figure 26.1 shows the outline of the function.

Is an ACK is expected from the other end?

81 `idle` is true if the maximum sequence number sent (`snd_max`) equals the oldest unacknowledged sequence number (`snd_una`), that is, if an ACK is not expected from the other end. In Figure 24.17 `idle` would be 0, since an ACK is expected for sequence numbers 4–6, which have been sent but not yet acknowledged.

Go back to slow start

69–88 If an ACK is not expected from the other end and a segment has not been received from the other end in one round-trip time, the congestion window is set to one segment (`t_maxseg` bytes). This forces slow start to occur for this connection the next time a segment is sent. When a significant pause occurs in the data transmission ("significant" being more than the RTT), the network conditions can change from what was previously measured on the connection. Net/3 assumes the worst and returns to slow start.

Send more than one segment

69–70 When `send` is jumped to, a single segment is sent by calling `ip_output`. But if `tcp_output` determines that more than one segment can be sent, `sendalot` is set to 1, and the function tries to send another segment. Therefore, one call to `tcp_output` can result in multiple segments being sent.

26.3 Determine if a Segment Should be Sent

Sometimes `tcp_output` is called but a segment is not generated. For example, the `PRU_RCVD` request is generated when the socket layer removes data from the socket's receive buffer, passing the data to a process. It is possible that the process removed enough data that TCP should send a segment to the other end with a new window advertisement, but this is just a possibility, not a certainty. The first half of `tcp_output` determines if there is a reason to send a segment to the other end. If not, the function returns without sending a segment.

```

43 int
44 tcp_output(tp)
45 struct tcpcb *tp;
46 {
47     struct socket *so = tp->t_inpcb->inp_socket;
48     long len, win;
49     int off, flags, error;
50     struct mbuf *m;
51     struct tcpiphdr *ti;
52     u_char opt[MAX_TCPOPTLEN];
53     unsigned optlen, hdrlen;
54     int idle, sendalot;
55
56     /*
57      * Determine length of data that should be transmitted
58      * and flags that will be used.
59      * If there are some data or critical controls (SYN, RST)
60      * to send, then transmit; otherwise, investigate further.
61      */
62     idle = (tp->snd_max == tp->snd_una);
63     if (idle && tp->t_idle >= tp->t_rxtcur)
64         /*
65          * We have been idle for "a while" and no acks are
66          * expected to clock out any data we send --
67          * slow start to get ack "clock" running again.
68          */
69         tp->snd_cwnd = tp->t_maxseg;
70
71     again:
72     sendalot = 0; /* set nonzero if more than one segment to output */
73
74     /* look for a reason to send a segment; */
75     /* goto send if a segment should be sent */
76
77     218 /*
78     219 * No reason to send a segment, just return.
79     220 */
80     221 return (0);
81
82     222 send:
83
84     /* form output segment, call ip_output() */
85
86     489 if (sendalot)
87     490     goto again;
88     491 return (0);
89     492 }

```

tcp_output.c

Figure 26.1 tcp_output function: overview.

Figure 26.2 shows the first of the tests to determine whether a segment should be sent.

```

71  off = tp->snd_nxt - tp->snd_una;                                     tcp_output.c
72  win = min(tp->snd_wnd, tp->snd_cwnd);
73  flags = tcp_outflags[tp->t_state];
74  /*
75   * If in persist timeout with window of 0, send 1 byte.
76   * Otherwise, if window is small but nonzero
77   * and timer expired, we will send what we can
78   * and go to transmit state.
79   */
80  if (tp->t_force) {
81      if (win == 0) {
82          /*
83           * If we still have some data to send, then
84           * clear the FIN bit. Usually this would
85           * happen below when it realizes that we
86           * aren't sending all the data. However,
87           * if we have exactly 1 byte of unsent data,
88           * then it won't clear the FIN bit below,
89           * and if we are in persist state, we wind
90           * up sending the packet without recording
91           * that we sent the FIN bit.
92           *
93           * We can't just blindly clear the FIN bit,
94           * because if we don't have any more data
95           * to send then the probe will be the FIN
96           * itself.
97           */
98           if (off < so->so_snd.sb_cc)
99               flags ^= TH_FIN;
100          win = 1;
101      } else {
102          tp->t_timer(TCPT_PERSIST) = 0;
103          tp->t_rxtshift = 0;
104      }
105  }

```

Figure 26.2 tcp_output function: data is being forced out.

71-72 off is the offset in bytes from the beginning of the send buffer of the first data byte to send. The first off bytes in the send buffer, starting with snd_una, have already been sent and are waiting to be ACKed.

win is the minimum of the window advertised by the receiver (snd_wnd) and the congestion window (snd_cwnd).

73 The tcp_outflags array was shown in Figure 24.16. The value of this array that is fetched and stored in flags depends on the current state of the connection. flags contains the combination of the TH_ACK, TH_FIN, TH_RST, and TH_SYN flag bits to send to the other end. The other two flag bits, TH_PUSH and TH_URG, will be logically ORED into flags if necessary before the segment is sent.

74-105 The flag `t_force` is set nonzero when the persist timer expires or when out-of-band data is being sent. These two conditions invoke `tcp_output` as follows:

```
tp->t_force = 1;
error = tcp_output(tp);
tp->t_force = 0;
```

This forces TCP to send a segment when it normally wouldn't send anything.

If `win` is 0, the connection is in the persist state (since `t_force` is nonzero). The FIN flag is cleared if there is more data in the socket's send buffer. `win` must be set to 1 byte to force out a single byte.

If `win` is nonzero, out-of-band data is being sent, so the persist timer is cleared and the exponential backoff index, `t_rxtshift`, is set to 0.

Figure 26.3 shows the next part of `tcp_output`, which calculates how much data to send.

```

106     len = min(so->so_snd.sb_cc, win) - off;
107     if (len < 0) {
108         /*
109          * If FIN has been sent but not acked,
110          * but we haven't been called to retransmit,
111          * len will be -1. Otherwise, window shrank
112          * after we sent into it. If window shrank to 0,
113          * cancel pending retransmit and pull snd_next
114          * back to (closed) window. We will enter persist
115          * state below. If the window didn't close completely,
116          * just wait for an ACK.
117          */
118         len = 0;
119         if (win == 0) {
120             tp->t_timer[TCPT_REXMT] = 0;
121             tp->snd_next = tp->snd_una;
122         }
123     }
124     if (len > tp->t_maxseg) {
125         len = tp->t_maxseg;
126         sendalot = 1;
127     }
128     if (SEQ_LT(tp->snd_next + len, tp->snd_una + so->so_snd.sb_cc))
129         flags |= TH_FIN;
130     win = sbospace(&so->su_fcvt);

```

tcp_output.c

Figure 26.3 `tcp_output` function: calculate how much data to send.

Calculate amount of data to send

106 `len` is the minimum of the number of bytes in the send buffer and `win` (which is the minimum of the receiver's advertised window and the congestion window, perhaps 1 byte if output is being forced). `off` is subtracted because that many bytes at the beginning of the send buffer have already been sent and are awaiting acknowledgment.

Check for window shrink

107-117

One way for `len` to be less than 0 occurs if the receiver *shrinks* the window, that is, the receiver moves the right edge of the window to the left. The following example demonstrates how this can happen. First the receiver advertises a window of 6 bytes and TCP transmits a segment with bytes 4, 5, and 6. TCP immediately transmits another segment with bytes 7, 8, and 9. Figure 26.4 shows the status of our end after the two segments are sent.

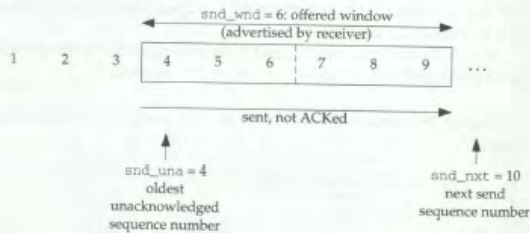


Figure 26.4 Send buffer after bytes 4 through 9 are sent.

Then an ACK is received with an acknowledgment field of 7 (acknowledging all data up through and including byte 6) but with a window of 1. The receiver has shrunk the window, as shown in Figure 26.5.

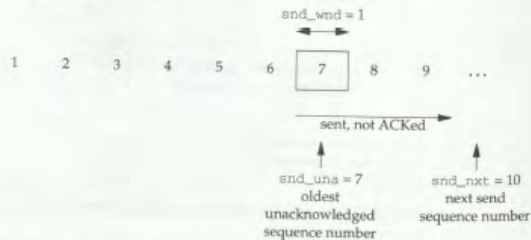


Figure 26.5 Send buffer after receiving acknowledgment of bytes 4 through 6.

Performing the calculations in Figures 26.2 and 26.3, after the window is shrunk, we have

```
off = snd_nxt - snd_una = 10 - 7 = 3
win = 1
len = min(so_snd, sb_cc, win) - off = min(3, 1) - 3 = -2
```

assuming the send buffer contains only bytes 7, 8, and 9.

Both RFC 793 and RFC 1122 strongly discourage shrinking the window. Nevertheless, implementations must be prepared for this. Handling scenarios such as this comes under the *Robustness Principle*, first mentioned in RFC 791: "Be liberal in what you accept, and conservative in what you send."

Another way for `len` to be less than 0 occurs if the FIN has been sent but not acknowledged and not retransmitted. (See Exercise 26.2.) We show this in Figure 26.6.

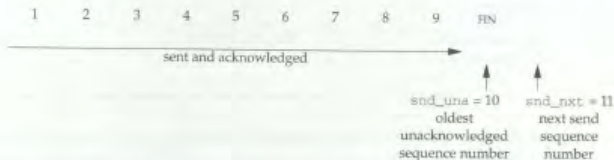


Figure 26.6 Bytes 1 through 9 have been sent and acknowledged, and then connection is closed.

This figure continues Figure 26.4, assuming the final segment with bytes 7, 8, and 9 is acknowledged, which sets `snd_una` to 10. The process then closes the connection, causing the FIN to be sent. We'll see later in this chapter that when the FIN is sent, `snd_next` is incremented by 1 (since the FIN takes a sequence number), which in this example sets `snd_next` to 11. The sequence number of the FIN is 10. Performing the calculations in Figures 26.2 and 26.3, we have

```
off = snd_next - snd_una = 11 - 10 = 1
win = 6
len = min(so_snd.sb_cc, win) - off = min(0, 6) - 1 = -1
```

We assume that the receiver advertises a window of 6, which makes no difference, since the number of bytes in the send buffer (0) is less than this.

Enter persist state

118-122 `len` is set to 0. If the advertised window is 0, any pending retransmission is canceled by setting the retransmission timer to 0. `snd_next` is also pulled to the left of the window by setting it to the value of `snd_una`. The connection will enter the persist state later in this function, and when the receiver finally opens its window, TCP starts retransmitting from the left of the window.

Send one segment at a time

124-127 If the amount of data to send exceeds one segment, `len` is set to a single segment and the `sendallot` flag is set to 1. As shown in Figure 26.1, this causes another loop through `tcp_output` after the segment is sent.

Turn off FIN flag if send buffer not emptied

128-129 If the send buffer is not being emptied by this output operation, the FIN flag must be cleared (in case it is set in `flags`). Figure 26.7 shows an example of this.

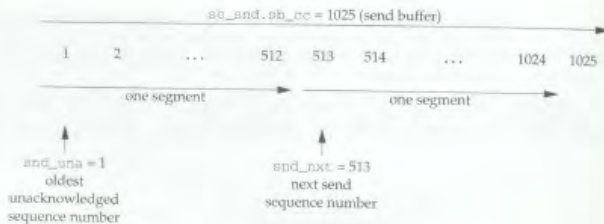


Figure 26.7 Example of send buffer not being emptied when FIN is set.

In this example the first 512-byte segment has already been sent (and is waiting to be acknowledged) and TCP is about to send the next 512-byte segment (bytes 512–1024). There is still 1 byte left in the send buffer (byte 1025) and the process closes the connection. `len` equals 512 (one segment), and the C expression becomes

```
SEQ_LT(1025, 1026)
```

which is true, so the FIN flag is cleared. If the FIN flag were mistakenly left on, TCP couldn't send byte 1025 to the receiver.

Calculate window advertisement

`win` is set to the amount of space available in the receive buffer, which becomes TCP's window advertisement to the other end. Be aware that this is the second use of this variable in this function. Earlier it contained the maximum amount of data TCP could send, but for the remainder of this function it contains the receive window advertised by this end of the connection.

The silly window syndrome (called SWS and described in Section 22.3 of Volume 1) occurs when small amounts of data, instead of full-sized segments, are exchanged across a connection. It can be caused by a receiver who advertises small windows and by a sender who transmits small segments. Correct avoidance of the silly window syndrome must be performed by both the sender and the receiver. Figure 26.8 shows silly window avoidance by the sender.

Sender silly window avoidance

142-143 If a full-sized segment can be sent, it is sent.

144-146 If an ACK is not expected (`idle` is true), or if the Nagle algorithm is disabled (`TF_NODELAY` is true) and TCP is emptying the send buffer, the data is sent. The Nagle algorithm (Section 19.4 of Volume 1) prevents TCP from sending less than a full-sized segment when an ACK is expected for the connection. It can be disabled using the `TCP_NODELAY` socket option. For a normal interactive connection (e.g., Telnet or Rlogin), if there is unacknowledged data, this `if` statement is false, since the Nagle algorithm is enabled by default.

147-148 If output is being forced by either the persist timer or sending out-of-band data, some data is sent.


```

131  /* tcp_output.c
132  * Sender silly window avoidance. If connection is idle
133  * and can send all data, a maximum segment,
134  * at least a maximum default-sized segment do it,
135  * or are forced, do it; otherwise don't bother.
136  * If peer's buffer is tiny, then send
137  * when window is at least half open.
138  * If retransmitting (possibly after persist timer forced us
139  * to send into a small window), then must resend.
140  */
141  if (len) {
142      if (len == tp->t_maxseg)
143          goto send;
144      if ((idle || tp->t_flags & TF_NODELAY) &&
145          len + off >= so->so_snd.sb_cc)
146          goto send;
147      if (tp->t_force)
148          goto send;
149      if (len >= tp->max_sndwnd / 2)
150          goto send;
151      if (SEQ_LT(tp->snd_nxt, tp->snd_max))
152          goto send;
153  }
tcp_output.c

```

Figure 26.8 tcp_output function: sender silly window avoidance.

143-150 If the receiver's window is at least half open, data is sent. This is to deal with peers that always advertise tiny windows, perhaps smaller than the segment size. The variable `max_sndwnd` is calculated by `tcp_input` as the largest window advertisement ever advertised by the other end. It is an attempt to guess the size of the other end's receive buffer and assumes the other end never reduces the size of its receive buffer.

151-152 If the retransmission timer expired, then a segment must be sent. `snd_max` is the highest sequence number that has been transmitted. We saw in Figure 25.26 that when the retransmission timer expires, `snd_nxt` is set to `snd_una`, that is, `snd_nxt` is moved to the left edge of the window, making it less than `snd_max`.

The next portion of `tcp_output`, shown in Figure 26.9, determines if TCP must send a segment just to advertise a new window to the other end. This is called a *window update*.

154-168 The expression

```
min|win, (long)TCP_MAXWIN << tp->rcv_scale)
```

is the smaller of the amount of available space in the socket's receive buffer (`win`) and the maximum size of the window allowed for this connection. This is the maximum window TCP can currently advertise to the other end. The expression

```
(tp->rcv_adv - tp->rcv_nxt)
```

is the number of bytes remaining in the last window advertisement that TCP sent to the other end. Subtracting this from the maximum window yields `adv`, the number of

```

154  /*
155  * Compare available window to amount of window
156  * known to peer (as advertised window less
157  * next expected input). If the difference is at least two
158  * max size segments, or at least 50% of the maximum possible
159  * window, then want to send a window update to peer.
160  */
161  if (win > 0) {
162      /*
163       * 'adv' is the amount we can increase the window,
164       * taking into account that we are limited by
165       * TCP_MAXWIN << tp->rcv_scale.
166       */
167      long adv = min(win, (long) TCP_MAXWIN << tp->rcv_scale) -
168      (tp->rcv_adv - tp->rcv_next);
169
170      if (adv >= (long) (2 * tp->t_maxseg))
171          goto send;
172      if (2 * adv >= (long) so->so_rcv.sb_hiwat)
173          goto send;
174  }

```

tcp_output.c

tcp_output.c

Figure 26.9 `tcp_output` function: check if a window update should be sent.

bytes by which the window has opened. `rcv_next` is incremented by `tcp_input` when data is received in sequence, and `rcv_adv` is incremented by `tcp_output` in Figure 26.32 when the edge of the advertised window moves to the right.

Consider Figure 24.18 and assume that a segment with bytes 4, 5, and 6 is received and that these three bytes are passed to the process. Figure 26.10 shows the state of the receive space at this point in `tcp_output`.

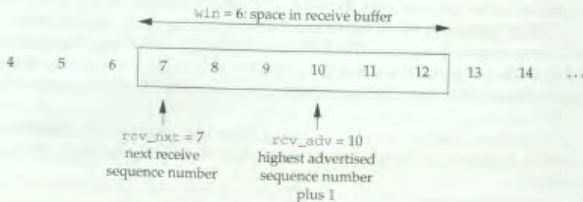


Figure 26.10 Transition from Figure 24.18 after bytes 4, 5, and 6 are received.

The value of `adv` is 3, since there are 3 more bytes of the receive space (bytes 10, 11, and 12) for the other end to fill.

169-170 If the window has opened by two or more segments, a window update is sent. When data is received as full-sized segments, this code causes every other received

segment to be acknowledged: TCP's ACK-every-other-segment property. (We show an example of this shortly.)

174-178 If the window has opened by at least 50% of the maximum possible window (the socket's receive buffer high-water mark), a window update is sent.

The next part of `tcp_output`, shown in Figure 26.11, checks whether various flags require TCP to send a segment.

```

174  /*----- tcp_output.c
175  * Send if we owe peer an ACK.
176  */
177  if (tp->t_flags & TF_ACKNOW)
178      goto send;
179  if (flags & (TH_SYN | TH_RST))
180      goto send;
181  if (SEQ_GT(tp->snd_up, tp->snd_una))
182      goto send;
183  /*
184  * If our state indicates that FIN should be sent
185  * and we have not yet done so, or we're retransmitting the FIN,
186  * then we need to send.
187  */
188  if (flags & TH_FIN &&
189      {(tp->t_flags & TF_SENTFIN) == 0 || tp->snd_next == tp->snd_una})
190      goto send;
----- tcp_output.c

```

Figure 26.11 `tcp_output` function: should a segment should be sent?

174-178 If an immediate ACK is required, a segment is sent. The `TF_ACKNOW` flag is set by various functions: when the 200-ms delayed ACK timer expires, when a segment is received out of order (for the fast retransmit algorithm), when a SYN is received during the three-way handshake, when a persist probe is received, and when a FIN is received.

179-180 If `flags` specifies that a SYN or RST should be sent, a segment is sent.

181-182 If the urgent pointer, `snd_up`, is beyond the start of the send buffer, a segment is sent. The urgent pointer is set by the `PRU_SENDOOB` request (Figure 30.9).

183-190 If `flags` specifies that a FIN should be sent, a segment is sent only if the FIN has not already been sent, or if the FIN is being retransmitted. The flag `TF_SENTFIN` is set later in this function when the FIN is sent.

At this point in `tcp_output` there is no need to send a segment. Figure 26.12 shows the final piece of code before `tcp_output` returns.

191-217 If there is data in the send buffer to send (`so_snd.sb_cc` is nonzero) and both the retransmission timer and the persist timer are off, turn the persist timer on. This scenario happens when the window advertised by the other end is too small to receive a full-sized segment, and there is no other reason to send a segment.

218-221 `tcp_output` returns, since there is no reason to send a segment.

```

191  /* tcp_output.c
192  * TCP window updates are not reliable, rather a polling protocol
193  * using 'persist' packets is used to ensure receipt of window
194  * updates. The three 'states' for the output side are:
195  *   idle           not doing retransmits or persists
196  *   persisting     to move a small or zero window
197  *   (re)transmitting and thereby not persisting
198  *
199  * tp->t_timer[TCPT_PERSIST]
200  *   is set when we are in persist state.
201  * tp->t_force
202  *   is set when we are called to send a persist packet.
203  * tp->t_timer[TCPT_REXMT]
204  *   is set when we are retransmitting
205  * The output side is idle when both timers are zero.
206  *
207  * If send window is too small, there is data to transmit, and no
208  * retransmit or persist is pending, then go to persist state.
209  * If nothing happens soon, send when timer expires:
210  * If window is nonzero, transmit what we can,
211  * otherwise force out a byte.
212  */
213  if (so->so_snd.sb_cc && tp->t_timer[TCPT_REXMT] == 0 &&
214      tp->t_timer[TCPT_PERSIST] == 0) {
215      tp->t_rxtshift = 0;
216      tcp_setpersist(tp);
217  }
218  /*
219  * No reason to send a segment, just return.
220  */
221  return (0);
tcp_output.c

```

Figure 26.12 tcp_output function: enter persist state.

Example

A process writes 100 bytes, followed by a write of 50 bytes, on an idle connection. Assume a segment size of 512 bytes. When the first write occurs, the code in Figure 26.8 (lines 144–146) sends a segment with 100 bytes of data since the connection is idle and TCP is emptying the send buffer.

When 50-byte write occurs, the code in Figure 26.8 does not send a segment: the amount of data is not a full-sized segment, the connection is not idle (assume TCP is awaiting the ACK for the 100 bytes that it just sent), the Nagle algorithm is enabled by default, `t_force` is not set, and assuming a typical receive window of 4096, 50 is not greater than or equal to 2048. These 50 bytes remain in the send buffer, probably until the ACK for the 100 bytes is received. This ACK will probably be delayed by the other end, causing more delay in sending the final 50 bytes.

This example shows the timing delays that can occur when sending less than full-sized segments with the Nagle algorithm enabled. See also Exercise 26.12.

Example

This example demonstrates the ACK-every-other-segment property of TCP. Assume a connection is established with a segment size of 1024 bytes and a receive buffer size of 4096. There is no data to send—TCP is just receiving.

A window of 4096 is advertised in the ACK of the SYN, and Figure 26.13 shows the two variables `rcv_nxt` and `rcv_adv`. The receive buffer is empty.

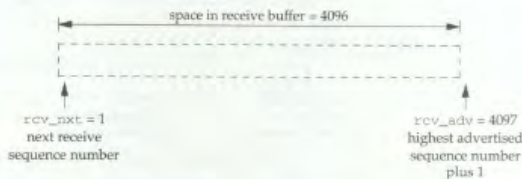


Figure 26.13 Receiver advertising a window of 4096.

The other end sends a segment with bytes 1–1024. `tcp_input` processes the segment, sets the delayed-ACK flag for the connection, and appends the 1024 bytes of data to the socket's receiver buffer (Figure 28.13). `rcv_nxt` is updated as shown in Figure 26.14.

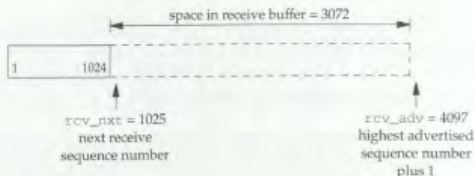


Figure 26.14 Transition from Figure 26.13 after bytes 1–1024 received.

The process reads the 1024 bytes in its socket receive buffer. We'll see in Figure 30.6 that the resulting `PRU_RCVD` request causes `tcp_output` to be called, because a window update might need to be sent after the process reads data from the receive buffer. When `tcp_output` is called, the two variables still have the values shown in Figure 26.14 and the only difference is that the amount of space in the receive buffer has increased to 4096 since the process has read the first 1024 bytes. The calculations in Figure 26.9 are performed:

$$\begin{aligned} \text{adv} &= \min(4096, 65535) - (4097 - 1025) \\ &= 1024 \end{aligned}$$

`TCP_MAXWIN` is 65535 and we assume a receive window scale shift of 0. Since the window has increased by less than two segments (2048), nothing is sent. But the delayed-ACK flag is still set, so if the 200-ms timer expires, an ACK will be sent.

When TCP receives the next segment with bytes 1025–2048, `tcp_input` processes the segment, sets the delayed-ACK flag for the connection (which was already on), and appends the 1024 bytes of data to the socket's receiver buffer. `rcv_nxt` is updated as shown in Figure 26.15.

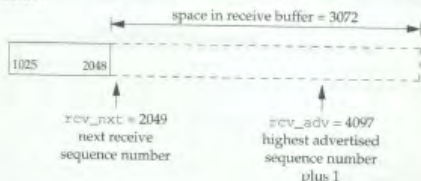


Figure 26.15 Transition from Figure 26.14 after bytes 1025–2048 received.

The process reads bytes 1025–2048 and `tcp_output` is called. The two variables still have the values shown in Figure 26.15, although the space in the receive buffer increases to 4096 when the process reads the 1024 bytes of data. The calculations in Figure 26.9 are performed:

$$\begin{aligned} \text{adv} &= \min(4096, 65535) - (4097 - 2049) \\ &= 2048 \end{aligned}$$

This value is now greater than or equal to two segments, so a segment is sent with an acknowledgment field of 2049 and an advertised window of 4096. This is a window update. The receiver is willing to receive bytes 2049 through 6145. We'll see later in this function that when this segment is sent, the value of `rcv_adv` also gets updated to 6145.

This example shows that when receiving data faster than the 200-ms delayed ACK timer, an ACK is sent when the receive window changes by more than two segments due to the process reading the data. If data is received for the connection but the process is not reading the data from the socket's receive buffer, the ACK-every-other-segment property won't occur. Instead the sender will only see the delayed ACKs, each advertising a smaller window, until the receive buffer is filled and the window goes to 0.

26.4 TCP Options

The TCP header can contain options. We digress to discuss these options since the next piece of `tcp_output` decides which options to send and constructs the options in the outgoing segment. Figure 26.16 shows the format of the options supported by Net/3.

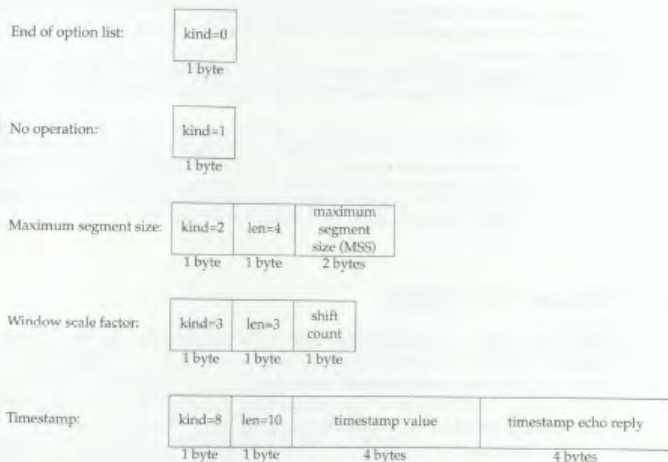


Figure 26.16 TCP options supported by Net/3.

Every option begins with a 1-byte *kind* that specifies the type of option. The first two options (with *kinds* of 0 and 1) are single-byte options. The other three are multi-byte options with a *len* byte that follows the *kind* byte. The length is the total length, including the *kind* and *len* bytes.

The multi-byte integers—the MSS and the two timestamp values—are stored in network byte order.

The final two options, window scale and timestamp, are new and therefore not supported by many systems. To provide interoperability with these older systems, the following rules apply.

1. TCP can send one of these options (or both) with the initial SYN segment corresponding to an active open (that is, a SYN without an ACK). Net/3 does this for both options if the global `tcp_do_rfc1323` is nonzero (it defaults to 1). This is done in `tcp_newtcpcb`.
2. The option is enabled only if the SYN reply from the other end also includes the desired option. This is handled in Figures 28.20 and 29.2.
3. If TCP performs a passive open and receives a SYN specifying the option, the response (the SYN plus ACK) must contain the option if TCP wants to enable the option. This is done in Figure 26.23.

Since a system must ignore options that it doesn't understand, the newer options are enabled by both ends only if both ends understand the option and both ends want the option enabled.

The processing of the MSS option is covered in Section 27.5. The next two sections summarize the Net/3 handling of the two newer options: window scale and timestamp.

Other options have been proposed. *kinds* of 4, 5, 6, and 7, called the selective-ACK and echo options, are defined in RFC 1072 [Jacobson and Braden 1988]. We don't show them in Figure 26.16 because the echo options were replaced with the timestamp option, and selective ACKs, as currently defined, are still under discussion and were not included in RFC 1323. Also, the T/TCP proposal for TCP transactions (RFC 1644 [Braden 1994], and Section 24.7 of Volume 1) specifies three options with *kinds* of 11, 12, and 13.

26.5 Window Scale Option

The window scale option, defined in RFC 1323, avoids the limitation of a 16-bit window size field in the TCP header (Figure 24.10). Larger windows are required for what are called *long fat pipes*, networks with either a high bandwidth or a long delay (i.e., a long RTT). Section 24.3 of Volume 1 gives examples of current networks that require larger windows to obtain maximum TCP throughput.

The 1-byte shift count in Figure 26.16 is between 0 (no scaling performed) and 14. This maximum value of 14 provides a maximum window of 1,073,725,440 bytes (65535×2^{14}). Internally Net/3 maintains window sizes as 32-bit values, not 16-bit values.

The window scale option can only appear in a SYN segment; therefore the scale factor is fixed in each direction when the connection is established.

The two variables `snd_scale` and `rcv_scale` in the TCP control block specify the shift count for the send window and the receive window, respectively. Both default to 0 for no scaling. Every 16-bit advertised window received from the other end is left shifted by `snd_scale` bits to obtain the real 32-bit advertised window size (Figure 28.6). Every time TCP sends a window advertisement to the other end, the internal 32-bit window size is right shifted by `rcv_scale` bits to give the value that is placed into the TCP header (Figure 26.29).

When TCP sends a SYN, either actively or passively, it chooses the value of `rcv_scale` to request, based on the size of the socket's receive buffer (Figures 28.7 and 30.4).

26.6 Timestamp Option

The timestamp option is also defined in RFC 1323 and lets the sender place a timestamp in every segment. The receiver sends the timestamp back in the acknowledgment, allowing the sender to calculate the RTT for each received ACK. Figure 26.17 summarizes the timestamp option and the variables involved.

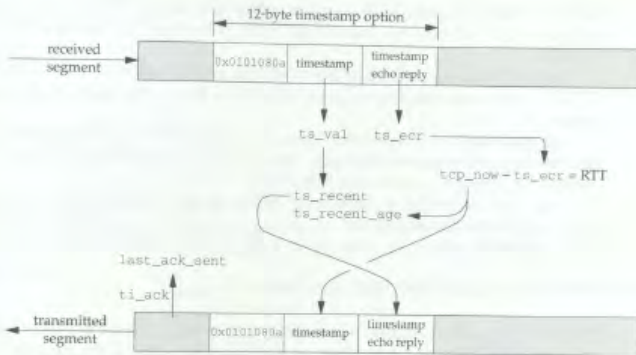


Figure 26.17 Summary of variables used with timestamp option.

The global variable `tcp_now` is the timestamp clock. It is initialized to 0 when the kernel is initialized and incremented by 1 every 500 ms (Figure 25.8). Three variables are maintained in the TCP control block for the timestamp option:

- `ts_recent` is a copy of the most-recent valid timestamp from the other end. (We describe shortly what makes a timestamp “valid.”)
- `ts_recent_age` is the value of `tcp_now` when `ts_recent` was last copied from a received segment.
- `last_ack_sent` is the value of the acknowledgment field (`ti_ack`) the last time a segment was sent (Figure 26.32). This is normally equal to `recv_nxt`, the next expected sequence number, unless ACKs are delayed.

The two variables `ts_val` and `ts_ecr` are local variables in the function `tcp_input` that contain the two values from the timestamp option.

- `ts_val` is the timestamp sent by the other end with its data.
- `ts_ecr` is the timestamp from the segment that is being acknowledged by the received segment.

In an outgoing segment, the first 4 bytes of the timestamp option are set to `0x0101080a`. This is the recommended value from Appendix A of RFC 1323. The 2 bytes of 1 are NOPs from Figure 26.16, followed by a *kind* of 8 and a *len* of 10, which identify the timestamp option. By placing two NOPs in front of the option, the two 32-bit timestamps in the option and the data that follows are aligned on 32-bit boundaries. Also, we show the received timestamp option in Figure 26.17 with the recommended 12-byte format (which Net/3 always generates), but the code that processes

received options (Figure 28.10) does not require this format. The 10-byte format shown in Figure 26.16, without two preceding NOPs, is handled fine on input (but see Exercise 28.4).

The RIT of a transmitted segment and its ACK is calculated as `tcp_now` minus `ts_echo`. The units are 500-ms clock ticks, since that is the units of the Net/3 timestamps.

The presence of the timestamp option also allows TCP to perform PAWS: protection against wrapped sequence numbers. We describe this algorithm in Section 28.7. The variable `ts_recent_age` is used with PAWS.

`tcp_output` builds a timestamp option in an outgoing segment by copying `tcp_now` into the timestamp and `ts_recent` into the echo reply (Figure 26.24). This is done for every segment when the option is in use, unless the RST flag is set.

Which Timestamp to Echo, RFC 1323 Algorithm

The test for a valid timestamp determines whether the value in `ts_recent` is updated, and since this value is always sent as the timestamp echo reply, the test for validity determines which timestamp gets echoed back to the other end. RFC 1323 specified the following test:

```
ti_seq <= last_ack_sent < ti_seq + ti_len
```

which is implemented in C as shown in Figure 26.18.

```
if (ts_present && SEQ_LEQ(ti->ti_seq, tp->last_ack_sent) &&
    SEQ_LT(tp->last_ack_sent, ti->ti_seq + ti->ti_len)) {
    tp->ts_recent_age = tcp_now;
    tp->ts_recent = ts_val;
}
```

Figure 26.18 Typical code to determine if received timestamp is valid.

The variable `ts_present` is true if a timestamp option was received in the segment. We encounter this code twice in `tcp_input`: Figure 28.11 does the test in the header prediction code, and Figure 28.35 does the test in the normal input processing.

To see what this test is doing, Figure 26.19 shows five different scenarios, corresponding to five different segments received on a connection. In each scenario `ti_len` is 3.

The left edge of the receive window begins with sequence number 4. In scenario 1 the segment contains completely duplicate data. The `SEQ_LEQ` test in Figure 28.11 is true, but the `SEQ_LT` test fails. For scenarios 2, 3, and 4, both the `SEQ_LEQ` and `SEQ_LT` tests are true because the left edge of the window is advanced by any one of these three segments, even though scenario 2 contains two duplicate bytes of data, and scenario 3 contains one duplicate byte of data. Scenario 5 fails the `SEQ_LEQ` test, because it doesn't advance the left edge of the window. This segment is one in the future that's not the next expected, implying that a previous segment was lost or reordered.

Unfortunately this test to determine whether to update `ts_recent` is flawed [Braden 1993]. Consider the following example.

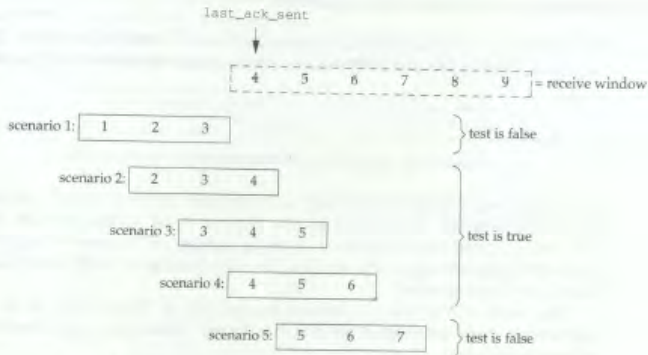


Figure 26.19 Example receive window and five different scenarios of received segment.

1. In Figure 26.19 a segment that we don't show arrives with bytes 1, 2, and 3. The timestamp in this segment is saved in `ts_recent` because `last_ack_sent` is 1. An ACK is sent with an acknowledgment field of 4, and `last_ack_sent` is set to 4 (the value of `rcv_next`). We have the receive window shown in Figure 26.19.
2. This ACK is lost.
3. The other end times out and retransmits the segment with bytes 1, 2, and 3. This segment arrives and is the one labeled "scenario 1" in Figure 26.19. Since the `SEQ_LT` test in Figure 26.18 fails, `ts_recent` is not updated with the value from the retransmitted segment.
4. A duplicate ACK is sent with an acknowledgment field of 4, but the timestamp echo reply is `ts_recent`, the value copied from the segment in step 1. But when the receiver calculates the RTT using this value, it will (incorrectly) take into account the original transmission, the lost ACK, the timeout, the retransmission, and the duplicate ACK.

For correct RTT estimation by the other end, the timestamp value from the retransmission should be returned in the duplicate ACK.

The tests in Figure 26.18 also fail to update `ts_recent` if the length of the received segment is 0, since the left edge of the window is not moved. This incorrect test can also lead to problems with long-lived (greater than 24 days, the PAWS limit described in Section 28.7), unidirectional connections (all the data flow is in one direction so the sender of the data always sends the same ACKs).

Which Timestamp to Echo, Corrected Algorithm

The algorithm we'll encounter in the Net/3 sources is from Figure 26.18. The correct algorithm given in [Braden 1993] replaces Figure 26.18 with the one in Figure 26.20.

```
if (ts_present && TSTMP_GEQ(ts_val, tp->ts_recent) &&
    SEQ_GEQ(ti->ti_seq, tp->last_ack_sent)) {
```

Figure 26.20 Correct code to determine if received timestamp is valid.

This doesn't test whether the left edge of the window moves or not, it just verifies that the new timestamp (`ts_val`) is greater than or equal to the previous timestamp (`ts_recent`), and that the starting sequence number of the received segment is not greater than the left edge of the window. Scenario 5 in Figure 26.20 would fail this new test since it is out of order.

The macro `TSTMP_GEQ` is identical to `SEQ_GEQ` in Figure 24.21. It is used with timestamps, since timestamps are 32-bit unsigned values that wrap around just like sequence numbers.

Timestamps and Delayed ACKs

It is constructive to see how timestamps and RTT calculations are affected by delayed ACKs. Recall from Figure 26.17 that the value saved by TCP in `ts_recent` becomes the echoed timestamp in segments that are sent, which are used by the other end in calculating its RTT. When ACKs are delayed, the delay time should be taken into account by the side that sees the delays, or else it might retransmit too quickly. In the example that follows we only consider the code in Figure 26.20, but the incorrect code in Figure 26.18 also handles delayed ACKs correctly.

Consider the receive sequence space in Figure 26.21 when the received segment contains bytes 4 and 5.

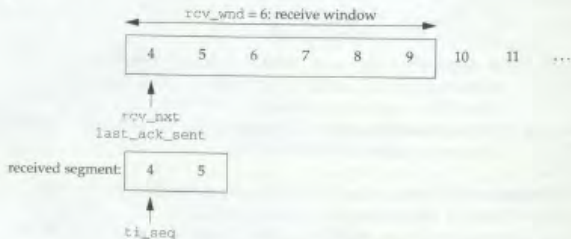


Figure 26.21 Receive sequence space when segment with bytes 4 and 5 arrives.

Since `ti_seq` is less than or equal to `last_ack_sent`, `ts_recent` is copied from the segment. `rcv_nxt` is also increased by 2.

Assume that the ACK for these 2 bytes is delayed, and before that delayed ACK is sent, the next in-order segment arrives. This is shown in Figure 26.22.

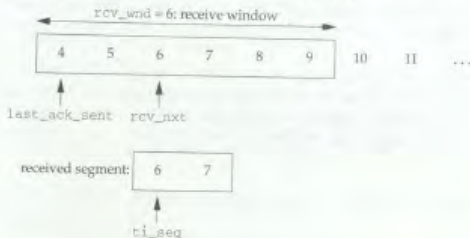


Figure 26.22 Receive sequence space when segment with bytes 6 and 7 arrives.

This time `ti_seq` is greater than `last_ack_sent`, so `ts_recent` is not updated. This is intentional. Assuming TCP now sends an ACK for sequence numbers 4–7, the other end's RTT will take into account the delayed ACK, since the echoed timestamp (Figure 26.24) is the one from the segment with sequence numbers 4 and 5. These figures also demonstrate that `rcv_nxt` equals `last_ack_sent` except when ACKs are delayed.

26.7 Send a Segment

The last half of `tcp_output` sends the segment—it fills in all the fields in the TCP header and passes the segment to IP for output.

Figure 26.23 shows the first part, which sends the MSS and window scale options with a SYN segment.

223-234 The TCP options are built in the array `opt`, and the integer `optlen` keeps a count of the number of bytes accumulated (since multiple options can be sent at once). If the SYN flag bit is set, `snd_nxt` is set to the initial send sequence number (`iss`). If TCP is performing an active open, `iss` is set by the `PRU_CONNECT` request when the TCP control block is created. If this is a passive open, `tcp_input` creates the TCP control block and sets `iss`. In both cases, `iss` is set from the global `tcp_iss`.

235 The flag `TF_NOOPT` is checked, but this flag is never enabled and there is no way to turn it on. Hence, the MSS option is always sent with a SYN segment.

In the Net/1 version of `tcp_newtcpcb`, the comment “send options!” appeared on the line that initialized `t_flags` to 0. The `TF_NOOPT` flag is probably a historical artifact from a pre-Net/1 system that had problems interoperating with other hosts when it sent the MSS option, so the default was to not send the option.

```

223  /*
224  * Before ESTABLISHED, force sending of initial options
225  * unless TCP set not to do any options.
226  * NOTE: we assume that the IP/TCP header plus TCP options
227  * always fit in a single mbuf, leaving room for a maximum
228  * link header, i.e.
229  * max_linkhdr + sizeof(struct tcpiphdr) + optlen <= MHLEN
230  */
231  optlen = 0;
232  hdrhlen = sizeof(struct tcpiphdr);
233  if (flags & TH_SYN) {
234      tp->snd_next = tp->iss;
235      if ((tp->t_flags & TF_NOOPT) == 0) {
236          u_short mss;
237
238          opt[0] = TCPOPT_MAXSEG;
239          opt[1] = 4;
240          mss = htons((u_short) tcp_mss(tp, 0));
241          bcopy((caddr_t) & mss, (caddr_t) (opt + 2), sizeof(mss));
242          optlen = 4;
243
244          if ((tp->t_flags & TF_REQ_SCALE) &&
245              ((flags & TH_ACK) == 0)) {
246              (tp->t_flags & TF_RCVD_SCALE)) {
247                  *((u_long *) (opt + optlen)) = htonl(TCPOPT_NOP << 24 |
248                                                         TCPOPT_WINDOW << 16 |
249                                                         TCPOLEN_WINDOW << 8 |
250                                                         tp->request_r_scale);
251                  optlen += 4;
252              }
253          }
254      }
255  }

```

Figure 26.23 tcp_output function: send options with first SYN segment.

Build MSS option

236-241 opt[0] is set to 2 (TCPOPT_MAXSEG) and opt[1] is set to 4, the length of the MSS option in bytes. The function `tcp_mss` calculates the MSS to announce to the other end; we cover this function in Section 27.5. The 16-bit MSS is stored in `opt[2]` and `opt[3]` by `bcopy` (Exercise 26.5). Notice that Net/3 always sends an MSS announcement with the SYN for a connection.

Should window scale option be sent?

242-244 If TCP is to request the window scale option, this option is sent only if this is an active open (`TH_ACK` is not set) or if this is a passive open and the window scale option was received in the SYN from the other end. Recall that `t_flags` was set to `TF_REQ_SCALE|TF_REQ_TSTMP` when the TCP control block was created in Figure 25.21, if the global variable `tcp_do_rfc1323` was nonzero (its default value).

Build window scale option

249-249 Since the window scale option occupies 3 bytes (Figure 26.16), a 1-byte NOP is stored before the option, forcing the option length to be 4 bytes. This causes the data in the segment that follows the options to be aligned on a 4-byte boundary. If this is an active open, `request_t_scale` is calculated by the `PRU_CONNECT` request. If this is a passive open, the window scale factor is calculated by `tcp_input` when the SYN is received.

RFC 1323 specifies that if TCP is prepared to scale windows it should send this option even if its own shift count is 0. This is because the option serves two purposes: to notify the other end that it supports the option, and to announce its shift count. Even though TCP may calculate its own shift count as 0, the other end might want to use a different value.

The next part of `tcp_output` is shown in Figure 26.24. It finishes building the options in the outgoing segment.

```

253      /*
254      * Send a timestamp and echo-reply if this is a SYN and our side
255      * wants to use timestamps (TF_REQ_TSTMP is set) or both our side
256      * and our peer have sent timestamps in our SYN's.
257      */
258      if ((tp->t_flags & (TF_REQ_TSTMP | TF_NOOPT)) == TF_REQ_TSTMP &&
259          (flegs & TH_RST) == 0 &&
260          ((flegs & (TH_SYN | TH_ACK)) == TH_SYN ||
261           {tp->t_flags & TF_RCVD_TSTMP})) {
262          u_long *lp = (u_long *) (opt + optlen);
263
264          /* Form timestamp option as shown in appendix A of RFC 1323. */
265          *lp++ = htonl(TCPOPT_TSTAMP_HDR);
266          *lp++ = htonl(tcp_now);
267          *lp = htonl(tp->ts_recent);
268          optlen += TCPOLEN_TSTAMP_APPA;
269      }
270      hdrlen += optlen;
271
272      /*
273      * Adjust data length if insertion of options will
274      * bump the packet length beyond the t_maxseg length.
275      */
276      if (len > tp->t_maxseg - optlen) {
277          len = tp->t_maxseg - optlen;
278          sendalot = 1;
279      }

```

Figure 26.24 `tcp_output` function: finish sending options.

Should timestamp option be sent?

253-261 If the following three conditions are all true, a timestamp option is sent: (1) TCP is configured to request the timestamp option, (2) the segment being formed does not contain the RST flag, and (3) either this is an active open (i.e., `flags` specifies the SYN flag

but not the ACK flag) or TCP has received a timestamp from the other end (`TF_RCVD_TSTMP`). Unlike the MSS and window scale options, a timestamp option can be sent with every segment once both ends agree to use the option.

Build timestamp option

283-287 The timestamp option (Section 26.6) consists of 12 bytes (`TCPOLEN_TSTAMP_APPA`). The first 4 bytes are `0x0101080a` (the constant `TCPOPT_TSTAMP_HDR`), as described with Figure 26.17. The timestamp value is taken from `tcp_now` (the number of 500-ms clock ticks since the system was initialized), and the timestamp echo reply is taken from `ts_recent`, which is set by `tcp_input`.

Check if options have overflowed segment

270-277 The size of the TCP header is incremented by the number of option bytes (`opt_len`). If the amount of data to send (`len`) exceeds the MSS minus the size of the options (`optlen`), the data length is decreased accordingly and the `sendatol` flag is set, to force another loop through this function after this segment is sent (Figure 26.1).

The MSS and window scale options only appear in SYN segments, which Net/3 always sends without data, so this adjustment of the data length doesn't apply. When the timestamp option is in use, however, it appears in all segments. This reduces the amount of data in each full-sized data segment from the announced MSS to the announced MSS minus 12 bytes.

The next part of `tcp_output`, shown in Figure 26.25, updates some statistics and allocates an mbuf for the IP and TCP headers. This code is executed when the segment being output contains some data (`len` is greater than 0).

Update statistics

284-292 If `tcp_force` is nonzero and TCP is sending a single byte of data, this is a window probe. If `snd_next` is less than `snd_max`, this is a retransmission. Otherwise, this is normal data transmission.

Allocate an mbuf for IP and TCP headers

293-297 An mbuf with a packet header is allocated by `MGETHDR`. This is for the IP and TCP headers, and possibly the data (if there's room). Although `tcp_output` is often called as part of a system call (e.g., `write`) it is also called at the software interrupt level by `tcp_input`, and as part of the timer processing. Therefore `M_DONTWAIT` is specified. If an error is returned, a jump is made to the label `out`. This label is near the end of the function, in Figure 26.32.

Copy data into mbuf

298-308 If the amount of data is less than 44 bytes (`100 - 40 - 16`, assuming no TCP options), the data is copied directly from the socket send buffer into the new packet header mbuf by `m_copydata`. Otherwise `m_copy` creates a new mbuf chain with the data from the socket send buffer and this chain is linked to the new packet header mbuf. Recall our description of `m_copy` in Section 2.9, where we showed that if the data is in a cluster, `m_copy` just references that cluster and doesn't make a copy of the data.


```

278  /*
279  * Grab a header mbuf, attaching a copy of data to
280  * be transmitted, and initialize the header from
281  * the template for sends on this connection.
282  */
283  if (len) {
284      if (tp->t_force && len == 1)
285          tcpstat.tcps_sndprobe++;
286      else if (SEQ_LT(tp->snd_next, tp->snd_max)) {
287          tcpstat.tcps_sndrexitpack++;
288          tcpstat.tcps_sndrexitbyte += len;
289      } else {
290          tcpstat.tcps_sndpack++;
291          tcpstat.tcps_sndbyte += len;
292      }
293      MGETHDR(m, M_DONTWAIT, MT_HEADER);
294      if (m == NULL) {
295          error = ENOBUFS;
296          goto out;
297      }
298      m->m_data += max_linkhdr;
299      m->m_len = hdrlen;
300      if (len <= MHLEN - hdrlen - max_linkhdr) {
301          m_copydata(so->so_snd.sb_mb, off, (int) len,
302                  mtod(m, caddr_t) + hdrlen);
303          m->m_len += len;
304      } else {
305          m->m_next = m_copy(so->so_snd.sb_mb, off, (int) len);
306          if (m->m_next == 0)
307              len = 0;
308      }
309      /*
310      * If we're sending everything we've got, set PUSH.
311      * (This will keep happy those implementations that
312      * give data to the user only when a buffer fills or
313      * a PUSH comes in.)
314      */
315      if (off + len == so->so_snd.sb_cc)
316          flags |= TH_PUSH;

```

Figure 26.25 tcp_output function: update statistics, allocate mbuf for IP and TCP headers.

Set PSH flag

309-316

If TCP is sending everything it has from the send buffer, the PSH flag is set. As the comment indicates, this is intended for receiving systems that only pass received data to an application when the PSH flag is received or when a buffer fills. We'll see in `tcp_input` that Net/3 never holds data in a socket receive buffer waiting for a received PSH flag.

The next part of `tcp_output`, shown in Figure 26.26, starts with the code that is executed when `len` equals 0: there is no data in the segment TCP is sending.

```

317 } else { /* len == 0 */
318     if (tp->t_flags & TF_ACKNOW)
319         tcpstat.tcps_sndacks++;
320     else if (flags & (TH_SYN | TH_FIN | TH_RST))
321         tcpstat.tcps_sndctrl++;
322     else if (SEQ_GT(tp->snd_up, tp->snd_una))
323         tcpstat.tcps_sndurg++;
324     else
325         tcpstat.tcps_sndwinup++;
326
327     MGETHDR(m, M_DONTWAIT, MT_HEADER);
328     if (m == NULL) {
329         error = ENOBUFS;
330         goto out;
331     }
332     m->m_data += max_linkhdr;
333     m->m_len = hdrlen;
334
335     m->m_pkthdr.rcvif = (struct ifnet *) 0;
336     ti = mtod(m, struct tcpiphdr *);
337     if (tp->t_template == 0)
338         panic("tcp_output");
339     bcopy((caddr_t) tp->t_template, (caddr_t) ti, sizeof(struct tcpiphdr));

```

Figure 26.26 `tcp_output` function: update statistics and allocate mbuf for IP and TCP headers.

Update statistics

318-325 Various statistics are updated: `TF_ACKNOW` and a length of 0 means this is an ACK-only segment. If any one of the flags SYN, FIN, or RST is set, this is a control segment. If the urgent pointer exceeds `snd_una`, the segment is being sent to notify the other end of the urgent pointer. If none of these conditions are true, this segment is a window update.

Get mbuf for IP and TCP headers

326-335 An mbuf with a packet header is allocated to contain the IP and TCP headers.

Copy IP and TCP header templates into mbuf

336-338 The template of the IP and TCP headers is copied from `t_template` into the mbuf by `bcopy`. This template was created by `tcp_template`.

Figure 26.27 shows the next part of `tcp_output`, which fills in some remaining fields in the TCP header.

Decrement `and_next` if FIN is being retransmitted

339-346 If TCP has already transmitted the FIN, the send sequence space appears as shown in Figure 26.28.

```

339  /*-----tcp_output.c
340  * Fill in fields, remembering maximum advertised
341  * window for use in delaying messages about window sizes.
342  * If resending a FIN, be sure not to use a new sequence number.
343  */
344  if (flags & TH_FIN && tp->t_flags & TF_SENTFIN &&
345      tp->snd_nxt == tp->snd_max)
346      tp->snd_nxt--;
347  /*
348  * If we are doing retransmissions, then snd_nxt will
349  * not reflect the first unsent octet. For ACK only
350  * packets, we do not want the sequence number of the
351  * retransmitted packet, we want the sequence number
352  * of the next unsent octet. So, if there is no data
353  * (and no SYN or FIN), use snd_max instead of snd_nxt
354  * when filling in ti_seq. But if we are in persist
355  * state, snd_max might reflect one byte beyond the
356  * right edge of the window, so use snd_nxt in that
357  * case, since we know we aren't doing a retransmission.
358  * (retransmit and persist are mutually exclusive...)
359  */
360  if ((len || (flags & (TH_SYN | TH_FIN)) || tp->t_timer(TCPT_PERSIST))
361      ti->ti_seq = htonl(tp->snd_nxt);
362  else
363      ti->ti_seq = htonl(tp->snd_max);
364  ti->ti_ack = htonl(tp->rcv_nxt);
365  if (optlen) {
366      bcopy((caddr_t) opt, (caddr_t) (ti + 1), optlen);
367      ti->ti_off = (sizeof(struct tcphdr) + optlen) >> 2;
368  }
369  ti->ti_flags = flags;
-----tcp_output.c

```

Figure 26.27 tcp_output function: set ti_seq, ti_ack, and ti_flags.

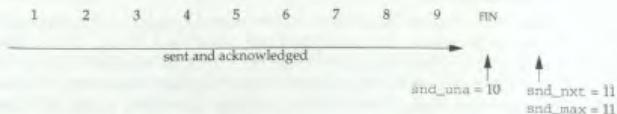


Figure 26.28 Send sequence space after FIN has been transmitted.

Therefore, if the FIN flag is set, and if the TF_SENTFIN flag is set, and if `snd_nxt` equals `snd_max`, TCP knows the FIN is being retransmitted. We'll see shortly (Figure 26.31) that when a FIN is sent, `snd_nxt` is incremented 1 one (since the FIN occupies a sequence number), so this piece of code decrements `snd_nxt` by 1.

Set sequence number field of segment

347-363 The sequence number field of the segment is normally set to `snd_nxt`, but is set to `snd_max` if (1) there is no data to send (`len` equals 0), (2) neither the SYN flag nor the FIN flag is set, and (3) the persist timer is not set.

Set acknowledgment field of segment

364 The acknowledgment field of the segment is always set to `rcv_nxt`, the next expected receive sequence number.

Set header length if options present

365-368 If TCP options are present (`opt_len` is greater than 0), the options are copied into the TCP header and the 4-bit header length in the TCP header (`th_off` in Figure 24.10) is set to the fixed size of the TCP header (20 bytes) plus the length of the options, divided by 4. This field is the number of 32-bit words in the TCP header, including options.

369 The flags field in the TCP header is set from the variable `flags`.

The next part of code, shown in Figure 26.29, fills in more fields in the TCP header and calculates the TCP checksum.

Don't advertise less than one full-sized segment

370-385 Avoidance of the silly window syndrome is performed, this time in calculating the window size that is advertised to the other end (`ti_win`). Recall that `win` was set at the end of Figure 26.3 to the amount of space in the socket's receive buffer. If `win` is less than one-fourth of the receive buffer size (`so_rcv.sb_hiwat`) and less than one full-sized segment, the advertised window will be 0. This is subject to the later test that prevents the window from shrinking. In other words, when the amount of available space reaches either one-fourth of the receive buffer size or one full-sized segment, the available space will be advertised.

Observe upper limit for advertised window on this connection

376-377 If `win` is larger than the maximum value for this connection, reduce it to its maximum value.

Do not shrink window

378-379 Recall from Figure 26.10 that `rcv_adv` minus `rcv_nxt` is the amount of space still available to the sender that was previously advertised. If `win` is less than this value, `win` is set to this value, because we must not shrink the window. This can happen when the available space is less than one full-sized segment (hence `win` was set to 0 at the beginning of this figure), but there is room in the receive buffer for some data. Figure 22.3 of Volume 1 shows an example of this scenario.

Set urgent offset

382-383 If the urgent pointer (`snd_up`) is greater than `snd_nxt`, TCP is in urgent mode. The urgent offset in the TCP header is set to the 16-bit offset of the urgent pointer from the starting sequence number of the segment, and the URG flag bit is set. TCP sends the urgent offset and the URG flag regardless of whether the referenced byte of urgent data is contained in this segment or not.


```

370 /* tcp_output.c
371  * Calculate receive window. Don't shrink window,
372  * but avoid silly window syndrome.
373  */
374 if (win < (long) (so->so_rcv.sb_hiwat / 4) && win < (long) tp->t_maxseg)
375     win = 0;
376 if (win > (long) TCP_MAXWIN << tp->rcv_scale)
377     win = (long) TCP_MAXWIN << tp->rcv_scale;
378 if (win < (long) (tp->rcv_adv - tp->rcv_next))
379     win = (long) (tp->rcv_adv - tp->rcv_next);
380 ti->ti_win = htons((u_short) (win >> tp->rcv_scale));
381 if (SEQ_GT(tp->snd_up, tp->snd_next)) {
382     ti->ti_urg = htons((u_short) (tp->snd_up - tp->snd_next));
383     ti->ti_flags |= TH_URG;
384 } else
385     /*
386     * If no urgent pointer to send, then we pull
387     * the urgent pointer to the left edge of the send window
388     * so that it doesn't drift into the send window on sequence
389     * number wraparound.
390     */
391     tp->snd_up = tp->snd_una; /* drag it along */
392 /*
393  * Put TCP length in extended header, and then
394  * checksum extended header and data.
395  */
396 if (len + optlen)
397     ti->ti_len = htons((u_short) (sizeof(struct tcphdr) +
398                                optlen + len));
399 ti->ti_sum = in_cksum(m, (int) (hdrlen + len));
tcp_output.c

```

Figure 26.29 tcp_output function: fill in more TCP header fields and calculate checksum.

Figure 26.30 shows an example of how the urgent offset is calculated, assuming the process executes

```
send(fd, buf, 3, MSG_OOB);
```

and the send buffer is empty when this call to `send` takes place. This shows that Berkeley-derived systems consider the urgent pointer to point to the first byte of data *after* the out-of-band byte. Recall our discussion after Figure 24.10 where we distinguished between the 32-bit *urgent pointer* in the data stream (`snd_up`), and the 16-bit *urgent offset* in the TCP header (`ti_urg`).

There is a subtle bug here. The bug occurs when the send buffer is larger than 65535, regardless of whether the window scale option is in use or not. If the send buffer is greater than 65535 and is nearly full, and the process sends out-of-band data, the offset of the urgent pointer from `snd_next` can exceed 65535. But the urgent pointer is a 16-bit unsigned value, and if the calculated value exceeds 65535, the 16 high-order bits are discarded, delivering a bogus urgent pointer to the other end. See Exercise 26.6 for a solution.

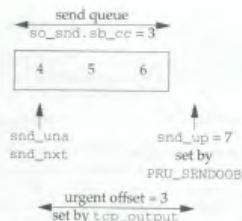


Figure 26.30 Example of urgent pointer and urgent offset calculation.

384-391 If TCP is not in urgent mode, the urgent pointer is moved to the left edge of the window (`snd_una`).

392-399 The TCP length is stored in the pseudo-header and the TCP checksum is calculated. All the fields in the TCP header have been filled in, and when the IP and TCP header template were copied from `t_template` (Figure 26.26), the fields in the IP header that are used as the pseudo-header were initialized (as shown in Figure 23.19 for the UDP checksum calculation).

The next part of `tcp_output`, shown in Figure 26.31, updates the sequence number if the SYN or FIN flags are set and initializes the retransmission timer.

Remember starting sequence number

400-405 If TCP is not in the persist state, the starting sequence number is saved in `startseq`. This is used later in Figure 26.31 if the segment is timed.

Increment `snd_nxt`

406-417 Since both the SYN and FIN flags take a sequence number, `snd_nxt` is incremented if either is set. TCP also remembers that the FIN has been sent, by setting the `TF_SENTFIN`. `snd_nxt` is then incremented by the number of bytes of data (`len`), which can be 0.

Update `snd_max`

418-419 If the new value of `snd_nxt` is larger than `snd_max`, this is not a retransmission. The new value of `snd_max` is stored.

420-428 If a segment is not currently being timed for this connection (`t_rtt` equals 0), the timer is started (`t_rtt` is set to 1) and the starting sequence number of the segment being timed is saved in `t_rtseq`. This sequence number is used by `tcp_input` to determine when the segment being timed is acknowledged, to update the RTT estimators. The sample code we discussed in Section 25.10 looked like

```
if (tp->t_rtt && SEQ_GT(ti->ti_ack, tp->t_rtseq))
    tcp_xmit_timer(tp, tp->t_rtt);
```

```

400  /* -tcp_output.c
401  * In transmit state, time the transmission and arrange for
402  * the retransmit. In persist state, just set snd_max.
403  */
404  if (tp->t_force == 0 || tp->t_timer[TCPT_PERSIST] == 0) {
405      tcp_seq startseq = tp->snd_nxt;

406      /*
407      * Advance snd_nxt over sequence space of this segment.
408      */
409      if (flags & (TH_SYN | TH_FIN)) {
410          if (flags & TH_SYN)
411              tp->snd_nxt++;
412          if (flags & TH_FIN) {
413              tp->snd_nxt++;
414              tp->t_flags |= TF_SENTFIN;
415          }
416      }
417      tp->snd_nxt += len;
418      if (SEQ_GT(tp->snd_nxt, tp->snd_max)) {
419          tp->snd_max = tp->snd_nxt;
420          /*
421          * Time this transmission if not a retransmission and
422          * not currently timing anything.
423          */
424          if (tp->t_rtt == 0) {
425              tp->t_rtt = 1;
426              tp->t_rtseq = startseq;
427              tcpstat.tcps_segstimed++;
428          }
429      }
430      /*
431      * Set retransmit timer if not currently set,
432      * and not doing an ack or a keepalive probe.
433      * Initial value for retransmit timer is smoothed
434      * round-trip time + 2 * round-trip time variance.
435      * Initialize counter which is used for backoff
436      * of retransmit time.
437      */
438      if (tp->t_timer[TCPT_REXMT] == 0 &&
439          tp->snd_nxt != tp->snd_una) {
440          tp->t_timer[TCPT_REXMT] = tp->t_rxtcur;
441          if (tp->t_timer[TCPT_PERSIST]) {
442              tp->t_timer[TCPT_PERSIST] = 0;
443              tp->t_rxtshift = 0;
444          }
445      }
446      } else if (SEQ_GT(tp->snd_nxt + len, tp->snd_max))
447      tp->snd_max = tp->snd_nxt + len;
-tcp_output.c

```

Figure 26.31 tcp_output function: fill in remaining fields in TCP header and calculate checksum.

Set retransmission timer

430-440 If the retransmission timer is not currently set, and if this segment contains data, the retransmission timer is set to `t_rxtcur`. Recall that `t_rxtcur` is set by `tcp_xmit_timer`, when an RTT measurement is made. This is an ACK-only segment if `snd_nxt` equals `snd_una` (since `len` was added to `snd_nxt` earlier in this figure), and the retransmission timer is set only for segments containing data.

441-444 If the persist timer is enabled, it is disabled. Either the retransmission timer or the persist timer can be enabled at any time for a given connection, but not both.

Persist state

446-447 The connection is in the persist state since `t_force` is nonzero and the persist timer is enabled. (This `else` clause is associated with the `if` at the beginning of the figure.) `snd_max` is updated, if necessary. In the persist state, `len` will be one.

The final part of `tcp_output`, shown in Figure 26.32 completes the formation of the outgoing segment and calls `ip_output` to send the datagram.

Add trace record for socket debugging

448-452 If the `SO_DEBUG` socket option is enabled, `tcp_trace` adds a record to TCP's circular trace buffer. We describe this function in Section 27.10.

Set IP length, TTL, and TOS

463-462 The final three fields in the IP header that must be set by the transport layer are stored: IP length, TTL, and TOS. These three fields are marked with an asterisk at the bottom of Figure 23.19.

The comments `xxx are` because the latter two fields normally remain constant for a connection and should be stored in the header template, instead of being assigned explicitly each time a segment is sent. But these two fields cannot be stored in the IP header until after the TCP checksum is calculated.

Pass datagram to IP

463-464 `ip_output` sends the datagram containing the TCP segment. The socket options are logically ANDed with `SO_DONTROUTE`, which means that the only socket option passed to `ip_output` is `SO_DONTROUTE`. The only other socket option examined by `ip_output` is `SO_BROADCAST`, so this logical AND turns off the `SO_BROADCAST` bit, if set. This means that a process cannot issue a `connect` to a broadcast address, even if it sets the `SO_BROADCAST` socket option.

467-470 The error `ENOBUFS` is returned if the interface queue is full or if IP needs to obtain an mbuf and can't. The function `tcp_quench` puts the connection into slow start, by setting the congestion window to one full-sized segment. Notice that `tcp_output` still returns 0 (OK) in this case, instead of the error, even though the datagram was discarded. This differs from `udp_output` (Figure 23.20), which returned the error. The difference is that UDP is unreliable, so the `ENOBUFS` error return is the only indication to the process that the datagram was discarded. TCP, however, will time out (if the segment contains data) and retransmit the datagram, and it is hoped that there will be space on the interface output queue or more available mbufs. If the TCP segment


```

448 /* tcp_output.c
449  * Trace.
450  */
451 if (so->so_options & SO_DEBUG)
452     tcp_trace(TA_OUTPUT, tp->t_state, tp, ti, 0);
453
454 /*
455  * Fill in IP length and desired time to live and
456  * send to IP level. There should be a better way
457  * to handle ttl and tos; we could keep them in
458  * the template, but need a way to checksum without them.
459  */
460 m->m_pkthdr.len = hdrlen + len;
461 ((struct ip *) ti)->ip_len = m->m_pkthdr.len;
462 ((struct ip *) ti)->ip_ttl = tp->t_inpcb->inp_ip.ip_ttl; /* XXX */
463 ((struct ip *) ti)->ip_tos = tp->t_inpcb->inp_ip.ip_tos; /* XXX */
464 error = ip_output(m, tp->t_inpcb->inp_options, &tp->t_inpcb->inp_route,
465                 so->so_options & SO_DONTROUTE, 0);
466 if (error) {
467     out:
468     if (error == ENOBUFS) {
469         tcp_quench(tp->t_inpcb, 0);
470         return (0);
471     }
472     if ((error == EHOSTUNREACH || error == ENETDOWN)
473         && TCPS_HAVERCVDSYN(tp->t_state)) {
474         tp->t_softerror = error;
475         return (0);
476     }
477     return (error);
478 }
479 tcpstat.tcps_sndtotal++;
480
481 /*
482  * Data sent (as far as we can tell).
483  * If this advertises a larger window than any other segment,
484  * then remember the size of the advertised window.
485  * Any pending ACK has now been sent.
486  */
487 if (win > 0 && SEQ_GT(tp->rcv_nxt + win, tp->rcv_adv))
488     tp->rcv_adv = tp->rcv_nxt + win;
489 tp->last_ack_sent = tp->rcv_nxt;
490 tp->t_flags &= ~(TF_ACKNOW | TF_DELACK);
491
492 if (sendalot)
493     goto again;
494 return (0);
495 }
tcp_output.c

```

Figure 26.32 tcp_output function: call ip_output to send segment.

doesn't contain data, the other end will time out when the ACK isn't received and will retransmit the data whose ACK was discarded.

471-475 If a route can't be located for the destination, and if the connection has received a SYN, the error is recorded as a soft error for the connection.

When `tcp_output` is called by `tcp_usrreq` as part of a system call by a process (Chapter 30, the `PRU_CONNECT`, `PRU_SEND`, `PRU_SENDOOB`, and `PRU_SHUTDOWN` requests), the process receives the return value from `tcp_output`. Other functions that call `tcp_output`, such as `tcp_input` and the fast and slow timeout functions, ignore the return value (because these functions don't return an error to a process).

Update `rcv_adv` and `last_ack_sent`

479-486 If the highest sequence number advertised in this segment (`rcv_nxt` plus `win`) is larger than `rcv_adv`, the new value is saved. Recall that `rcv_adv` was used in Figure 26.9 to determine how much the window had opened since the last segment that was sent, and in Figure 26.29 to make certain TCP was not shrinking the window.

487 The value of the acknowledgment field in the segment is saved in `last_ack_sent`. This variable is used by `tcp_input` with the timestamp option (Section 26.6).

488 Any pending ACK has been sent, so the `TF_ACKNOW` and `TF_DELACK` flags are cleared.

More data to send?

489-490 If the `sendatol` flag is set, a jump is made back to the label again (Figure 26.1). This occurs if the send buffer contains more than one full-sized segment that can be sent (Figure 26.3), or if a full-sized segment was being sent and TCP options were included that reduced the amount of data in the segment (Figure 26.24).

26.8 `tcp_template` Function

The function `tcp_newtcpcb` (from the previous chapter) is called when the socket is created, to allocate and partially initialize the TCP control block. When the first segment is sent or received on the socket (an active open is performed, the `PRU_CONNECT` request, or a SYN arrives for a listening socket), `tcp_template` creates a template of the IP and TCP headers for the connection. This minimizes the amount of work required by `tcp_output` when a segment is sent on the connection.

Figure 26.33 shows the `tcp_template` function.

Allocate `mbuf`

49-72 The template of the IP and TCP headers is formed in an `mbuf`, and a pointer to the `mbuf` is stored in the `t_template` member of the TCP control block. Since this function can be called at the software interrupt level, from `tcp_input`, the `M_DONTWAIT` flag is specified.

Initialize header fields

73-88 All the fields in the IP and TCP headers are set to 0 except as follows: `ti_pr` is set to the IP protocol value for TCP (6); `ti_len` is set to 20, the default length of the TCP

```

59 struct tcpiphdr *
60 tcp_template(tp)
61 struct tcpcb *tp;
62 {
63     struct inpcb *inp = tp->t_inpcb;
64     struct mbuf *m;
65     struct tcpiphdr *n;
66
67     if ((n = tp->t_template) == 0) {
68         m = m_get(M_DONTWAIT, MT_HEADER);
69         if (m == NULL)
70             return (0);
71         m->m_len = sizeof(struct tcpiphdr);
72         n = mtod(m, struct tcpiphdr *);
73     }
74     n->ti_next = n->ti_prev = 0;
75     n->ti_x1 = 0;
76     n->ti_pr = IPPROTO_TCP;
77     n->ti_len = htons(sizeof(struct tcpiphdr) - sizeof(struct ip));
78     n->ti_src = inp->inp_laddr;
79     n->ti_dst = inp->inp_faddr;
80     n->ti_sport = inp->inp_lport;
81     n->ti_dport = inp->inp_fport;
82     n->ti_seq = 0;
83     n->ti_ack = 0;
84     n->ti_x2 = 0;
85     n->ti_off = 5; /* 5 32-bit words = 20 bytes */
86     n->ti_flags = 0;
87     n->ti_wln = 0;
88     n->ti_sum = 0;
89     n->ti_urp = 0;
90     return (n);

```

Figure 26.33 tcp_template function: create template of IP and TCP headers.

header; and `ti_off` is set to 5, the number of 32-bit words in the 20-byte TCP header. Also the source and destination IP addresses and TCP port numbers are copied from the Internet PCB into the TCP header template.

Pseudo-header for TCP checksum computation

⁷³⁻⁸⁸ The initialization of many of the fields in the combined IP and TCP header simplifies the computation of the TCP checksum, using the same pseudo-header technique as discussed for UDP in Section 23.6. Examining the `udpiphdr` structure in Figure 23.19 shows why `tcp_template` initializes fields such as `ti_next` and `ti_prev` to 0.

26.9 tcp_respond Function

The function `tcp_respond` is a special-purpose function that also calls `ip_output` to send IP datagrams. `tcp_respond` is called in two cases:

1. by `tcp_input` to generate an RST segment, with or without an ACK, and
2. by `tcp_timers` to send a keepalive probe.

Instead of going through all the logic of `tcp_output` for these two cases, the special-purpose function `tcp_respond` is called. We also note that the function `tcp_drop` that we cover in the next chapter also generates RST segments by calling `tcp_output`. Not all RST segments are generated by `tcp_respond`.

Figure 26.34 shows the first half of `tcp_respond`.

```

104 void                                     tcp_subr.c
105 tcp_respond(tp, ti, m, ack, seq, flags)
106 struct tcpcb *tp;
107 struct tcphdr *ti;
108 struct mbuf *m;
109 tcp_seq ack, seq;
110 int flags;
111 {
112     int tlen;
113     int win = 0;
114     struct route *ro = 0;
115     if (tp) {
116         win = sbspace(&tp->t_inpcb->inp_socket->so_rcv);
117         ro = &tp->t_inpcb->inp_route;
118     }
119     if (m == 0) { /* generate keepalive probe */
120         m = m_gethdr(M_DONTWAIT, MT_HEADER);
121         if (m == NULL)
122             return;
123         tlen = 0; /* no data is sent */
124         m->m_data += MAX_linkhdr;
125         *mtod(m, struct tcphdr *) = *ti;
126         ti = mtod(m, struct tcphdr *);
127         flags = TH_ACK;
128     } else { /* generate RST segment */
129         m_freem(m->m_next);
130         m->m_next = 0;
131         m->m_data = (caddr_t) ti;
132         m->m_len = sizeof(struct tcphdr);
133         tlen = 0;
134 #define xchg(a,b,type) { type t; t=a; a=b; b=t; }
135 xchg(ti->ti_dst.s_addr, ti->ti_src.s_addr, u_long);
136 xchg(ti->ti_dport, ti->ti_sport, u_short);
137 #undef xchg
138     }

```

Figure 26.34 `tcp_respond` function: first half.

104-110 Figure 26.35 shows the different arguments to `tcp_respond` for the three cases in which it is called.

	Arguments					
generate RST without ACK	tp	ti	m	ack	seq	flags
generate RST with ACK	tp	ti	m	0	ti_ack	TH_RST TH_ACK
generate keepalive	tp	t_template	NULL	rcv_nxt	snd_una	0

Figure 26.35 Arguments to `tcp_respond`.

`tp` is a pointer to the TCP control block (possibly a null pointer); `ti` is a pointer to an IP/TCP header template; `m` is a pointer to the mbuf containing the segment causing the RST to be generated; and the last three arguments are the acknowledgment field, sequence number field, and flags field of the segment being generated.

113-116 It is possible for `tcp_input` to generate an RST when a segment is received that does not have an associated TCP control block. This happens, for example, when a segment is received that doesn't reference an existing connection (e.g., a SYN for a port without an associated listening server). In this case `tp` is null and the initial values for `win` and `ro` are used. If `tp` is not null, the amount of space in the receive buffer will be sent as the advertised window, and the pointer to the cached route is saved in `ro` for the call to `ip_output`.

Send keepalive probe when keepalive timer expires

113-127 The argument `m` is a pointer to the mbuf chain for the received segment. But a keepalive probe is sent in response to the keepalive timer expiring, not in response to a received TCP segment. Therefore `m` is null and `m_gethdr` allocates a packet header mbuf to contain the IP and TCP headers. `ti_len`, the length of the TCP data, is set to 0, since the keepalive probe doesn't contain any data.

Some older implementations based on 4.2BSD do not respond to these keepalive probes unless the segment contains data. Net/3 can be configured to send 1 garbage byte of data in the probe to elicit the response by defining the name `TCP_COMPAT_42` when the kernel is compiled. This assigns 1, instead of 0, to `ti_len`. The garbage byte causes no harm, because it is not the expected byte (it is a byte that the receiver has previously received and acknowledged), so it is thrown away by the receiver.

The assignment of `*ti` copies the TCP header template structure pointed to by `ti` into the data portion of the mbuf. The pointer `ti` is then set to point to the header template in the mbuf.

Send RST segment in response to received segment

128-138 An RST segment is being sent by `tcp_input` in response to a received segment. The mbuf containing the input segment is reused for the response. All the mbufs on the chain are released by `m_free` except the first mbuf (the packet header), since the segment generated by `tcp_respond` consists of only an IP header and a TCP header. The source and destination IP address and port numbers are swapped in the IP and TCP header.

Figure 26.36 shows the final half of `tcp_respond`.

```

239  ti->ti_len = htons((u_short) (sizeof(struct tcphdr) + tlen)); tcp_subrc
240  tlen += sizeof(struct tcpiphdr);
241  m->m_len = tlen;
242  m->m_pkthdr.len = tlen;
243  m->m_pkthdr.rcvif = (struct ifnet *) 0;
244  ti->ti_next = ti->ti_prev = 0;
245  ti->ti_x1 = 0;
246  ti->ti_seq = htonl(seq);
247  ti->ti_ack = htonl(ack);
248  ti->ti_x2 = 0;
249  ti->ti_off = sizeof(struct tcphdr) >> 2;
250  ti->ti_flags = flags;
251  if (tp)
252      ti->ti_win = htons((u_short) (win >> tp->rcv_scale));
253  else
254      ti->ti_win = htons((u_short) win);
255  ti->ti_urp = 0;
256  ti->ti_sum = 0;
257  ti->ti_sum = in_cksum(m, tlen);
258  ((struct ip *) ti->ip_len = tlen;
259  ((struct ip *) ti->ip_ttl = ip_defttl;
260  (void) ip_output(m, NULL, ro, 0, NULL);
261 }
tcp_subrc

```

Figure 26.36 `tcp_respond` function: second half.

139-157 The fields in the IP and TCP headers must be initialized for the TCP checksum computation. These statements are similar to the way `tcp_template` initializes the `t_template` field. The sequence number and acknowledgment fields are passed by the caller as arguments. Finally `ip_output` sends the datagram.

26.10 Summary

This chapter has looked at the general-purpose function that generates most TCP segments (`tcp_output`) and the special-purpose function that generates RST segments and keepalive probes (`tcp_respond`).

Many factors determine whether TCP can send a segment or not: the flags in the segment, the window advertised by the other end, the amount of data ready to send, whether unacknowledged data already exists for the connection, and so on. Therefore the logic of `tcp_output` determines whether a segment can be sent (the first half of the function), and if so, what values to set all the TCP header fields to (the last half of the function). If a segment is sent, the TCP control block variables for the send sequence space must be updated.

One segment at a time is generated by `tcp_output`, and at the end of the function a check is made of whether more data can still be sent. If so, the function loops around and tries to send another segment. This looping continues until there is no more data to

send, or until some other condition (e.g., the receiver's advertised window) stops the transmission.

A TCP segment can also contain options. The options supported by Net/3 specify the maximum segment size, a window scale factor, and a pair of timestamps. The first two can only appear with SYN segments, while the timestamp option (if supported by both ends) normally appears in every segment. Since the window scale and timestamp options are newer and optional, if the first end to send a SYN wants to use the option, it sends the option with its SYN and uses the option only if the other end's SYN also contains the option.

Exercises

- 26.1 Slow start is resumed in Figure 26.1 when there is a pause in the *sending* of data, yet the amount of idle time is calculated as the amount of time since the last segment was *received* on the connection. Why doesn't TCP calculate the idle time as the amount of time since the last segment was *sent* on the connection?
- 26.2 With Figure 26.6 we said that `len` is less than 0 if the FIN has been sent but not acknowledged and not retransmitted. What happens if the FIN is retransmitted?
- 26.3 Net/3 always sends the window scale and timestamp options with an active open. Why does the global variable `tcp_do_rfc1323` exist?
- 26.4 In Figure 25.28, which did not use the timestamp option, the RTT estimators are updated eight times. If the timestamp option had been used in this example, how many times would the RTT estimators have been updated?
- 26.5 In Figure 26.23 `bcopy` is called to store the received MSS in the variable `mss`. Why not cast the pointer to `opt[2]` into a pointer to an unsigned short and perform an assignment?
- 26.6 After Figure 26.29 we described a bug in the code, which can cause a bogus urgent offset to be sent. Propose a solution. (*Hint*: What is the largest amount of TCP data that can be sent in a segment?)
- 26.7 With Figure 26.32 we mentioned that an error of `ENOBUFS` is not returned to the process because (1) if the discarded segment contained data, the retransmission timer will expire and the data will be retransmitted, or (2) if the discarded segment was an ACK-only segment, the other end will retransmit its data when it doesn't receive the ACK. What if the discarded segment contains an RST?
- 26.8 Explain the settings of the PSH flag in Figure 20.3 of Volume 1.
- 26.9 Why does Figure 26.36 use the value of `ip_defrttl` for the TTL, while Figure 26.32 uses the value in the PCB?
- 26.10 Describe what happens with the mbuf allocated in Figure 26.25 when IP options are specified by the process for the TCP connection. Implement a better solution.
- 26.11 `tcp_output` is a long function (about 500 lines, including comments), which can appear to be inefficient. But lots of the code handles special cases. Assume the function is called with a full-sized segment ready to be sent, and no special cases: no IP options and no special flags such as SYN, FIN, or URG. About how many lines of C code are actually executed? How many functions are called before the segment is passed to `ip_output`?

- 26.12 In the example at the end of Section 26.3 in which the application did a write of 100 bytes followed by a write of 50 bytes, would anything change if the application called `writen` once for both buffers, instead of calling `write` twice? Does anything change with `writen` if the two buffer lengths are 200 and 300, instead of 100 and 50?
- 26.13 The timestamp that is sent in the timestamp option is taken from the global `tcp_now`, which is incremented every 500 ms. Modify TCP to use a higher resolution timestamp value.

TCP Functions

27.1 Introduction

This chapter presents numerous TCP functions that we need to cover before discussing TCP input in the next two chapters:

- `tcp_drain` is the protocol's drain function, called when the kernel is out of mbufs. It does nothing.
- `tcp_drop` aborts a connection by sending an RST.
- `tcp_close` performs the normal TCP connection termination: send a FIN and wait for the four-way exchange to complete. Section 18.2 of Volume 1 talks about the four packets that are exchanged when a connection is closed.
- `tcp_mss` processes a received MSS option and calculates the MSS to announce when TCP sends an MSS option of its own.
- `tcp_ctlinput` is called when an ICMP error is received in response to a TCP segment, and it calls `tcp_notify` to process the ICMP error. `tcp_quench` is a special case function that handles ICMP source quench errors.
- The `TCP_REASS` macro and the `tcp_reass` function manipulate segments on TCP's reassembly queue for a given connection. This queue handles the receipt of out-of-order segments, some of which might overlap.
- `tcp_trace` adds records to the kernel's circular debug buffer for TCP (the `SO_DEBUG` socket option) that can be printed with the `trpt(8)` program.

27.2 tcp_drain Function

The simplest of all the TCP functions is `tcp_drain`. It is the protocol's `pr_drain` function, called by `m_reclaim` when the kernel runs out of mbufs. We saw in Figure 10.32 that `ip_drain` discards all the fragments on its reassembly queue, and UDP doesn't define a drain function. Although TCP holds onto mbufs—segments that have arrived out of order, but within the receive window for the socket—the Net/3 implementation of TCP does not discard these pending mbufs if the kernel runs out of space. Instead, `tcp_drain` does nothing, on the assumption that a received (but out-of-order) TCP segment is “more important” than an IP fragment.

27.3 tcp_drop Function

`tcp_drop` is called from numerous places to drop a connection by sending an RST and to report an error to the process. This differs from closing a connection (the `tcp_disconnect` function), which sends a FIN to the other end and follows the connection termination steps in the state transition diagram.

Figure 27.1 shows the seven places where `tcp_drop` is called and the `errno` argument.

Function	errno	Description
<code>tcp_input</code>	ENOBUFS	SYN arrives on listening socket, but kernel out of mbufs for <code>t_template</code> .
<code>tcp_input</code>	ECONNREFUSED	RST received in response to SYN.
<code>tcp_input</code>	ECONNRESET	RST received on existing connection.
<code>tcp_timers</code>	ETIMEDOUT	Retransmission timer has expired 13 times in a row with no ACK from other end (Figure 25.25).
<code>tcp_timers</code>	ETIMEDOUT	Connection-establishment timer has expired (Figure 25.16), or keepalive timer has expired with no response to nine consecutive probes (Figure 25.17).
<code>tcp_usrreq</code>	ECONNABORTED	PRU_ABORT request.
<code>tcp_usrreq</code>	0	Socket closed and <code>SO_LINGER</code> socket option set with linger time of 0.

Figure 27.1 Calls to `tcp_drop` and `errno` argument.

Figure 27.2 shows the `tcp_drop` function.

302-213 If TCP has received a SYN, the connection is synchronized and an RST must be sent to the other end. This is done by setting the state to CLOSED and calling `tcp_output`. In Figure 24.16 the value of `tcp_outflags` for the CLOSED state includes the RST flag.

214-216 If the error is ETIMEDOUT but a soft error was received on the connection (e.g., EHOSTUNREACH), the soft error becomes the socket error, instead of the less specific ETIMEDOUT.

217 `tcp_close` finishes closing the socket.

```

202 struct tcpcb *                                tcp_subr.c
203 tcp_drop(tp, errno)
204 struct tcpcb *tp;
205 int      errno;
206 {
207     struct socket *so = tp->t_inpcb->inp_socket;
208     if (TCPS_HAVERCVDSYN(tp->t_state)) {
209         tp->t_state = TCPS_CLOSED;
210         (void) tcp_output(tp);
211         tcpstat.tcps_drops++;
212     } else
213         tcpstat.tcps_connndrops++;
214     if (errno == ETIMEDOUT && tp->t_softerror)
215         errno = tp->t_softerror;
216     so->so_error = errno;
217     return (tcp_close(tp));
218 }

```

Figure 27.2 tcp_drop function.

27.4 tcp_close Function

`tcp_close` is normally called by `tcp_input` when the process has done a passive close and the ACK is received in the `LAST_ACK` state, and by `tcp_timers` when the 2MSL timer expires and the socket moves from the `TIME_WAIT` to `CLOSED` state. It is also called in other states, possibly after an error has occurred, as we saw in the previous section. It releases the memory occupied by the connection (the IP and TCP header template, the TCP control block, the Internet PCB, and any out-of-order segments remaining on the connection's reassembly queue) and updates the route characteristics.

We describe this function in three parts, the first two dealing with the route characteristics and the final part showing the release of resources.

Route Characteristics

Nine variables are maintained in the `rt_metrics` structure (Figure 18.26), six of which are used by TCP. Eight of these can be examined and changed with the `route(8)` command (the ninth, `rmx_pkssent` is never used): these variables are shown in Figure 27.3.

Additionally, the `-lock` modifier can be used with the `route` command to set the corresponding `RTV_xxx` bit in the `rmx_locks` member (Figure 20.13). Setting the `RTV_xxx` bit tells the kernel not to update that metric.

When a TCP socket is closed, `tcp_close` updates three of the routing metrics—the smoothed RTT estimator, the smoothed mean deviation estimator, and the slow start threshold—but only if enough data was transferred on the connection to yield meaningful statistics and the variable is not locked.

Figure 27.4 shows the first part of `tcp_close`.

rt_metrics member	saved by tcp_close?	used by tcp_mss?	route(8) modifier
rmx_expire			-expire
rmx_hopcount			-hopcount
rmx_mtu		•	-mtu
rmx_recvpipe		•	-recvpipe
rmx_rtt	•	•	-rtt
rmx_rttvar	•	•	-rttvar
rmx_sendpipe		•	-sendpipe
rmx_ssthresh	•	•	-ssthresh

Figure 27.3 Members of the `rt_metrics` structure used by TCP.

Check if enough data sent to update statistics

234-248 The default send buffer size is 8192 bytes (`sb_hiwat`), so the first test is whether 131,072 bytes (16 full buffers) have been transferred across the connection. The initial send sequence number is compared to the maximum sequence number sent on the connection. Additionally, the socket must have a cached route and that route cannot be the default route. (See Exercise 19.2.)

Notice there is a small chance for an error in the first test, because of sequence number wrap, if the amount of data transferred is within $N \times 2^{32}$ and $N \times 2^{32} + 131072$, for any N greater than 1. But few connections (today) transfer 4 gigabytes of data.

Despite the prevalence of default routes in the Internet, this information is still useful to maintain in the routing table. If a host continually exchanges data with another host (or network), even if a default route can be used, a host specific or network-specific route can be entered into the routing table with the `route` command just to maintain this information across connections. (See Exercise 19.2.) This information is lost when the system is rebooted.

250 The administrator can lock any of the variables from Figure 27.3, preventing them from being updated by the kernel, so before modifying each variable this lock must be checked.

Update RTT

251-264 `t_srtt` is stored as ticks $\times 8$ (Figure 25.19) and `rmx_rtt` is stored as microseconds. So `t_srtt` is multiplied by 1,000,000 (`RTM_RTTUNIT`) and then divided by 2 (ticks/second) times 8. If a value for `rmx_rtt` already exists, the new value is one-half the old value plus one-half the new value. Otherwise the new value is stored in `rmx_rtt`.

Update mean deviation

265-273 The same algorithm is applied to the mean deviation estimator. It too is stored as microseconds, requiring a conversion from the `t_rttvar` units of ticks $\times 4$.


```

225 struct tcpcb *                                     tcp_subrc
226 tcp_close(tp)
227 struct tcpcb *tp;
228 {
229     struct tepiphdr *t;
230     struct inpcb *inp = tp->t_inpcb;
231     struct socket *so = inp->inp_socket;
232     struct mbuf *m;
233     struct rtentry *rt;
234     /*
235     * If we sent enough data to get some meaningful characteristics,
236     * save them in the routing entry. 'Enough' is arbitrarily
237     * defined as the sendpipesize (default 8K) * 16. This would
238     * give us 16 rtt samples assuming we only get one sample per
239     * window (the usual case on a long haul net). 16 samples is
240     * enough for the rtt filter to converge to within 5% of the correct
241     * value; fewer samples and we could save a very bogus rtt.
242     *
243     * Don't update the default route's characteristics and don't
244     * update anything that the user 'locked'.
245     */
246     if (SBQ_LT(tp->iss + so->so_snd.sb_hiwat * 16, tp->snd_max) &&
247         (rt = inp->inp_route.ro_rt) &&
248         ((struct sockaddr_in *) rt_key(rt))->sin_addr.s_addr != INADDR_ANY) {
249         u_long i;
250         if ((rt->rt_rmx.rmx_locks & RTV_RTT) == 0) {
251             i = tp->t_srtt *
252                 (RTM_RTTUNIT / (PR_SLOWHZ * TCP_RTT_SCALE));
253             if (rt->rt_rmx.rmx_rtt && i)
254                 /*
255                  * filter this update to half the old & half
256                  * the new values, converting scale.
257                  * See route.h and tcp_var.h for a
258                  * description of the scaling constants.
259                  */
260                 rt->rt_rmx.rmx_rtt =
261                     (rt->rt_rmx.rmx_rtt + i) / 2;
262             else
263                 rt->rt_rmx.rmx_rtt = i;
264         }
265         if ((rt->rt_rmx.rmx_locks & RTV_RTTVAR) == 0) {
266             i = tp->t_rttvar *
267                 (RTM_RTTUNIT / (PR_SLOWHZ * TCP_RTTVAR_SCALE));
268             if (rt->rt_rmx.rmx_rttvar && i)
269                 rt->rt_rmx.rmx_rttvar =
270                     (rt->rt_rmx.rmx_rttvar + i) / 2;
271             else
272                 rt->rt_rmx.rmx_rttvar = i;
273         }

```

tcp_subrc

Figure 27.4 tcp_close function: update RTT and mean deviation.

Figure 27.5 shows the next part of `tcp_close`, which updates the slow start threshold for the route.

```

274      /*
275      * update the pipelimit (ssthresh) if it has been updated
276      * already or if a pipesize was specified & the threshold
277      * got below half the pipesize. I.e., wait for bad news
278      * before we start updating, then update on both good
279      * and bad news.
280      */
281      if ((rt->rt_rmx.rmx_locks & RTV_SSTHRESH) == 0 &&
282          (i = tp->snd_ssthresh) && rt->rt_rmx.rmx_ssthresh ||
283          i < (rt->rt_rmx.rmx_sendpipe / 2)) {
284          /*
285          * convert the limit from user data bytes to
286          * packets then to packet data bytes.
287          */
288          i = (i + tp->t_maxseg / 2) / tp->t_maxseg;
289          if (i < 2)
290              i = 2;
291          i *= (u_long) (tp->t_maxseg + sizeof(struct tciphdr));
292          if (rt->rt_rmx.rmx_ssthresh)
293              rt->rt_rmx.rmx_ssthresh =
294                  (rt->rt_rmx.rmx_ssthresh + i) / 2;
295          else
296              rt->rt_rmx.rmx_ssthresh = i;
297      }
298  }

```

tcp_subr.c

Figure 27.5 `tcp_close` function: update slow start threshold.

274-283 The slow start threshold is updated only if (1) it has been updated already (`rmx_ssthresh` is nonzero) or (2) `rmx_sendpipe` is specified by the administrator and the new value of `snd_ssthresh` is less than one-half the value of `rmx_sendpipe`. As the comment in the code indicates, TCP does not update the value of `rmx_ssthresh` until it is forced to because of packet loss; from that point on it considers itself free to adjust the value either up or down.

284-290 The variable `snd_ssthresh` is maintained in bytes. The first conversion divides this variable by the MSS (`t_maxseg`), yielding the number of segments. The addition of one-half `t_maxseg` rounds the integer result. The lower bound on this result is two segments.

291-297 The size of the IP and TCP headers (40) is added to the MSS and multiplied by the number of segments. This value updates `rmx_ssthresh`, using the same filtering as in Figure 27.4 (one-half the old plus one-half the new).

Resource Release

The final part of `tcp_close`, shown in Figure 27.6, releases the memory resources held by the socket.

```

299      /* free the reassembly queue, if any */
300      t = tp->seg_next;
301      while (t != (struct tcpiphdr *) tp) {
302          t = (struct tcpiphdr *) t->ti_next;
303          m = REASS_MDUP((struct tcpiphdr *) t->ti_prev);
304          remque(t->ti_prev);
305          m_freem(m);
306      }
307      if (tp->t_template)
308          (void) m_free(dtom(tp->t_template));
309      free(tp, M_PCB);
310      inp->inp_ppcb = 0;
311      sodisconnected(so);
312      /* clobber input pcb cache if we're closing the cached connection */
313      if (inp == tcp_last_inpcb)
314          tcp_last_inpcb = stcb;
315      in_pcbdetach(inp);
316      tcpstat.tcps_closed++;
317      return ((struct tcpcb *) 0);
318 }

```

Figure 27.6 tcp_close function: release connection resources.

Release any mbufs on reassembly queue

299-306 If any segments are left on the connection's reassembly queue, they are discarded. This queue is for segments that arrive out of order but within the receive window. They are held in a reassembly queue until the required "earlier" segments are received, at which time they are reassembled and passed to the application in the correct order. We discuss this in more detail in Section 27.9.

Release header template and TCP control block

307-309 The template of the IP and TCP headers is released by `m_free` and the TCP control block is released by `free`. The socket is marked as disconnected by `sodisconnect`, which issues the `PRU_DISCONNECT` request.

Release PCB

310-318 If the Internet PCB for this socket is the one currently cached by TCP, the cache is marked as empty by setting `tcp_last_inpcb` to the head of TCP's PCB list. The PCB is then detached, which releases the memory used by the PCB.

27.5 tcp_mss Function

The `tcp_mss` function is called from two other functions:

1. from `tcp_output`, when a SYN segment is being sent, to include an MSS option, and
2. from `tcp_input`, when an MSS option is received in a SYN segment.

The `tcp_mss` function checks for a cached route to the destination and calculates the MSS to use for this connection.

Figure 27.7 shows the first part of `tcp_mss`, which acquires a route to the destination if one is not already held by the PCB.

```

1391 int                                     tcp_input.c
1392 tcp_mss(tp, offer)
1393 struct tcpcb *tp;
1394 u_int offer;
1395 {
1396     struct route *ro;
1397     struct rtentry *rt;
1398     struct ifnet *ifp;
1399     int rtt, mss;
1400     u_long bufsize;
1401     struct inpcb *inp;
1402     struct socket *so;
1403     extern int tcp_mssdflt;

1404     inp = tp->t_inpcb;
1405     ro = &inp->inp_route;

1406     if ((rt = ro->ro_rt) == (struct rtentry *) 0) {
1407         /* No route yet, so try to acquire one */
1408         if (inp->inp_faddr.sa_addr != INADDR_ANY) {
1409             ro->ro_dst.sa_family = AF_INET;
1410             ro->ro_dst.sa_len = sizeof(ro->ro_dst);
1411             ((struct sockaddr_in *) &ro->ro_dst)->sin_addr =
1412                 inp->inp_faddr;
1413             rtalloc(ro);
1414         }
1415         if ((rt = ro->ro_rt) == (struct rtentry *) 0)
1416             return (tcp_mssdflt);
1417     }
1418     ifp = rt->rt_ifp;
1419     so = inp->inp_socket;

```

Figure 27.7 `tcp_mss` function: acquire a route if one is not held by the PCB.

Acquire a route if necessary

1391-1417 If the socket does not have a cached route, `rtalloc` acquires one. The interface pointer associated with the outgoing route is saved in `ifp`. Knowing the outgoing interface is important, since its associated MTU can affect the MSS announced by TCP. If a route is not acquired, the default of 512 (`tcp_mssdflt`) is returned immediately.

The next part of `tcp_mss`, shown in Figure 27.8, checks whether the route has metrics associated with it; if so, the variables `t_rttmin`, `t_srtt`, and `t_rttvar` can be initialized from the metrics.


```

1420  /*                                     tcp_input.c
1421  * While we're here, check if there's an initial rtt
1422  * or rttvar. Convert from the route-table units
1423  * to scaled multiples of the slow timeout timer.
1424  */
1425  if (tp->t_srtt == 0 && (rtt = rt->rt_rmx.rmx_rtt)) {
1426      /*
1427       * XXX the lock bit for RTT indicates that the value
1428       * is also a minimum value; this is subject to time.
1429       */
1430      if (rt->rt_rmx.rmx_locks & RTV_RTT)
1431          tp->t_rttmin = rtt / (RTM_RTTUNIT / PR_SLOWHZ);
1432      tp->t_srtt = rtt / (RTM_RTTUNIT / (PR_SLOWHZ * TCP_RTT_SCALE));
1433
1434      if (rt->rt_rmx.rmx_rttvar)
1435          tp->t_rttvar = rt->rt_rmx.rmx_rttvar /
1436              (RTM_RTTUNIT / (PR_SLOWHZ * TCP_RTTVAR_SCALE));
1437      else
1438          /* default variation is +/- 1 rtt */
1439          tp->t_rttvar =
1440              TCP_RTTVAR_SCALE * TCP_RTTVAR_SCALE / TCP_RTT_SCALE;
1441
1442      TCPT_RANGESET(tp->t_rttcur,
1443                  ((tp->t_srtt >> 2) + tp->t_rttvar) >> 1,
1444                  tp->t_rttmin, TCPTV_REXMTMAX);
1445  }

```

Figure 27.8 tcp_rss function: check if the route has an associated RTT metric.

Initialize smoothed RTT estimator

1420-1432 If there are no RTT measurements yet for the connection (t_srtt is 0) and rmx_rtt is nonzero, the latter initializes the smoothed RTT estimator t_srtt . If the RTV_RTT bit in the routing metric lock flag is set, it indicates that rmx_rtt should also be used to initialize the minimum RTT for this connection (t_rttmin). We saw that $tcp_newtcpcb$ initializes t_rttmin to 2 ticks.

rmx_rtt (in units of microseconds) is converted to t_srtt (in units of ticks $\times 8$). This is the reverse of the conversion done in Figure 27.4. Notice that t_rttmin is set to one-eighth the value of t_srtt , since the former is not divided by the scale factor TCP_RTT_SCALE .

Initialize smoothed mean deviation estimator

1433-1439 If the stored value of rmx_rttvar is nonzero, it is converted from units of microseconds into ticks $\times 4$ and stored in t_rttvar . But if the value is 0, t_rttvar is set to t_rtt , that is, the variation is set to the mean. This defaults the variation to ± 1 RTT. Since the units of the former are ticks $\times 4$ and the units of the latter are ticks $\times 8$, the value of t_srtt is converted accordingly.

Calculate initial RTO

1440-1442 The current RTO is calculated and stored in `t_rttcur`, using the unscaled equation

$$RTO = srtt + 2 \times rttvar$$

A multiplier of 2, instead of 4, is used to calculate the first RTO. This is the same equation that was used in Figure 25.21. Substituting the scaling relationships we get

$$\begin{aligned} RTO &= \frac{t_srtt}{8} + 2 \times \frac{t_rttvar}{4} \\ &= \frac{t_srtt}{4} + t_rttvar \\ &= \frac{t_srtt}{2} \end{aligned}$$

which is the second argument to `TCPT_RANGESET`.

The next part of `tcp_mss`, shown in Figure 27.9, calculates the MSS.

```

1444 /*
1445  * if there's an mtu associated with the route, use it
1446  */
1447 if (rt->rt_rmx.rmx_mtu)
1448     mss = rt->rt_rmx.rmx_mtu - sizeof(struct tcpiphdr);
1449 else {
1450     mss = ifp->if_mtu - sizeof(struct tcpiphdr);
1451 #if (MCLBYTES & (MCLBYTES - 1)) == 0
1452     if (mss > MCLBYTES)
1453         mss &= ~(MCLBYTES - 1);
1454 #else
1455     if (mss > MCLBYTES)
1456         mss = mss / MCLBYTES * MCLBYTES;
1457 #endif
1458     if (!in_localaddr(inp->inp_faddr))
1459         mss = min(mss, tcp_mssdfilt);
1460 }

```

tcp_input.c

Figure 27.9 `tcp_mss` function: calculate MSS.

Use MSS from routing table MTU

1444-1450 If the MTU is set in the routing table, `mss` is set to that value. Otherwise `mss` starts at the value of the outgoing interface MTU minus 40 (the default size of the IP and TCP headers). For an Ethernet, `mss` would start at 1460.

Round MSS down to multiple of MCLBYTES

1451-1457 The goal of these lines of code is to reduce the value of `mss` to the next-lower multiple of the mbuf cluster size, if `mss` exceeds `MCLBYTES`. If the value of `MCLBYTES` (typically 1024 or 2048) logically ANDed with the value minus 1 equals 0, then `MCLBYTES` is a power of 2. For example, 1024 (0x400) logically ANDed with 1023 (0x3ff) is 0.

The value of `mss` is reduced to the next-lower multiple of `MCLBYTES` by clearing the appropriate number of low-order bits: if the cluster size is 1024, logically ANDing `mss` with the one's complement of 1023 (`0xfffffc00`) clears the low-order 10 bits. For an Ethernet, this reduces `mss` from 1460 to 1024. If the cluster size is 2048, logically ANDing `mss` with the one's complement of 2047 (`0xffff8000`) clears the low-order 11 bits. For a token ring with an MTU of 4464, this reduces the value of `mss` from 4424 to 4096. If `MCLBYTES` is not a power of 2, the rounding down to the next-lower multiple of `MCLBYTES` is done with an integer division followed by a multiplication.

Check if destination local or nonlocal

1458-1459 If the foreign IP address is not local (`in_localaddr` returns 0), and if `mss` is greater than 512 (`tcp_mssdflt`), it is set to 512.

Whether an IP address is "local" or not depends on the value of the global `subnetsarelocal`, which is initialized from the symbol `SUBNETSARELOCAL` when the kernel is compiled. The default value is 1, meaning that an IP address with the same network ID as one of the host's interfaces is considered local. If the value is 0, an IP address must have the same network ID and the same subnet ID as one of the host's interfaces to be considered local.

This minimization for nonlocal hosts is an attempt to avoid fragmentation across wide-area networks. It is a historical artifact from the ARPANET when the MTU across most WAN links was 1006. As discussed in Section 11.7 of Volume 1, most WANs today support an MTU of 1500 or greater. See also the discussion of the path MTU discovery feature (RFC 1191 [Mogul and Deering 1990]), in Section 24.2 of Volume 1. Net/3 does not support path MTU discovery.

The final part of `tcp_mss` is shown in Figure 27.10.

Other end's MSS is upper bound

1461-1472 The argument `offer` is nonzero when this function is called from `tcp_input`, and its value is the MSS advertised by the other end. If the value of `mss` is greater than the value advertised by the other end, it is set to the value of `offer`. For example, if the function calculates an `mss` of 1024 but the advertised value from the other end is 512, `mss` must be set to 512. Conversely, if `mss` is calculated as 536 (say the outgoing MTU is 576) and the other end advertises an MSS of 1460, TCP will use 536. TCP can always use a value less than the advertised MSS, but it can't exceed the advertised value. The argument `offer` is 0 when this function is called by `tcp_output` to send an MSS option. The value of `mss` is also lower-bounded by 32.

1473-1483 If the value of `mss` has decreased from the default set by `tcp_newtcpcb` in the variable `t_maxseg` (512), or if TCP is processing a received MSS option (`offer` is nonzero), the following steps occur. First, if the value of `rmx_sendpipe` has been stored for the route, its value will be used as the send buffer high-water mark (Figure 16.4). If the buffer size is less than `mss`, the smaller value is used. This should never happen unless the application explicitly sets the send buffer size to a small value, or the administrator sets `rmx_sendpipe` to a small value, since the high-water mark of the send buffer defaults to 8192, larger than most values for the MSS.

```

1461      /*
1462      * The current mss, t_maxseg, was initialized to the default value
1463      * of 512 (tcp_mssdflt) by tcp_newtcpcb().
1464      * If we compute a smaller value, reduce the current mss.
1465      * If we compute a larger value, return it for use in sending
1466      * a max seg size option, but don't store it for use
1467      * unless we received an offer at least that large from peer.
1468      * However, do not accept offers under 32 bytes.
1469      */
1470      if (offer)
1471          mss = min(mss, offer);
1472      mss = max(mss, 32); /* sanity */
1473      if (mss < tp->t_maxseg || offer != 0) {
1474          /*
1475           * If there's a pipesize, change the socket buffer
1476           * to that size. Make the socket buffers an integral
1477           * number of mss units; if the mss is larger than
1478           * the socket buffer, decrease the mss.
1479           */
1480           if ((bufsize = rt->rt_rmx.rmx_sendpipe) == 0)
1481               bufsize = so->so_snd.sb_hiwat;
1482           if (bufsize < mss)
1483               mss = bufsize;
1484           else {
1485               bufsize = roundup(bufsize, mss);
1486               if (bufsize > sb_max)
1487                   bufsize = sb_max;
1488               (void) sbreserve(&so->so_snd, bufsize);
1489           }
1490           tp->t_maxseg = mss;
1491           if ((bufsize = rt->rt_rmx.rmx_rcvpipe) == 0)
1492               bufsize = so->so_rcv.sb_hiwat;
1493           if (bufsize > mss) {
1494               bufsize = roundup(bufsize, mss);
1495               if (bufsize > sb_max)
1496                   bufsize = sb_max;
1497               (void) sbreserve(&so->so_rcv, bufsize);
1498           }
1499       }
1500       tp->snd_cwnd = mss;
1501       if (rt->rt_rmx.rmx_ssthresh) {
1502           /*
1503            * There's some sort of gateway or interface
1504            * buffer limit on the path. Use this to set
1505            * the slow start threshold, but set the
1506            * threshold to no less than 2*mss.
1507            */
1508           tp->snd_ssthresh = max(2 * mss, rt->rt_rmx.rmx_ssthresh);
1509       }
1510       return (mss);
1511     }

```

Figure 27.10 tcp_mss function: complete processing.

Round buffer sizes to multiple of MSS

1484-1489 The send buffer size is rounded up to the next integral multiple of the MSS, bounded by the value of `sb_max` (262,144 on Net/3, which is 256×1024). The socket's high-water mark is set by `sb_reserve`. For example, the default high-water mark is 8192, but for a local TCP connection on an Ethernet with a cluster size of 2048 (i.e., an MSS of 1460) this code increases the high-water mark to 8760 (which is 6×1460). But for a nonlocal connection with an MSS of 512, the high-water mark is left at 8192.

1490 The value of `t_maxseg` is set, either because it decreased from the default (512) or because an MSS option was received from the other end.

1491-1499 The same logic just applied to the send buffer is also applied to the receive buffer.

Initialize congestion window and slow start threshold

1500-1509 The value of the congestion window, `snd_cwnd`, is set to one segment. If the `rmx_ssthresh` value in the routing table is nonzero, the slow start threshold (`snd_ssthresh`) is set to that value, but the value must not be less than two segments.

1510 The value of `mss` is returned by the function. `tcp_input` ignores this value in Figure 28.10 (since it received an MSS from the other end), but `tcp_output` sends this value as the announced MSS in Figure 26.23.

Example

Let's go through an example of a TCP connection establishment and the operation of `tcp_mss`, since it can be called twice: once when the SYN is sent and once when a SYN is received with an MSS option.

1. The socket is created and `tcp_newtcpcb` sets `t_maxseg` to 512.
2. The process calls `connect`, and `tcp_output` calls `tcp_mss` with an `offer` argument of 0, to include an MSS option with the SYN. Assuming a local destination, an Ethernet LAN, and an `mbuf` cluster size of 2048, `mss` is set to 1460 by the code in Figure 27.9. Since `offer` is 0, Figure 27.10 leaves the value as 1460 and this is the function's return value. The buffer sizes aren't modified, since 1460 is larger than the default (512) and a value hasn't been received from the other end yet. `tcp_output` sends an MSS option announcing a value of 1460.
3. The other end replies with its SYN, announcing an MSS of 1024. `tcp_input` calls `tcp_mss` with an `offer` argument of 1024. The logic in Figure 27.9 still yields a value of 1460 for `mss`, but the call to `min` at the beginning of Figure 27.10 reduces this to 1024. Since the value of `offer` is nonzero, the buffer sizes are rounded up to the next integral multiple of 1024 (i.e., they're left at 8192). `t_maxseg` is set to 1024.

It might appear that the logic of `tcp_mss` is flawed: TCP announces an MSS of 1460 but receives an MSS of 1024 from the other end. While TCP is restricted to sending 1024-byte segments, the other end is free to send 1460-byte segments. We might think that the send buffer should be a multiple of 1024, but the receive buffer should be a multiple of 1460. Yet the code in Figure 27.10 sets both buffer sizes based on the received MSS. The reasoning is that even if TCP announces an MSS of 1460, since it receives an MSS of 1024 from the other end, the other end probably won't send 1460-byte segments, but will restrict itself to 1024-byte segments.

27.6 tcp_ctlinput Function

Recall from Figure 22.32 that `tcp_ctlinput` processes five types of ICMP errors: destination unreachable, parameter problem, source quench, time exceeded, and redirects. All redirects are passed to both TCP and UDP. For the other four errors, `tcp_ctlinput` is called only if a TCP segment caused the error.

`tcp_ctlinput` is shown in Figure 27.11. It is similar to `udp_ctlinput`, shown in Figure 23.30.

```

355 void                                     tcp_subr.c
356 tcp_ctlinput(cmd, sa, ip)
357 int    cmd;
358 struct sockaddr *sa;
359 struct ip *ip;
360 {
361     struct tcphdr *th;
362     extern struct in_addr zeroin_addr;
363     extern u_char inetctlerrmap[];
364     void (*notify) (struct inpcb *, int) = tcp_notify;

365     if (cmd == PRC_QUENCH)
366         notify = tcp_quench;
367     else if (!PRC_IS_REDIRECT(cmd) &&
368             ((unsigned) cmd > PRC_NCMSDS || inetctlerrmap[cmd] == 0))
369         return;
370     if (ip) {
371         th = (struct tcphdr *) ((caddr_t) ip + (ip->ip_hl << 2));
372         in_pcbnotify(&tc, sa, th->th_dport, ip->ip_src, th->th_sport,
373                    cmd, notify);
374     } else
375         in_pcbnotify(&tc, sa, 0, zeroin_addr, 0, cmd, notify);
376 }

```

Figure 27.11 `tcp_ctlinput` function.

365-366 The only difference in the logic from `udp_ctlinput` is how an ICMP source quench error is handled. UDP ignores these errors since the `PRC_QUENCH` entry of `inetctlerrmap` is 0. TCP explicitly checks for this error, changing the `notify` function from its default of `tcp_notify` to `tcp_quench`.

27.7 tcp_notify Function

`tcp_notify` is called by `tcp_ctlinput` to handle destination unreachable, parameter problem, time exceeded, and redirect errors. This function is more complicated than its UDP counterpart, since TCP must intelligently handle soft errors for an established connection. Figure 27.12 shows the `tcp_notify` function.

```

328 void
329 tcp_notify(inp, error)
330 struct inpcb *inp;
331 int error;
332 {
333     struct tcpcb *tp = (struct tcpcb *) inp->inp_ppcb;
334     struct socket *so = inp->inp_socket;
335     /*
336      * Ignore some errors if we are hooked up.
337      * If connection hasn't completed, has retransmitted several times,
338      * snd receives a second error, give up now. This is better
339      * than waiting a long time to establish a connection that
340      * can never complete.
341      */
342     if (tp->t_state == TCPS_ESTABLISHED &&
343         (error == EHOSTUNREACH || error == ENETUNREACH ||
344          error == EHOSTDOWN)) {
345         return;
346     } else if (tp->t_state < TCPS_ESTABLISHED && tp->t_extshift > 3 &&
347                tp->t_softerror)
348         so->so_error = error;
349     else
350         tp->t_softerror = error;
351     wakeup((caddr_t) & so->so_timeo);
352     sowakeup(so);
353     sowakeup(so);
354 }

```

Figure 27.12 tcp_notify function.

328-345 If the connection is ESTABLISHED, the errors EHOSTUNREACH, ENETUNREACH, and EHOSTDOWN are ignored.

This handling of these three errors is new with 4.4BSD. Net/2 and earlier releases recorded these errors in the connection's soft error variable (`t_softerror`), and the error was reported to the process should the connection eventually fail. Recall that `tcp_xmit_timer` resets this variable to 0 when an ACK is received for a segment that hasn't been retransmitted.

346-353 If the connection is not yet established, TCP has retransmitted the current segment four or more times, and an error has already been recorded in `t_softerror`, the current error is recorded in the socket's `so_error` variable. By setting this socket variable, the socket becomes readable and writable if the process calls `select`. Otherwise the current error is just saved in `t_softerror`. We saw that `tcp_drop` sets the socket error to this saved value if the connection is subsequently dropped because of a timeout. Any processes waiting to receive or send on the socket are then awakened to receive the error.

27.8 tcp_quench Function

`tcp_quench`, which is shown in Figure 27.13, is called by `tcp_ctlinput` when a source quench is received for the connection, and by `tcp_output` (Figure 26.32) when `ip_output` returns `ENOBUFS`.

```

381 void                                     tcp_subr.c
382 tcp_quench(inp, errno)
383 struct inpcb *inp;
384 int      errno;
385 {
386     struct tcpcb *tp = inrtotpcb(inp);
387     if (tp)
388         tp->snd_cwnd = tp->t_maxseg;
389 }

```

Figure 27.13 `tcp_quench` function.

The congestion window is set to one segment, causing slow start to take over. The slow start threshold is not changed (as it is when `tcp_timers` handles a retransmission timeout), so the window will open up exponentially until `snd_ssthresh` is reached, or congestion occurs.

27.9 TCP_REASS Macro and `tcp_reass` Function

TCP segments can arrive out of order, and it is TCP's responsibility to place the misordered segments into the correct order for presentation to the process. For example, if a receiver advertises a window of 4096 with byte number 0 as the next expected byte, and receives a segment with bytes 0–1023 (an in-order segment) followed by a segment with bytes 2048–3071, this second segment is out of order. TCP does not discard the out-of-order segment if it is within the receive window. Instead it places the segment on the reassembly list for the connection, waiting for the missing segment to arrive (with bytes 1024–2047), at which time it can acknowledge bytes 1024–3071 and pass these 2048 bytes to the process. In this section we examine the code that manipulates the TCP reassembly queue, before discussing `tcp_input` in the next two chapters.

If we assume that a single mbuf contains the IP header, TCP header, and 4 bytes of TCP data (as shown in the left half of Figure 2.14) we would have the arrangement shown in Figure 27.14. We also assume the data bytes are sequence numbers 7, 8, 9, and 10.

The `ipovly` and `tcphdr` structures form the `tcpiphdr` structure, which we showed in Figure 24.12. We showed a picture of the `tcphdr` structure in Figure 24.10. In Figure 27.14 we show only the variables used in the reassembly: `ti_next`, `ti_prev`, `ti_len`, `ti_sport`, `ti_dport`, and `ti_seq`. The first two are pointers that form a doubly linked list of all the out-of-order segments for a given connection. The head of this list is the TCP control block for the connection: the `seg_next` and `seg_prev` members, which are the first two members of the structure. The `ti_next` and `ti_prev`

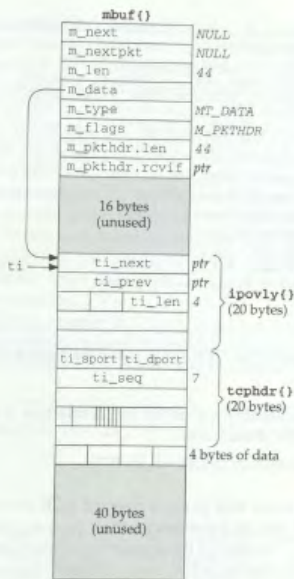


Figure 27.14 Example mbuf with IP and TCP headers and 4 bytes of data.

pointers overlay the first 8 bytes of the IP header, which aren't needed once the datagram reaches TCP. `ti_len` is the length of the TCP header, and is calculated and stored by TCP before verifying the TCP checksum.

TCP_REASS Macro

When data is received by `tcp_input`, the macro `TCP_REASS`, shown in Figure 27.15, is invoked to place the data onto the connection's reassembly queue. This macro is called from only one place: see Figure 29.22.

⁵⁴⁻⁶³ `tp` is a pointer to the TCP control block for the connection and `ti` is a pointer to the `tcphdr` structure for the received segment. If the following three conditions are all true:

1. this segment is in-order (the sequence number `ti_seq` equals the next expected sequence number for the connection, `rcv_nxt`), and

```

53 #define TCP_REASS(tp, ti, m, so, flags) { \
54     if ((ti)->ti_seq == (tp)->rcv_nxt && \
55         (tp)->seg_next == (struct tcpiphdr *) (tp) && \
56         (tp)->t_state == TCPS_ESTABLISHED) { \
57         tp->t_flags |= TF_DELACK; \
58         (tp)->rcv_nxt += (ti)->ti_len; \
59         flags = (ti)->ti_flags & TH_FIN; \
60         tcpstat.tcps_rcvpack++; \
61         tcpstat.tcps_rcvbyte += (ti)->ti_len; \
62         sbappend((so)->so_rcv, (m)); \
63         sorwakeup(so); \
64     } else { \
65         (flags) = tcp_reass((tp), (ti), (m)); \
66         tp->t_flags |= TF_ACKNOW; \
67     } \
68 }

```

tcp_input.c

tcp_input.c

Figure 27.15 TCP_REASS macro: add data to reassembly queue for connection.

2. the reassembly queue for the connection is empty (*seg_next* points to itself, not some mbuf), and
3. the connection is ESTABLISHED.

the following steps take place: a delayed ACK is scheduled, *rcv_nxt* is updated with the amount of data in the segment, the *flags* argument is set to *TH_FIN* if the FIN flag is set in the TCP header of the segment, two statistics are updated, the data is appended to the socket's receive buffer, and any receiving processes waiting for the socket are awakened.

The reason all three conditions must be true is that, first, if the data is out of order, it must be placed onto the connection's reassembly queue and the "preceding" segments must be received before anything can be passed to the process. Second, even if the data is in order, if there is out-of-order data already on the reassembly queue, there's a chance that the new segment might fill a hole, allowing the received segment and one or more segments on the queue to all be passed to the process. Third, it is OK for data to arrive with a SYN segment that establishes a connection, but that data cannot be passed to the process until the connection is ESTABLISHED—any such data is just added to the reassembly queue when it arrives.

64-67 If these three conditions are not all true, the `TCP_REASS` macro calls the function `tcp_reass` to add the segment to the reassembly queue. Since the segment is either out of order, or the segment might fill a hole from previously received out-of-order segments, an immediate ACK is scheduled. One important feature of TCP is that a receiver should generate an immediate ACK when an out-of-order segment is received. This aids the *fast retransmit* algorithm (Section 29.4).

Before looking at the code for the `tcp_reass` function, we need to explain what's done with the two port numbers in the TCP header in Figure 27.14, *ti_sport* and

`ti_dport`. Once the TCP control block is located and `tcp_reass` is called, these two port numbers are no longer needed. Therefore, when a TCP segment is placed on a reassembly queue, the address of the corresponding mbuf is stored over these two port numbers. In Figure 27.14 this isn't needed, because the IP and TCP headers are in the data portion of the mbuf, so the `dtom` macro works. But recalling our discussion of `m_pullup` in Section 2.6, if the IP and TCP headers are in a cluster (as in Figure 2.16, which is the normal case for a full-sized TCP segment), the `dtom` macro doesn't work. We mentioned in that section that TCP stores its own back pointer from the TCP header to the mbuf, and that back pointer is stored over the two TCP port numbers.

Figure 27.16 shows an example of this technique with two out-of-order segments for a connection, each segment stored in an mbuf cluster. The head of the doubly linked list of out-of-order segments is the `seg_next` member of the control block for this connection. To simplify the figure we don't show the `seg_prev` pointer and the `ti_next` pointer of the last segment on the list.

The next expected sequence number is 1 (`rcv_nxt`) but we assume that segment was lost. The next two segments have been received, containing bytes 1461–4380, but they are out of order. The segments were placed into clusters by `m_devget`, as shown in Figure 2.16.

The first 32 bits of the TCP header contain a back pointer to the corresponding mbuf. This back pointer is used in the `tcp_reass` function, shown next.

tcp_reass Function

Figure 27.17 shows the first part of the `tcp_reass` function. The arguments are: `tp`, a pointer to the TCP control block for the received segment; `ti`, a pointer to the IP and TCP headers of the received segment; and `m`, a pointer to the mbuf chain for the received segment. As mentioned earlier, `ti` can point into the data area of the mbuf pointed to by `m`, or `ti` can point into a cluster.

69-83 We'll see that `tcp_input` calls `tcp_reass` with a null `ti` pointer when a SYN is acknowledged (Figures 28.20 and 29.2). This means the connection is now established, and any data that might have arrived with the SYN (which `tcp_reass` had to queue earlier) can now be passed to the application. Data that arrives with a SYN cannot be passed to the process until the connection is established. The label `present` is in Figure 27.23.

84-90 Go through the list of segments for this connection, starting at `seg_next`, to find the first one with a sequence number that is greater than the received sequence number (`ti_seq`). Note that the `if` statement is the entire body of the `for` loop.

Figure 27.18 shows an example with two out-of-order segments already on the queue when a new segment arrives. We show the pointer `q` pointing to the next segment on the list, the one with bytes 10–15. In this figure we also show the two pointers `ti_next` and `ti_prev`, the starting sequence number (`ti_seq`), the length (`ti_len`), and the sequence numbers of the data bytes. With the small segments we show, each segment is probably in a single mbuf, as in Figure 27.14.

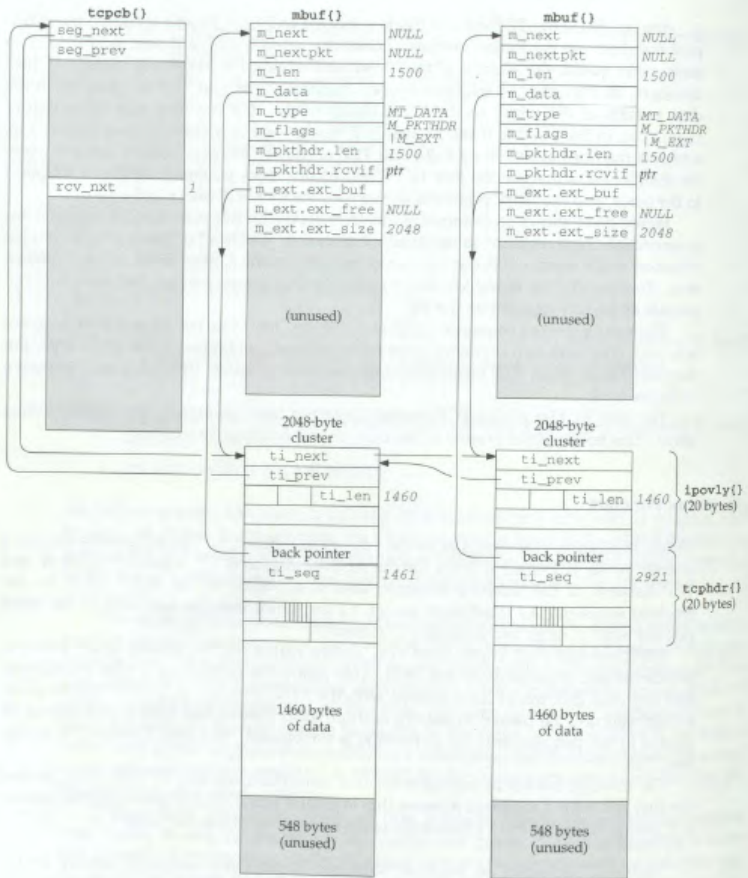


Figure 27.16 Two out-of-order TCP segments stored in mbuf clusters.


```

69 int                                     tcp_input.c
70 tcp_reass(tp, ti, m)
71 struct tcpcb *tp;
72 struct tcpiphdr *ti;
73 struct mhuf *m;
74 {
75     struct tcpiphdr *q;
76     struct socket *so = tp->t_inpcb->inp_socket;
77     int flags;
78     /*
79     * Call with ti==0 after become established to
80     * force pre-ESTABLISHED data up to user socket.
81     */
82     if (ti == 0)
83         goto present;
84     /*
85     * Find a segment that begins after this one does.
86     */
87     for (q = tp->seg_next; q != (struct tcpiphdr *) tp;
88         q = (struct tcpiphdr *) q->ti_next)
89         if (SEQ_GT(q->ti_seq, ti->ti_seq))
90             break;

```

Figure 27.17 tcp_reass function: first part.

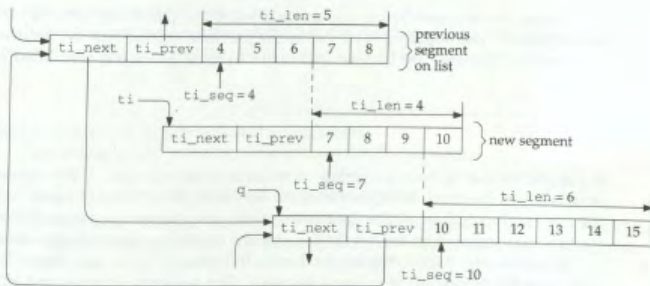


Figure 27.18 Example of TCP reassembly queue with overlapping segments.

The next part of `tcp_reass` is shown in Figure 27.19.

```

91  /* tcp_input.c
92  * If there is a preceding segment, it may provide some of
93  * our data already. If so, drop the data from the incoming
94  * segment. If it provides all of our data, drop us.
95  */
96  if ((struct tcpiphdr *) q->ti_prev != (struct tcpiphdr *) tp) {
97      int i;
98      q = (struct tcpiphdr *) q->ti_prev;
99      /* conversion to int (in i) handles seq wraparound */
100     i = q->ti_seq + q->ti_len - ti->ti_seq;
101     if (i > 0) {
102         if (i >= ti->ti_len) {
103             tcpstat.tcps_rcvduppack++;
104             tcpstat.tcps_rcvdupbyte += ti->ti_len;
105             m_freem(m);
106             return (0);
107         }
108         m_adj(m, i);
109         ti->ti_len -= i;
110         ti->ti_seq -= i;
111     }
112     q = (struct tcpiphdr *) (q->ti_next);
113 }
114 tcpstat.tcps_rcvooopack++;
115 tcpstat.tcps_rcvooobyte += ti->ti_len;
116 REASS_MBUF(ti) = m; /* XXX */
tcp_input.c

```

Figure 27.19 `tcp_reass` function: second part.

91-107 If there is a segment before the one pointed to by `q`, that segment may overlap the new segment. The pointer `q` is moved to the previous segment on the list (the one with bytes 4-8 in Figure 27.18) and the number of bytes of overlap is calculated and stored in `i`:

```

i = q->ti_seq + q->ti_len - ti->ti_seq;
  = 4 + 5 - 7
  = 2

```

If `i` is greater than 0, there is overlap, as we have in our example. If the number of bytes of overlap in the previous segment on the list (`i`) is greater than or equal to the size of the new segment, then all the data bytes in the new segment are already contained in the previous segment on the list. In this case the duplicate segment is discarded.

108-112 If there is only partial overlap (as there is in Figure 27.18), `m_adj` discards `i` bytes of data from the beginning of the new segment. The sequence number and length of the new segment are updated accordingly. `q` is moved to the next segment on the list. Figure 27.20 shows our example at this point.

116 The address of the mbuf `m` is stored in the TCP header, over the source and destination TCP ports. We mentioned earlier in this section that this provides a back pointer

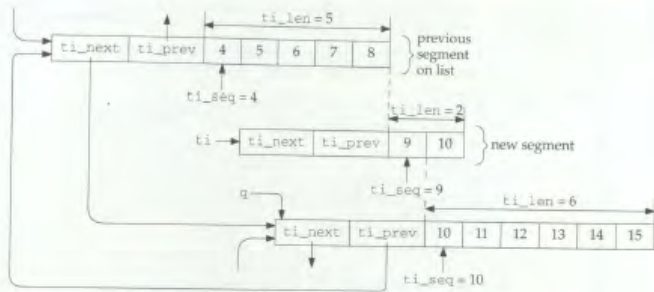


Figure 27.20 Update of Figure 27.18 after bytes 7 and 8 have been removed from new segment.

from the TCP header to the mbuf, in case the TCP header is stored in a cluster, meaning that the macro `atom` won't work. The macro `REASS_MBUF` is

```
#define REASS_MBUF(ti) (*(struct mbuf **)&(ti->ti_t))
```

`ti_t` is the `tcphdr` structure (Figure 24.12) and the first two members of the structure are the two 16-bit port numbers. The comment `XXX` in Figure 27.19 is because this hack assumes that a pointer fits in the 32 bits occupied by the two port numbers.

The third part of `tcp_reass` is shown in Figure 27.21. It removes any overlap from the next segment in the queue.

117-135 If there is another segment on the list, the number of bytes of overlap between the new segment and that segment is calculated in `i`. In our example we have

$$\begin{aligned} i &= 9 + 2 - 10 \\ &= 1 \end{aligned}$$

since byte number 10 overlaps the two segments.

Depending on the value of `i`, one of three conditions exists:

1. If `i` is less than or equal to 0, there is no overlap.
2. If `i` is less than the number of bytes in the next segment (`q->ti_len`), there is partial overlap and `m_adj` removes the first `i` bytes from the next segment on the list.
3. If `i` is greater than or equal to the number of bytes in the next segment, there is complete overlap and that next segment on the list is deleted.

136-139 The new segment is inserted into the reassembly list for this connection by `insque`. Figure 27.22 shows the state of our example at this point.

```

117  /*
118  * While we overlap succeeding segments trim them or,
119  * if they are completely covered, dequeue them.
120  */
121  while (q != (struct tcpiphdr *) tp) {
122      int i = (ti->ti_seq + ti->ti_len) - q->ti_seq;
123      if (i <= 0)
124          break;
125      if (i < q->ti_len) {
126          q->ti_seq += i;
127          q->ti_len -= i;
128          m_adj(REASS_MBUF(q), i);
129          break;
130      }
131      q = (struct tcpiphdr *) q->ti_next;
132      m = REASS_MBUF((struct tcpiphdr *) q->ti_prev);
133      remque(q->ti_prev);
134      m_freem(m);
135  }
136  /*
137  * Stick new segment in its place.
138  */
139  insque(ti, q->ti_prev);

```

Figure 27.21 tcp_reass function: third part.

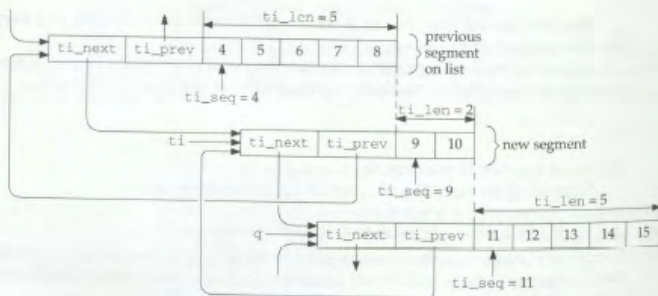


Figure 27.22 Update of Figure 27.20 after removal of all overlapping bytes.

Figure 27.23 shows the final part of `tcp_reass`. It passes the data to the process, if possible.


```

140 present:                                     tcp_input.c
141 /*
142  * Present data to user, advancing rcv_nxt through
143  * completed sequence space.
144  */
145 if (TCPS_HAVERCVDSYN(tp->t_state) == 0)
146     return (0);
147 ti = tp->seg_next;
148 if (ti == (struct tcpiphdr *) tp || ti->ti_seq != tp->rcv_nxt)
149     return (0);
150 if (tp->t_state == TCPS_SYN_RECEIVED && ti->ti_len)
151     return (0);
152 do {
153     tp->rcv_nxt += ti->ti_len;
154     flags = ti->ti_flags & TR_FIN;
155     remque(ti);
156     m = REASS_MBUF(ti);
157     ti = (struct tcpiphdr *) ti->ti_next;
158     if (so->so_state & SS_CANTRCVMORE)
159         m_freem(m);
160     else
161         sbappend(&so->so_rcv, m);
162 } while (ti != (struct tcpiphdr *) tp && ti->ti_seq == tp->rcv_nxt);
163 sorwakeup(so);
164 return (flags);
165 }

```

Figure 27.23 tcp_reass function: fourth part.

145-146 If the connection has not received a SYN (i.e., it is in the LISTEN or SYN_SENT state), data cannot be passed to the process and the function returns. When this function is called by TCP_REASS, the return value of 0 is stored in the flags argument to the macro. This can have the side effect of clearing the FIN flag. We'll see that this side effect is a possibility when TCP_REASS is invoked in Figure 29.22, and the received segment contains a SYN, FIN, and data (not a typical segment, but valid).

147-149 ti starts at the first segment on the list. If the list is empty, or if the starting sequence number of the first segment on the list (ti->ti_seq) does not equal the next receive sequence number (rcv_nxt), the function returns a value of 0. If the second condition is true, there is still a hole in the received data starting with the next expected sequence number. For instance, in our example (Figure 27.22), if the segment with bytes 4-8 is the first on the list but rcv_nxt equals 2, bytes 2 and 3 are still missing, so bytes 4-15 cannot be passed to the process. The return of 0 turns off the FIN flag (if set), because one or more data segments are still missing, so a received FIN cannot be processed yet.

150-151 If the state is SYN_RCVD and the length of the segment is nonzero, the function returns a value of 0. If both of these conditions are true, the socket is a listening socket that has received in-order data with the SYN. The data is left on the connection's queue, waiting for the three-way handshake to complete.

152-164 This loop starts with the first segment on the list (which is known to be in order) and appends it to the socket's receive buffer. `rcv_nxt` is incremented by the number of bytes in the segment. The loop stops when the list is empty or when the sequence number of the next segment on the list is out of order (i.e., there is a hole in the sequence space). When the loop terminates, the `Flags` variable (which becomes the return value of the function) is 0 or `TH_FIN`, depending on whether the final segment placed in the socket's receive buffer has the FIN flag set or not.

After all the mbufs have been placed onto the socket's receive buffer, `sockwakep` wakes any process waiting for data to be received on the socket.

27.10 tcp_trace Function

In `tcp_output`, before sending a segment to IP for output, we saw the following call to `tcp_trace` in Figure 26.32:

```
if (so->so_options & SO_DEBUG)
    tcp_trace(TA_OUTPUT, tp->t_state, tp, ti, 0);
```

This call adds a record to a circular buffer in the kernel that can be examined with the `trpt(8)` program. Additionally, if the kernel is compiled with `TCPDEBUG` defined, and if the variable `tcpconsdebug` is nonzero, information is output on the system console.

Any process can set the `SO_DEBUG` socket option for a TCP socket, causing the information to be stored in the kernel's circular buffer. But `trpt` must read the kernel memory (`/dev/kmem`) to fetch this information, and this often requires special privileges.

The `SO_DEBUG` socket option can be set for any type of socket (e.g., UDP or raw IP), but TCP is the only protocol that looks at the option.

The information saved by the kernel is a `tcp_debug` structure, shown in Figure 27.24.

```

35 struct tcp_debug {
36     u_time  td_time;           /* lptime(): ms since midnight, UTC */
37     short  td_act;           /* TA_XXX value (Figure 27.25) */
38     short  td_ostate;        /* old state */
39     caddr_t td_tcb;          /* addr of TCP connection block */
40     struct tcpiphdr td_ti;    /* IP and TCP headers */
41     short  td_req;           /* PRU_XXX value for TA_USER */
42     struct tcpcb td_cb;       /* TCP connection block */
43 };
44
45 #define TCP_NDEBUG 100
46 struct tcp_debug tcp_debug[TCP_NDEBUG];
47 int    tcp_debug;

```

tcp_debug.h

Figure 27.24 `tcp_debug` structure.

35-43 This is a large structure (196 bytes), since it contains two other structures: the `tcpiphdr` structure with the IP and TCP headers; and the `tcpcb` structure, the entire TCP control block. Since the entire TCP control block is saved, any variable in the

control block can be printed by `trpt`. Also, if `trpt` doesn't print the variable we're interested in, we can modify the source code (it is available with the Net/3 release) to print whatever information we would like from the control block. The RTT variables in Figure 25.28 were obtained using this technique.

53-55 We also show the declaration of the array `tcp_debug`, which is used as the circular buffer. The index into the array (`tcp_debx`) is initialized to 0. This array occupies almost 20,000 bytes.

There are only four calls to `tcp_trace` in the kernel. Each call stores a different value in the `td_act` member of the structure, as shown in Figure 27.25.

td_act	Description	Reference
TA_DROP	from <code>tcp_input</code> , when input segment is dropped	Figure 29.27
TA_INPUT	after input processing complete, before call to <code>tcp_output</code>	Figure 29.26
TA_OUTPUT	before calling <code>ip_output</code> to send segment	Figure 26.32
TA_USER	from <code>tcp_usrreq</code> , after processing <code>PRU_XXX</code> request	Figure 30.1

Figure 27.25 `td_act` values and corresponding call to `tcp_trace`

Figure 27.26 shows the main body of the `tcp_trace` function. We omit the code that outputs directly to the console.

48-137 `ostate` is the old state of the connection, when the function was called. By saving this value and the new state of the connection (which is in the control block) we can see the state transition that occurred. In Figure 27.25, `TA_OUTPUT` doesn't change the state of the connection, but the other three calls can change the state.

Sample Output

Figure 27.27 shows the first four lines of `tcpdump` output corresponding to the three-way handshake and the first data segment from the example in Section 25.12. (Appendix A of Volume 1 provides additional details on the `tcpdump` output format.)

```

1 0.0          bsd1.1025 > vangogh.discard: S 20288001:20288001(0)
                               win 4096 <mss 512>
2 0.362719 (0.3627)  vangogh.discard > bsd1.1025: S 3202722817:3202722817(0)
                               ack 20288002 win 8192
                               <mss 512>
3 0.364316 (0.0016)  bsd1.1025 > vangogh.discard: . ack 1 win 4096
4 0.415859 (0.0515)  bsd1.1025 > vangogh.discard: . 1513(512) ack 1 win 4096

```

Figure 27.27 `tcpdump` output from example in Figure 25.28.

Figure 27.28 shows the corresponding output from `trpt`.

This output contains a few changes from the normal `trpt` output. The 32-bit decimal sequence numbers are printed as unsigned values (`trpt` incorrectly prints them as signed numbers). Some values printed by `trpt` in hexadecimal have been output in decimal. The values from `t_rtt` through `t_rxtcur` were added to `trpt` by the authors, for Figure 25.28.

```

48 void
49 tcp_trace(act, ostate, tp, ti, req)
50 short act, ostate;
51 struct tcpcb *tp;
52 struct tcphdr *ti;
53 int req;
54 {
55     tcp_seq seq, ack;
56     int len, flags;
57     struct tcp_debug *td = &tcp_debug[tcp_debug++];
58     if (tcp_debug == TCP_NDEBUG)
59         tcp_debug = 0; /* circle back to start */
60     td->td_time = iptime();
61     td->td_act = act;
62     td->td_ostate = ostate;
63     td->td_tcb = (caddr_t) tp;
64     if (tp)
65         td->td_cb = *tp; /* structure assignment */
66     else
67         bzero((caddr_t) &td->td_cb, sizeof(*tp));
68     if (ti)
69         td->td_ti = *ti; /* structure assignment */
70     else
71         bzero((caddr_t) &td->td_ti, sizeof(*ti));
72     td->td_req = req;
73 #ifdef TCPDEBUG
74     if (tcpconsdebug == 0)
75         return;
76
77     /* output information on console */
78
79 #endif
80 }

```

Figure 27.26 tcp_trace function: save information in kernel's circular buffer.

At time 953738 the SYN is sent. Notice that only the lower 6 digits of the millisecond time are output—it would take 8 digits to represent 1 minute before midnight. The ending sequence number that is output is wrong (20288005). Four bytes are sent with the SYN, but these are the MSS option, not data. The retransmit timer is 6 seconds (REXMT) and the keepalive timer is 75 seconds (KEEP). These timer values are in 500-ms ticks. The value of 1 for `t_rtt` means this segment is being timed for an RTT measurement.

This SYN segment is sent in response to the process calling `connect`. One millisecond later the trace record for this system call is added to the kernel's buffer. Even though the call to `connect` generates the SYN segment, since the call to `tcp_trace`


```

953738 SYN_SENT: output 20288001:20288005(4) #0 (win=4096)
<SYN> -> SYN_SENT
rcv_nxt 0, rcv_wnd 0
snd_una 20288001, snd_nxt 20288002, snd_max 20288002
snd_wll 0, snd_wl2 0, snd_wnd 0
REXMT=12 (t_rxtshift=0), KEEP=150
t_rtt=1, t_srtt=0, t_rttvar=24, t_rxtcur=12

953739 CLOSED: user CONNECT -> SYN_SENT
rcv_nxt 0, rcv_wnd 0
snd_una 20288001, snd_nxt 20288002, snd_max 20288002
snd_wll 0, snd_wl2 0, snd_wnd 0
REXMT=12 (t_rxtshift=0), KEEP=150
t_rtt=1, t_srtt=0, t_rttvar=24, t_rxtcur=12

954103 SYN_SENT: input 3202722817:3202722817(0) #20288002 (win=8192)
<SYN,ACK> -> ESTABLISHED
rcv_nxt 3202722818, rcv_wnd 4096
snd_una 20288002, snd_nxt 20288002, snd_max 20288002
snd_wll 3202722818, snd_wl2 20288002, snd_wnd 8192
KEEP=14400
t_rtt=0, t_srtt=16, t_rttvar=4, t_rxtcur=6

954103 ESTABLISHED: output 20288002:20288002(0) #3202722818 (win=4096)
<ACK> -> ESTABLISHED
rcv_nxt 3202722818, rcv_wnd 4096
snd_una 20288002, snd_nxt 20288002, snd_max 20288002
snd_wll 3202722818, snd_wl2 20288002, snd_wnd 8192
KEEP=14400
t_rtt=0, t_srtt=16, t_rttvar=4, t_rxtcur=6

954153 ESTABLISHED: output 20288002:20288514(512) #3202722818 (win=4096)
<ACK> -> ESTABLISHED
rcv_nxt 3202722818, rcv_wnd 4096
snd_una 20288002, snd_nxt 20288514, snd_max 20288514
snd_wll 3202722818, snd_wl2 20288002, snd_wnd 8192
REXMT=6 (t_rxtshift=0), KEEP=14400
t_rtt=1, t_srtt=16, t_rttvar=4, t_rxtcur=6

```

Figure 27.28 `trpt` output from example in Figure 25.28.

appears after processing the `PRU_CONNECT` request, the two trace records appear backward in the buffer. Also, when the process called `connect`, the connection state was `CLOSED`, and it changes to `SYN_SENT`. Nothing else changes from the first trace record to this one.

The third trace record, at time 954103, occurs 365 ms after the first. (`tcpdump` shows a 362.7 ms difference.) This is how the values in the column “actual delta (ms)” in Figure 25.28 were computed. The connection state changes from `SYN_SENT` to `ESTABLISHED` when the segment with a `SYN` and an `ACK` is received. The `RTT` estimators are updated because the segment being timed was acknowledged.

The fourth trace record is the third segment of the three-way handshake: the `ACK` of the other end’s `SYN`. Since this segment contains no data, it is not timed (`rtt` is 0).

After the ACK has been sent at time 954103, the `connect` system call returns to the process, which then calls `write` to send data. This generates TCP output, shown in trace record 5 at time 954153, 50 ms after the three-way handshake is complete. 512 bytes of data are sent, starting with sequence number 20288002. The retransmission timer is set to 3 seconds and the segment is timed.

This output is caused by an application `write`. Although we don't show any more trace records, the next four are from `PRU_SEND` requests. The first `PRU_SEND` request generates the output of the first 512-byte segment that we show, but the other three do not cause output, since the connection has just started and is in slow start. Four trace records are generated because the system used for this example uses a TCP send buffer of 4096 and a cluster size of 1024. Once the send buffer is full, the process is put to sleep.

27.11 Summary

This chapter has covered a wide range of TCP functions that we'll encounter in the following chapters.

TCP connections can be aborted by sending an RST or they can be closed down gracefully, by sending a FIN and waiting for the four-way exchange of segments to complete.

Eight variables are stored in each routing table entry, three of which are updated when a connection is closed and six of which can be used later when a new connection is established. This lets the kernel keep track of certain variables, such as the RTT estimators and the slow start threshold, between successive connections to the same destination. The system administrator can also set and lock some of these variables, such as the MTU, receive pipe size, and send pipe size, that affect TCP connections to that destination.

TCP is tolerant of received ICMP errors—none cause Net/3 to terminate an established connection. This handling of ICMP errors by Net/3 differs from earlier Berkeley releases.

Received TCP segments can arrive out of order and can contain duplicate data, and TCP must handle these anomalies. We saw that a reassembly queue is maintained for each connection, and this holds the out-of-order segments along with segments that arrive before they can be passed to the application.

Finally we looked at the type of information saved by the kernel when the `SO_DEBUG` socket option is enabled for a TCP socket. This trace information can be a useful diagnostic tool in addition to programs such as `tcpdump`.

Exercises

- 27.1 Why is the `errno` value 0 for the last row in Figure 27.1?
- 27.2 What is the maximum value that can be stored in `rmx_rtt`?
- 27.3 To save the route information in Figure 27.3 for a given host, we enter a route into the routing table by hand for this destination. We then run the FTP client to send data to this host, making certain we send enough data, as described with Figure 27.4. But after terminating the FTP client we look at the routing table, and all the values for this host are still 0. What's happening?

TCP Input

27.4 Introduction

FTP is a protocol for the transfer of files from a host to another host. It is a client-server protocol. The client is the user's computer, and the server is the computer that stores the files. The client sends a request to the server to get a file, and the server sends the file back to the client. The client then saves the file on its own disk.

The server is a computer that stores files. It is a host. The client is a computer that sends requests to the server. It is also a host. The client and server are both hosts.

The client sends a request to the server to get a file. The server sends the file back to the client. The client then saves the file on its own disk.

The client sends a request to the server to get a file. The server sends the file back to the client. The client then saves the file on its own disk.

The client sends a request to the server to get a file. The server sends the file back to the client. The client then saves the file on its own disk.

The client sends a request to the server to get a file. The server sends the file back to the client. The client then saves the file on its own disk.

The first part of the document discusses the background and the purpose of the study. It mentions that the study was conducted to determine the effectiveness of the program in improving the performance of the participants. The study was conducted over a period of six months, and the results were analyzed using statistical methods. The findings of the study indicate that the program was effective in improving the performance of the participants, and that the improvements were maintained over the long term. The study also identified some limitations and areas for future research.

2.2. Summary

The second part of the document provides a summary of the findings and conclusions of the study. It states that the program was effective in improving the performance of the participants, and that the improvements were maintained over the long term. The study also identified some limitations and areas for future research. The findings of the study have important implications for the design and implementation of similar programs. The study also identified some limitations and areas for future research. The findings of the study have important implications for the design and implementation of similar programs.

TCP Input

28.1 Introduction

TCP input processing is the largest piece of code that we examine in this text. The function `tcp_input` is about 1100 lines of code. The processing of incoming segments is not complicated, just long and detailed. Many implementations, including the one in Net/3, closely follow the input event processing steps in RFC 793, which spell out in detail how to respond to the various input segments, based on the current state of the connection.

The `tcp_input` function is called by `ipintr` (through the `pr_input` function in the protocol switch table) when a datagram is received with a protocol field of TCP. `tcp_input` executes at the software interrupt level.

The function is so long that we divide its discussion into two chapters. Figure 28.1 outlines the processing steps in `tcp_input`. This chapter discusses the steps through RST processing, and the next chapter starts with ACK processing.

The first few steps are typical: validate the input segment (checksum, length, etc.) and locate the PCB for this connection. Given the length of the remainder of the function, however, an attempt is made to bypass all this logic with an algorithm called *header prediction* (Section 28.4). This algorithm is based on the assumption that segments are not typically lost or reordered, hence for a given connection TCP can often guess what the next received segment will be. If the header prediction algorithm works, notice that the function returns. This is the fast path through `tcp_input`.

The slow path through the function ends up at the label `do_data`, which tests a few flags and calls `tcp_output` if a segment should be sent in response to the received segment.

```

void
tcp_input()
{
    checksum TCP header and data;
findpcb:
    locate PCB for segment;
    if (not found)
        goto dropwithreset;
    reset idle time to 0 and keepalive timer to 2 hours;
    process options if not LISTEN state;
    if (packet matched by header prediction) {
        completely process received segment;
        return;
    }
    switch (tp->t_state) {
    case TCPS_LISTEN:
        if SYN flag set, accept new connection request;
        goto trimthenstep6;
    case TCPS_SYN_SENT:
        if ACK of our SYN, connection completed;
        trimthenstep6;
        trim any data not within window;
        goto step6;
    }
    process RFC 1323 timestamp;
    check if some data bytes are within the receive window;
    trim data segment to fit within window;
    if (RST flag set) {
        process depending on state;
        goto drop;
    }
    if (ACK flag set) {
        /* Chapter 28 finishes here */
        if (SYN_RCVD state) /* Chapter 29 starts here */
            simultaneous open complete;
        if (duplicate ACK)
            fast recovery algorithm;
        update RTT estimators if segment timed;
        open congestion window;
        remove ACKed data from send buffer;
        change state if in FIN_WAIT_1, CLOSING, or LAST_ACK state;
    }
    step6:
        update window information;
        process URG flag;
}

```

```

do data:
    process data in segment, add to reassembly queue;
    if (FIN flag is set)
        process depending on state;
    if (SO_DEBUG socket option)
        tcp_trace(TA_INPUT);
    if (need output || ACK now)
        tcp_output();
    return;
dropafterack:
    tcp_output() to generate ACK;
    return;
dropwithreset:
    tcp_respond() to generate RST;
    return;
drop:
    if (SO_DEBUG socket option)
        tcp_trace(TA_DROP);
    return;
}

```

Figure 28.1 Summary of TCP input processing steps.

There are also three labels at the end of the function that are jumped to when errors occur: *dropafterack*, *dropwithreset*, and *drop*. The term *drop* means to drop the segment being processed, not drop the connection, but when an RST is sent by *dropwithreset* it normally causes the connection to be dropped.

The only other branching in the function occurs when a valid SYN is received in either the LISTEN or SYN_SENT states, at the switch following header prediction. When the code at *trimthenstep6* finishes, it jumps to *step6*, which continues the normal flow.

28.2 Preliminary Processing

Figure 28.2 shows the declarations and the initial processing of the received TCP segment.

Get IP and TCP headers in first mbuf

170-204 The argument *iphlen* is the length of the IP header, including possible IP options. If the length is greater than 20 bytes, options are present, and *ip_stripoptions* discards the options. TCP ignores all IP options other than a source route, which is saved specially by IP (Section 9.6) and fetched later by TCP in Figure 28.7. If the number of bytes in the first mbuf in the chain is less than the size of the combined IP/TCP header (40 bytes), *m_pullup* moves the first 40 bytes into the first mbuf.

```

170 void                                     tcp_input.c
171 tcp_input(m, iphlen)
172 struct mbuf *m;
173 int iphlen;
174 {
175     struct tcpiphdr *ti;
176     struct inpcb *inp;
177     caddr_t optp = NULL;
178     int optlen;
179     int len, tlen, off;
180     struct tcpcb *tp = 0;
181     int tiflags;
182     struct socket *so;
183     int todrop, acked, ourfinisacked, needoutput = 0;
184     short ostate;
185     struct in_addr laddr;
186     int dropsoccket = 0;
187     int iss = 0;
188     u_long tiwin, ts_val, ts_ecr;
189     int ts_present = 0;
190     tcpstat.tcps_rcvtotal++;
191     /*
192     * Get IP and TCP header together in first mbuf.
193     * Note: IP leaves IP header in first mbuf.
194     */
195     ti = mtod(m, struct tcpiphdr *);
196     if (iphlen > sizeof(struct ip))
197         ip_stripoptions(m, (struct mbuf *) 0);
198     if (m->m_len < sizeof(struct tcpiphdr)) {
199         if ((m = m_pullup(m, sizeof(struct tcpiphdr))) == 0) {
200             tcpstat.tcps_rcvshort++;
201             return;
202         }
203         ti = mtod(m, struct tcpiphdr *);
204     }

```

Figure 28.2 tcp_input function: declarations and preliminary processing.

The next piece of code, shown in Figure 28.3, verifies the TCP checksum and offset field.

Verify TCP checksum

205-217 tlen is the TCP length, the number of bytes following the IP header. Recall that IP has already subtracted the IP header length from ip_len. The variable len is then set to the length of the IP datagram, the number of bytes to be checksummed, including the pseudo-header. The fields in the pseudo-header are set, as required for the checksum calculation, as shown in Figure 23.19.

Verify TCP offset field

218-228 The TCP offset field, ti_off, is the number of 32-bit words in the TCP header, including any TCP options. It is multiplied by 4 (to become the byte offset of the first


```

205  /*
206  * Checksum extended TCP header and data.
207  */
208  tlen = ((struct ip *) ti)->ip_len;
209  len = sizeof(struct ip) + tlen;
210  ti->ti_next = ti->ti_prev = 0;
211  ti->ti_xl = 0;
212  ti->ti_len = (u_short) tlen;
213  HTONS(ti->ti_len);
214  if (ti->ti_sum = in_cksum(m, len)) {
215      tcpstat.tcps_rcvbadsum++;
216      goto drop;
217  }
218  /*
219  * Check that TCP offset makes sense,
220  * pull out TCP options and adjust length.    XXX
221  */
222  off = ti->ti_off << 2;
223  if (off < sizeof(struct tcphdr) || off > tlen) {
224      tcpstat.tcps_rcvbadoff++;
225      goto drop;
226  }
227  tlen -= off;
228  ti->ti_len = tlen;

```

Figure 28.3 tcp_input function: verify TCP checksum and offset field.

data byte in the TCP segment) and checked for sanity. It must be greater than or equal to the size of the standard TCP header (20) and less than or equal to the TCP length.

The byte offset of the first data byte is subtracted from the TCP length, leaving `tlen` with the number of bytes of data in the segment (possibly 0). This value is stored back into the TCP header, in the variable `ti_len`, and will be used throughout the function.

Figure 28.4 shows the next part of processing: handling of certain TCP options.

Get headers plus option into first mbuf

230-236 If the byte offset of the first data byte is greater than 20, TCP options are present. `m_pullup` is called, if necessary, to place the standard IP header, standard TCP header, and any TCP options in the first mbuf in the chain. Since the maximum size of these three pieces is 80 bytes (20 + 20 + 40), they all fit into the first packet header mbuf on the chain.

Since the only way `m_pullup` can fail here is when fewer than 20 plus `off` bytes are in the IP datagram, and since the TCP checksum has already been verified, we expect this call to `m_pullup` never to fail. Unfortunately the counter `tcps_rcvshort` is also shared by the call to `m_pullup` in Figure 28.2, so looking at the counter doesn't tell us which call failed. Nevertheless, Figure 24.5 shows that after receiving almost 9 million TCP segments, this counter is 0.

```

229     if (off > sizeof(struct tcphdr)) {
230         if (m->m_len < sizeof(struct ip) + off) {
231             if ((m = m_pullup(m, sizeof(struct ip) + off)) == 0) {
232                 tcpstat.tcps_rcvshort++;
233                 return;
234             }
235             ti = mtod(m, struct tcpiphdr *);
236         }
237         optlen = off - sizeof(struct tcphdr);
238         optp = mtod(m, caddr_t) + sizeof(struct tcpiphdr);
239         /*
240          * Do quick retrieval of timestamp options ("options
241          * prediction?"). If timestamp is the only option and it's
242          * formatted as recommended in RFC 1323 Appendix A, we
243          * quickly get the values now and not bother calling
244          * tcp_dooptions(), etc.
245          */
246         if ((optlen == TCPOLEN_TSTAMP_APPA ||
247             (optlen > TCPOLEN_TSTAMP_APPA &&
248              optp[TCPOLEN_TSTAMP_APPA] == TCPOPT_EOL)) &&
249             *(u_long *) optp == htonl(TCPOPT_TSTAMP_HDR) &&
250             (ti->ti_flags & TH_SYN) == 0) {
251             ts_present = 1;
252             ts_val = ntohl(*(u_long *) (optp + 4));
253             ts_ecr = ntohl(*(u_long *) (optp + 8));
254             optp = NULL; /* we've parsed the options */
255         }
256     }

```

Figure 28.4 tcp_input function: handle certain TCP options.

Process timestamp option quickly

237-255 optlen is the number of bytes of options, and optp is a pointer to the first option byte. If the following three conditions are all true, only the timestamp option is present and it is in the desired format:

1. (a) The TCP option length equals 12 (TCPOLEN_TSTAMP_APPA), or (b) the TCP option length is greater than 12 and optp[12] equals the end-of-option byte.
2. The first 4 bytes of options equals 0x0101080a (TCPOPT_TSTAMP_HDR, which we described in Section 26.6).
3. The SYN flag is not set (i.e., this segment is for an established connection, hence if a timestamp option is present, we know both sides have agreed to use the option).

If all three conditions are true, ts_present is set to 1; the two timestamp values are fetched and stored in ts_val and ts_ecr; and optp is set to null, since all the options have been parsed. The benefit in recognizing the timestamp option this way is to avoid calling the general option processing function tcp_dooptions later in the code. The general option processing function is OK for the other options that appear only with the

SYN segment that creates a connection (the MSS and window scale options), but when the timestamp option is being used, it will appear with almost every segment on an established connection, so the faster it can be recognized, the better.

The next piece of code, shown in Figure 28.5, locates the Internet PCB for the segment.

```

257     tiflags = ti->ti_flags;                                     tcp_input.c
258     /*
259     * Convert TCP protocol specific fields to host format.
260     */
261     NTOHL(ti->ti_seq);
262     NTOHL(ti->ti_ack);
263     NTOHS(ti->ti_win);
264     NTOHS(ti->ti_urp);
265     /*
266     * Locate pcb for segment.
267     */
268     findpcb:
269     inp = tcp_last_inpcb;
270     if (inp->inp_lport != ti->ti_dport ||
271         inp->inp_fport != ti->ti_sport ||
272         inp->inp_faddr.s_addr != ti->ti_src.s_addr ||
273         inp->inp_laddr.s_addr != ti->ti_dst.s_addr) {
274         inp = in_pcblookup(stcb, ti->ti_src, ti->ti_sport,
275                         ti->ti_dst, ti->ti_dport, INPLOOKUP_WILDCARD);
276         if (inp)
277             tcp_last_inpcb = inp;
278         ++tcpstat.tcps_pcbcachehits;
279     }

```

Figure 28.5 tcp_input function: locate Internet PCB for segment.

Save input flags and convert fields to host byte order

257-264 The received flags (SYN, FIN, etc.) are saved in the local variable `tiflags`, since they are referenced throughout the code. Two 16-bit values and the two 32-bit values in the TCP header are converted from network byte order to host byte order. The two 16-bit port numbers are left in network byte order, since the port numbers in the Internet PCB are in that order.

Locate Internet PCB

265-279 TCP maintains a one-behind cache (`tcp_last_inpcb`) containing the address of the PCB for the last received TCP segment. This is the same technique used by UDP. The comparison of the four elements in the socket pair is in the same order as done by `udp_input`. If the cache entry does not match, `in_pcblookup` is called, and the cache is set to the new PCB entry.

TCP does not have the same problem that we encountered with UDP: wildcard entries in the cache causing a high miss rate. The only time a TCP socket has a wildcard entry is for a server listening for connection requests. Once a connection is made, all

four entries in the socket pair contain nonwildcard values. In Figure 24.5 we see a cache hit rate of almost 80%.

Figure 28.6 shows the next piece of code.

```

280  /*
281  * If the state is CLOSED (i.e., TCB does not exist) then
282  * all data in the incoming segment is discarded.
283  * If the TCB exists but is in CLOSED state, it is embryonic,
284  * but should either do a listen or a connect soon.
285  */
286  if (inp == 0)
287      goto dropwithreset;
288  tp = intotcpb(inp);
289  if (tp == 0)
290      goto dropwithreset;
291  if (tp->t_state == TCPS_CLOSED)
292      goto drop;
293  /* Unscale the window into a 32-bit value, */
294  if ((tiflags & TH_SYN) == 0)
295      tiwin = ti->ti_win << tp->snd_scale;
296  else
297      tiwin = ti->ti_win;

```

Figure 28.6 tcp_input function: check if segment should be dropped.

Drop segment and generate RST

280-287 If the PCB was not found, the input segment is dropped and an RST is sent as a reply. This is how TCP handles SYN's that arrive for a server that doesn't exist, for example. Recall that UDP sends an ICMP port unreachable in this case.

288-289 If the PCB exists but a corresponding TCP control block does not exist, the socket is probably being closed (`tcp_close` releases the TCP control block first, and then releases the PCB), so the input segment is dropped and an RST is sent as a reply.

Silently drop segment

291-292 If the TCP control block exists, but the connection state is CLOSED, the socket has been created and a local address and local port may have been assigned, but neither `connect` nor `listen` has been called. The segment is dropped but nothing is sent as a reply. This scenario can happen if a client catches a server between the server's call to `bind` and `listen`. By silently dropping the segment and not replying with an RST, the client's connection request should time out, causing the client to retransmit the SYN.

Unscale advertised window

293-297 If window scaling is to take place for this connection, both ends must specify their send scale factor using the window scale option when the connection is established. If the segment contains a SYN, the window scale factor has not been established yet, so `tiwin` is copied from the value in the TCP header. Otherwise the 16-bit value in the header is left shifted by the send scale factor into a 32-bit value.

The next piece of code, shown in Figure 28.7, does some preliminary processing if the socket debug option is enabled or if the socket is listening for incoming connection requests.

```

298     so = inp->inp_socket;
299     if (so->so_options & (SO_DEBUG | SO_ACCEPTCONN)) {
300         if (so->so_options & SO_DEBUG) {
301             ostate = tp->t_state;
302             tcp_saveti = *ti;
303         }
304         if (so->so_options & SO_ACCEPTCONN) {
305             so = sonewconn(so, 0);
306             if (so == 0)
307                 goto drop;
308             /*
309              * This is ugly, but ....
310              *
311              * Mark socket as temporary until we're
312              * committed to keeping it. The code at
313              * 'drop' and 'dropwithreset' check the
314              * flag dropsocket to see if the temporary
315              * socket created here should be discarded.
316              * We mark the socket as discardable until
317              * we're committed to it below in TCPS_LISTEN.
318              */
319             dropsocket++;
320             inp = (struct inpcb *) so->so_pcb;
321             inp->inp_laddr = ti->ti_daddr;
322             inp->inp_lport = ti->ti_dport;
323             #if BSD>=43
324                 inp->inp_options = ip_srcroute();
325             #endif
326             tp = intotcpcb(inp);
327             tp->t_state = TCPS_LISTEN;
328             /* Compute proper scaling value from buffer space */
329             while (tp->request_r_scale < TCP_MAX_WINSHIFT &&
330                 TCP_MAXWIN << tp->request_r_scale < so->so_rcv.sb_hiwat)
331                 tp->request_r_scale++;
332         }
333     }

```

Figure 28.7 tcp_input function: handle debug option and listening sockets.

Save connection state and IP/TCP headers if socket debug option enabled

100-103 If the `SO_DEBUG` socket option is enabled the current connection state is saved (`ostate`) as well as the IP and TCP headers (`tcp_saveti`). These become arguments to `tcp_trace` when it is called at the end of the function (Figure 29.26).

Create new socket if segment arrives for listening socket

104-119 When a segment arrives for a listening socket (`SO_ACCEPTCONN` is enabled by `listen`), a new socket is created by `sonewconn`. This issues the protocol's

PRU_ATTACH request (Figure 30.2), which allocates an Internet PCB and a TCP control block. But more processing is needed before TCP commits to accept the connection request (such as the fundamental question of whether the segment contains a SYN or not), so the flag `dropsocket` is set, to cause the code at the labels `drop` and `dropwithreset` to discard the new socket if an error is encountered. If the received segment is OK, `dropsocket` is set back to 0 in Figure 28.17.

320-326 `inp` and `tp` point to the new socket that has been created. The local address and local port are copied from the destination address and destination port of the IP and TCP headers. If the input datagram contained a source route, it was saved by `save_rte`. TCP calls `ip_srcroute` to fetch that source route, saving a pointer to the mbuf containing the source route option in `inp_options`. This option is passed to `ip_output` by `tcp_output`, and the reverse route is used for datagrams sent on this connection.

327 The state of the new socket is set to LISTEN. If the received segment contains a SYN, the code in Figure 28.16 completes the connection request.

Compute window scale factor

328-331 The window scale factor that will be requested is calculated from the size of the receive buffer. 65535 (`TCP_MAXWIN`) is left shifted until the result exceeds the size of the receive buffer, or until the maximum window scale factor is encountered (14, `TCP_MAX_WNSHIFT`). Notice that the requested window scale factor is chosen based on the size of the listening socket's receive buffer. This means the process must set the `SO_RCVBUF` socket option before listening for incoming connection requests or it inherits the default value in `tcp_recvspace`.

The maximum scale factor is 14, and 65535×2^{14} is 1,073,725,440. This is far greater than the maximum size of the receive buffer (262,144 in Net/3), so the loop should always terminate with a scale factor much less than 14. See Exercises 28.1 and 28.2.

Figure 28.8 shows the next part of TCP input processing.

```

334 /*
335  * Segment received on connection.
336  * Reset idle time and keepalive timer.
337  */
338 tp->t_idle = 0;
339 tp->t_timer[TCP_T_KEEPC] = tcp_keeptime;
340 /*
341  * Process options if not in LISTEN state,
342  * else do it below (after getting remote address).
343  */
344 if (optp && tp->t_state != TCPS_LISTEN)
345     tcp_dooptions(tp, optp, optlen, ti,
346                 &ts_present, &ts_val, &ts_ecr);

```

Figure 28.8 `tcp_input` function: reset idle time and keepalive timer, process options.

Reset idle time and keepalive timer

134-139 `t_idle` is set to 0 since a segment has been received on the connection. The keep-alive timer is also reset to 2 hours.

Process TCP options if not in LISTEN state

340-346 If options are present in the TCP header, and if the connection state is not LISTEN, `tcp_dooptions` processes the options. Recall that if only a timestamp option appears for an established connection, and that option is in the format recommended by Appendix A of RFC 1323, it was already processed in Figure 28.4 and `opt` was set to a null pointer. If the socket is in the LISTEN state, `tcp_dooptions` is called in Figure 28.17 after the peer's address has been recorded in the PCB, because processing the MSS option requires knowledge of the route that will be used to this peer.

28.3 tcp_dooptions Function

This function processes the five TCP options supported by Net/3 (Section 26.4): the EOL, NOP, MSS, window scale, and timestamp options. Figure 28.9 shows the first part of this function.

```

1213 void                                     tcp_input.c
1214 tcp_dooptions(tp, cp, cnt, ti, ts_present, ts_val, ts_ecr)
1215 struct tcpcb *tp;
1216 u_char *cp;
1217 int cnt;
1218 struct tcphdr *ti;
1219 int *ts_present;
1220 u_long *ts_val, *ts_ecr;
1221 {
1222     u_short mss;
1223     int opt, optlen;
1224
1225     for (; cnt > 0; cnt -= optlen, cp += optlen) {
1226         opt = cp[0];
1227         if (opt == TCPOPT_EOL)
1228             break;
1229         if (opt == TCPOPT_NOP)
1230             optlen = 1;
1231         else {
1232             optlen = cp[1];
1233             if (optlen <= 0)
1234                 break;
1235         }
1236         switch (opt) {
1237             default:
1238                 continue;

```

Figure 28.9 `tcp_dooptions` function: handle EOL and NOP options.

Fetch option type and length

1225-1229 The options are scanned and an EOL (end-of-options) terminates the processing, causing the function to return. The length of a NOP is set to 1, since this option is not followed by a length byte (Figure 26.16). The NOP will be ignored via the default in the switch statement.

1230-1234 All other options have a length byte that is stored in `optlen`.

Any new options that are not understood by this implementation of TCP are also ignored. This occurs because:

1. Any new options defined in the future will have an option length (NOP and EOL are the only two without a length), and the for loop skips `optlen` bytes each time around the loop.
2. The default in the switch statement ignores unknown options.

The final part of `tcp_dooptions`, shown in Figure 28.10, handles the MSS, window scale, and timestamp options.

MSS option

1218-1246 If the length is not 4 (`TCPOLEN_MAXSEG`), or the segment does not have the SYN flag set, the option is ignored. Otherwise the 2 MSS bytes are copied into a local variable, converted to host byte order, and processed by `tcp_mss`. This has the side effect of setting the variable `t_maxseg` in the control block, the maximum number of bytes that can be sent in a segment to the other end.

Window scale option

1247-1254 If the length is not 3 (`TCPOLEN_WINDOW`), or the segment does not have the SYN flag set, the option is ignored. Net/3 remembers that it received a window scale request, and the scale factor is saved in `requested_s_scale`. Since only 1 byte is referenced by `cp[2]`, there can't be alignment problems. When the ESTABLISHED state is entered, if both ends requested window scaling, it is enabled.

Timestamp option

1255-1293 If the length is not 10 (`TCPOLEN_TIMESTAMP`), the segment is ignored. Otherwise the flag pointed to by `ts_present` is set to 1, and the two timestamps are saved in the variables pointed to by `ts_val` and `ts_ecr`. If the received segment contains the SYN flag, Net/3 remembers that a timestamp request was received. `ts_recent` is set to the received timestamp and `ts_recent_age` is set to `tcp_now`, the counter of the number of 500-ms clock ticks since the system was initialized.

28.4 Header Prediction

We now continue with the code in `tcp_input`, from where we left off in Figure 28.8.

Header prediction was put into the 4.3BSD Reno release by Van Jacobson. The only description of the algorithm, other than the source code we're about to examine, is in [Jacobson 1990b], which is a copy of three slides showing the code.


```

1238     case TCPOPT_MAXSEG:                                     tcp_input.c
1239         if (optlen != TCPOLEN_MAXSEG)
1240             continue;
1241         if (!(ti->ti_flags & TH_SYN))
1242             continue;
1243         bcopy((char *) cp + 2, (char *) &mss, sizeof(mss));
1244         NTOHS(mss);
1245         (void) tcp_mss(tp, mss); /* sets t_maxseg */
1246         break;
1247     case TCPOPT_WINDOW:
1248         if (optlen != TCPOLEN_WINDOW)
1249             continue;
1250         if (!(ti->ti_flags & TH_SYN))
1251             continue;
1252         tp->t_flags |= TF_RCVD_SCALE;
1253         tp->requested_s_scale = min(cp[2], TCP_MAX_WINSHIFT);
1254         break;
1255     case TCPOPT_TIMESTAMP:
1256         if (optlen != TCPOLEN_TIMESTAMP)
1257             continue;
1258         *ts_present = 1;
1259         bcopy((char *) cp + 2, (char *) ts_val, sizeof(*ts_val));
1260         NTOHL(*ts_val);
1261         bcopy((char *) cp + 6, (char *) ts_ecr, sizeof(*ts_ecr));
1262         NTOHL(*ts_ecr);
1263         /*
1264          * A timestamp received in a SYN makes
1265          * it ok to send timestamp requests and replies.
1266          */
1267         if (ti->ti_flags & TH_SYN) {
1268             tp->t_flags |= TF_RCVD_TSTMP;
1269             tp->ts_recent = *ts_val;
1270             tp->ts_recent_age = tcp_now;
1271         }
1272         break;
1273     }
1274 }
1275 }

```

Figure 28.10 tcp_dooptions function: process MSS, window scale, and timestamp options.

Header prediction helps unidirectional data transfer by handling the two common cases.

1. If TCP is sending data, the next expected segment for this connection is an ACK for outstanding data.
2. If TCP is receiving data, the next expected segment for this connection is the next in-sequence data segment.

In both cases a small set of tests determines if the next expected segment has been received, and if so, it is handled in-line, faster than the general processing that follows later in this chapter and the next.

[Partridge 1993] shows an even faster version of TCP header prediction from a research implementation developed by Van Jacobson.

Figure 28.11 shows the first part of header prediction.

```

347  /*
348  * Header prediction: check for the two common cases
349  * of a uni-directional data xfer. If the packet has
350  * no control flags, is in-sequence, the window didn't
351  * change and we're not retransmitting, it's a
352  * candidate. If the length is zero and the ack moved
353  * forward, we're the sender side of the xfer. Just
354  * free the data acked & wake any higher-level process
355  * that was blocked waiting for space. If the length
356  * is non-zero and the ack didn't move, we're the
357  * receiver side. If we're getting packets in order
358  * (the reassembly queue is empty), add the data to
359  * the socket buffer and note that we need a delayed ack.
360  */
361  if (tp->t_state == TCPS_ESTABLISHED &&
362      (tiflags & (TH_SYN | TH_FIN | TH_RST | TH_URG | TH_ACK)) == TH_ACK &&
363      ((ts_present || TSTMP_GEQ(ts_val, tp->ts_recent)) &&
364       ti->ti_seq == tp->rcv_nxt &&
365       tiwin && tiwin == tp->snd_wnd &&
366       tp->snd_nxt == tp->snd_max) {
367
368      /*
369      * If last ACK falls within this segment's sequence numbers,
370      * record the timestamp.
371      */
372      if (ts_present && SEQ_LEQ(ti->ti_seq, tp->last_ack_sent) &&
373          SEQ_LT(tp->last_ack_sent, ti->ti_seq + ti->ti_len)) {
374          tp->ts_recent_age = tcp_now;
375          tp->ts_recent = ts_val;
376      }
377  }

```

Figure 28.11 tcp_input function: header prediction, first part.

Check if segment is the next expected

347-366 The following six conditions must *all* be true for the segment to be the next expected data segment or the next expected ACK:

1. The connection state must be ESTABLISHED.
2. The following four control flags must not be on: SYN, FIN, RST, or URG. The ACK flag must be on. In other words, of the six TCP control flags, the ACK flag must be set, the four just listed must be cleared, and it doesn't matter whether

PSH is set or cleared. (Normally in the ESTABLISHED state the ACK flag is always on unless the RST flag is on.)

3. If the segment contains a timestamp option, the timestamp value from the other end (`ts_val`) must be greater than or equal to the previous timestamp received for this connection (`ts_recent`). This is basically the PAWS test, which we describe in detail in Section 28.7. If `ts_val` is less than `ts_recent`, this segment is out of order because it was sent before the most previous segment received on this connection. Since the other end always sends its timestamp clock (the global variable `tcp_now` in Net/3) as its timestamp value, the received timestamps of in-order segments always form a monotonic increasing sequence.

The timestamp need not increase with every in-order segment. Indeed, on a Net/3 system that increments the timestamp clock (`tcp_now`) every 500 ms, multiple segments are often sent on a connection before that clock is incremented. Think of the timestamp and sequence number as forming a 64-bit value, with the sequence number in the low-order 32 bits and the timestamp in the high-order 32 bits. This 64-bit value always increases by at least 1 for every in-order segment (taking into account the modulo arithmetic).

4. The starting sequence number of the segment (`ti_seq`) must equal the next expected receive sequence number (`rcv_nxt`). If this test is false, then the received segment is either a retransmission or a segment beyond the one expected.
5. The window advertised by the segment (`tiwin`) must be nonzero, and must equal the current send window (`snd_wnd`). This means the window has not changed.
6. The next sequence number to send (`snd_nxt`) must equal the highest sequence number sent (`snd_max`). This means the last segment sent by TCP was not a retransmission.

Update `ts_recent` from received timestamp

367-375 If a timestamp option is present and if its value passes the test described with Figure 26.18, the received timestamp (`ts_val`) is saved in `ts_recent`. Also, the current time (`tcp_now`) is recorded in `ts_recent_age`.

Recall our discussion with Figure 26.18 on how this test for a valid timestamp is flawed, and the correct test presented in Figure 26.20. In this header prediction code the `TSTMP_GEQ` test in Figure 26.20 is redundant, since it was already done as step 3 of the `if` test at the beginning of Figure 28.11.

The next part of the header prediction code, shown in Figure 28.12, is for the sender of unidirectional data: process an ACK for outstanding data.

Test for pure ACK

376-379 If the following four conditions are all true, this segment is a pure ACK.

```

376     if (ti->ti_len == 0) {
377         if (SEQ_GT(ti->ti_ack, tp->snd_una) &&
378             SEQ_LEQ(ti->ti_ack, tp->snd_max) &&
379             tp->snd_cwnd >= tp->snd_wnd) {
380             /*
381              * this is a pure ack for outstanding data.
382              */
383             ++tcpstat.tcps_predack;
384             if (ts_present)
385                 tcp_xmit_timer(tp, tcp_now - ts_eck + 1);
386             else if (tp->t_rtt &&
387                 SEQ_GT(ti->ti_ack, tp->t_rtseq))
388                 tcp_xmit_timer(tp, tp->t_rtt);
389             acked = ti->ti_ack - tp->snd_una;
390             tcpstat.tcps_rcvackpack++;
391             tcpstat.tcps_rcvackbyte += acked;
392             sbdrop(&so->so_snd, acked);
393             tp->snd_una = ti->ti_ack;
394             m_freem(m);
395             /*
396              * If all outstanding data is acked, stop
397              * retransmit timer, otherwise restart timer
398              * using current (possibly backed-off) value.
399              * If process is waiting for space,
400              * wakeup/selwakeup/signal.  If data
401              * is ready to send, let tcp_output
402              * decide between more output or persist.
403              */
404             if (tp->snd_una == tp->snd_max)
405                 tp->t_timer[TCPT_REXMT] = 0;
406             else if (tp->t_timer[TCPT_PERSIST] == 0)
407                 tp->t_timer[TCPT_REXMT] = tp->t_rxtcur;
408             if (so->so_snd.sb_flags & SB_NOTIFY)
409                 sowwakeup(so);
410             if (so->so_snd.sb_cc)
411                 (void) tcp_output(tp);
412             return;
413         }

```

Figure 28.12 tcp_input function: header prediction, sender processing.

1. The segment contains no data (`ti_len` is 0).
2. The acknowledgment field in the segment (`ti_ack`) is greater than the largest unacknowledged sequence number (`snd_una`). Since this test is "greater than" and not "greater than or equal to," it is true only if some positive amount of data is acknowledged by the ACK.
3. The acknowledgment field in the segment (`ti_ack`) is less than or equal to the maximum sequence number sent (`snd_max`).

4. The congestion window (`snd_cwnd`) is greater than or equal to the current send window (`snd_wnd`). This test is true only if the window is fully open, that is, the connection is not in the middle of slow start or congestion avoidance.

Update RTT estimators

384-388 If the segment contains a timestamp option, or if a segment was being timed and the acknowledgment field is greater than the starting sequence number being timed, `tcp_xmit_timer` updates the RTT estimators.

Delete acknowledged bytes from send buffer

382-394 `acked` is the number of bytes acknowledged by the segment. `sbdrop` deletes those bytes from the send buffer. The largest unacknowledged sequence number (`snd_una`) is set to the acknowledgment field and the received mbuf chain is released. (Since the length is 0, there should be just a single mbuf containing the headers.)

Stop retransmit timer

395-407 If the received segment acknowledges all outstanding data (`snd_una` equals `snd_max`), the retransmission timer is turned off. Otherwise, if the persist timer is off, the retransmit timer is restarted using `t_rxtout` as the timeout.

Recall that when `tcp_output` sends a segment, it sets the retransmit timer only if the timer is not currently enabled. If two segments are sent one right after the other, the timer is set when the first is sent, but not touched when the second is sent. But if an ACK is received only for the first segment, the retransmit timer must be restarted, in case the second was lost.

Awaken waiting processes

408-409 If a process must be awakened when the send buffer is modified, `sowakeup` is called. From Figure 16.5, `SB_NOTIFY` is true if a process is waiting for space in the buffer, if a process is `selecting` on the buffer, or if a process wants the `SIGIO` signal for this socket.

Generate more output

410-411 If there is data in the send buffer, `tcp_output` is called because the sender's window has moved to the right. `snd_una` was just incremented and `snd_wnd` did not change, so in Figure 24.17 the entire window has shifted to the right.

The next part of header prediction, shown in Figure 28.13, is the receiver processing when the segment is the next in-sequence data segment.

Test for next in-sequence data segment

414-416 If the following four conditions are all true, this segment is the next expected data segment for the connection, and there is room in the socket buffer for the data.

1. The amount of data in the segment (`ti_len`) is greater than 0. This is the `else` portion of the `if` at the beginning of Figure 28.12.
2. The acknowledgment field (`ti_ack`) equals the largest unacknowledged sequence number. This means no data is acknowledged by this segment.

```

414     } else if (ti->ti_ack == tp->end_una &&                                     tcp_input.c
415                tp->seg_next == (struct tcpiphdr *) tp &&
416                ti->ti_len <= sbpace(&so->so_rcv)) {
417         /*
418          * this is a pure, in-sequence data packet;
419          * with nothing on the reassembly queue and
420          * we have enough buffer space to take it.
421          */
422         ++tcpstat.tcps_preddat;
423         tp->rcv_next += ti->ti_len;
424         tcpstat.tcps_rcvpack++;
425         tcpstat.tcps_rcvbyte += ti->ti_len;
426         /*
427          * Drop TCP, IP headers and TCP options then add data
428          * to socket buffer.
429          */
430         m->m_data += sizeof(struct tcpiphdr) + off - sizeof(struct tcphdr);
431         m->m_len -= sizeof(struct tcpiphdr) + off - sizeof(struct tcphdr);
432         sbappend(&so->so_rcv, m);
433         sorwakeup(so);
434         tp->t_flags |= TF_DELACK;
435         return;
436     }
437 }

```

Figure 28.13 tcp_input function: header prediction, receiver processing.

3. The reassembly list of out-of-order segments for the connection is empty (seg_next equals tp).
4. There is room in the receive buffer for the data in the segment.

Complete processing of received data

423-435 The next expected receive sequence number (rcv_next) is incremented by the number of bytes of data. The IP header, TCP header, and any TCP options are dropped from the mbuf, and the mbuf chain is appended to the socket's receive buffer. The receiving process is awakened by sorwakeup. Notice that this code avoids calling the TCP_REASS macro, since the tests performed by that macro have already been performed by the header prediction tests. The delayed-ACK flag is set and the input processing is complete.

Statistics

How useful is header prediction? A few simple unidirectional transfers were run across a LAN (between bsd1 and svr4, in both directions) and across a WAN (between vangogh.cs.berkeley.edu and ftp.uu.net in both directions). The netstat output (Figure 24.5) shows the two header prediction counters.

On the LAN, with no packet loss but a few duplicate ACKs, header prediction worked between 97 and 100% of the time. Across the WAN, however, the header prediction percentages dropped slightly to between 83 and 99%.

Realize that header prediction works on a per-connection basis, regardless how much additional TCP traffic is being received by the host, while the PCB cache works on a per-host basis. Even though lots of TCP traffic can cause PCB cache misses, if packets are not lost on a given connection, header prediction still works on that connection.

28.5 TCP Input: Slow Path Processing

We continue with the code that's executed if header prediction fails, the slow path through `tcp_input`. Figure 28.14 shows the next piece of code, which prepares the received segment for input processing.

```

438      /* tcp_input.c
439      * Drop TCP, IP headers and TCP options.
440      */
441      m->m_data += sizeof(struct tcpiphdr) + off - sizeof(struct tcphdr);
442      m->m_len -= sizeof(struct tcpiphdr) + off - sizeof(struct tcphdr);
443
444      /*
445      * Calculate amount of space in receive window,
446      * and then do TCP input processing.
447      * Receive window is amount of space in rcv queue,
448      * but not less than advertised window.
449      */
450      {
451          int win;
452
453          win = sbspace(&so->so_rcv);
454          if (win < 0)
455              win = 0;
456          tp->rcv_wnd = max(win, (int) (tp->rcv_adv - tp->rcv_nxt));
457      }
tcp_input.c

```

Figure 28.14 `tcp_input` function: drop IP and TCP headers.

Drop IP and TCP headers, including TCP options

438-442 The data pointer and length of the first mbuf in the chain are updated to skip over the IP header, TCP header, and any TCP options. Since `off` is the number of bytes in the TCP header, including options, the size of the normal TCP header (20) must be subtracted from the expression.

Calculate receive window

443-455 `win` is set to the number of bytes available in the socket's receive buffer. `rcv_adv` minus `rcv_nxt` is the current advertised window. The receive window is the maximum of these two values. The `max` is taken to ensure that the value is not less than the currently advertised window. Also, if the process has taken data out of the socket

receive buffer since the window was last advertised, `win` could exceed the advertised window, so TCP accepts up to `win` bytes of data (even though the other end should not be sending more than the advertised window).

This value is calculated now, since the code later in this function must determine how much of the received data (if any) fits within the advertised window. Any received data outside the advertised window is dropped: data to the left of the window is duplicate data that has already been received and acknowledged, and data to the right should not be sent by the other end.

28.6 Completion of Passive Open or Active Open

If the state is `LISTEN` or `SYN_SENT`, the code shown in this section is executed. The expected segment in these two states is a `SYN`, and we'll see that any other received segment is dropped.

Completion of Passive Open

Figure 28.15 shows the processing when the connection is in the `LISTEN` state. In this code the variables `tp` and `inp` refer to the *new* socket that was created in Figure 28.7, not the server's listening socket.

```

456  switch (tp->t_state) {
457      /*
458       * If the state is LISTEN then ignore segment if it contains an RST.
459       * If the segment contains an ACK then it is bad and send an RST.
460       * If it does not contain a SYN then it is not interesting; drop it.
461       * Don't bother responding if the destination was a broadcast.
462       * Otherwise initialize tp->rcv_nxt, and tp->irs, select an initial
463       * tp->iss, and send a segment:
464       *     <SEQ=ISS>-<ACK=RCV_NXT>-<CTL=SYN,ACK>
465       * Also initialize tp->snd_nxt to tp->iss+1 and tp->snd_una to tp->iss
466       * Fill in remote peer address fields if not previously specified.
467       * Enter SYN_RECEIVED state, and process any other fields of this
468       * segment in this state.
469       */
470  case TCPS_LISTEN: {
471      struct mbuf *am;
472      struct sockaddr_in *sin;
473
474      if (tiflags & TH_RST)
475          goto drop;
476      if (tiflags & TH_ACK)
477          goto dropwithreset;
478      if ((tiflags & TH_SYN) == 0)
479          goto drop;

```

Figure 28.15 `tcp_input` function: check if SYN received for listening socket.

Drop if RST, ACK, or no SYN

473-478 If the received segment contains the RST flag, it is dropped. If it contains an ACK, it is dropped and an RST is sent as the reply. (The initial SYN to open a connection is one of the few segments that does not contain an ACK.) If the SYN flag is not set, the segment is dropped. The remaining code for this case handles the reception of a SYN for a connection in the LISTEN state. The new state will be SYN_RCVD.

Figure 28.16 shows the next piece of code for this case.

```

479          /*                                     tcp_input.c
480          * RFC1122 4.2.3.10, p. 104; discard bcst/mcast SYN
481          * in_broadcast() should never return true on a received
482          * packet with M_BCAST not set.
483          */
484          if (m->m_flags & (M_BCAST | M_MCAST) ||
485              IN_MULTICAST(ti->ti_dst.e_addr))
486              goto drop;
487
488          am = m_get(M_DONTWAIT, MT_SONAME); /* XXX */
489          if (am == NULL)
490              goto drop;
491          am->m_len = sizeof(struct sockaddr_in);
492          sin = mtod(am, struct sockaddr_in *);
493          sin->sin_family = AF_INET;
494          sin->sin_len = sizeof(*sin);
495          sin->sin_addr = ti->ti_src;
496          sin->sin_port = ti->ti_sport;
497          bzero((caddr_t) sin->sin_zero, sizeof(sin->sin_zero));
498
499          laddr = inp->inp_laddr;
500          if (inp->inp_laddr.e_addr == INADDR_ANY)
501              inp->inp_laddr = ti->ti_dst;
502          if (in_pcbconnect(inp, am)) {
503              inp->inp_laddr = laddr;
504              (void) m_free(am);
505              goto drop;
506          }
507          (void) m_free(am);

```

Figure 28.16 tcp_input function: process SYN for listening socket.

Drop if broadcast or multicast

479-486 If the packet was sent to a broadcast or multicast address, it is dropped. TCP is defined only for unicast applications. Recall that the M_BCAST and M_MCAST flags were set by ether_input, based on the destination hardware address of the frame. The IN_MULTICAST macro tests whether the IP address is a class D address.

The comment reference to in_broadcast is because the Net/1 code (which did not support multicasting) called that function here, to check whether the destination IP address was a broadcast address. The setting of the M_BCAST and M_MCAST flags by ether_input, based on the destination hardware address, was introduced with Net/2.

This Net/3 code tests only whether the destination hardware address is a broadcast address, and does not call `in_broadcast` to test whether the destination IP address is a broadcast address, on the assumption that a packet should never be received with a destination IP address that is a broadcast address unless the packet was sent to the hardware broadcast address. This assumption is made to avoid calling `in_broadcast`. Nevertheless, if a Net/3 system receives a SYN destined for a broadcast IP address but a unicast hardware address, that segment will be processed by the code in Figure 28.16.

The destination address argument to `IN_MULTICAST` needs to be converted to host byte order.

Get mbuf for client's IP address and port

487-496 An mbuf is allocated to hold a `sockaddr_in` structure, and the structure is filled in with the client's IP address and port number. The IP address is copied from the source address in the IP header and the port number is copied from the source port number in the TCP header. This structure is used shortly to connect the server's PCB to the client, and then the mbuf is released.

The XXX comment is probably because of the cost associated with obtaining an mbuf just for the call to `in_pcbconnect` that follows. But this is the slow processing path for TCP input. Figure 24.5 shows that less than 2% of all received segments execute this code.

Set local address in PCB

497-499 `laddr` is the local address bound to the socket. If the server bound the wildcard address to the socket (the normal scenario), the destination address from the IP header becomes the local address in the PCB. Note that the destination address from the IP header is used, regardless of which local interface the datagram was received on.

Notice that `laddr` cannot be the wildcard address, because in Figure 28.7 it is explicitly set to the destination IP address from the received datagram.

Connect PCB to peer

500-505 `in_pcbconnect` connects the server's PCB to the client. This fills in the foreign address and foreign process in the PCB. The mbuf is then released.

The next piece of code, shown in Figure 28.17 completes the processing for this case.

Allocate and initialize IP and TCP header template

506-511 A template of the IP and TCP headers is created by `tcp_template`. The call to `sonewconn` in Figure 28.7 allocated the PCB and TCP control block for the new connection, but not the header template.

Process any TCP options

512-514 If TCP options are present, they are processed by `tcp_dooptions`. The call to this function in Figure 28.8 was done only if the connection was not in the LISTEN state. This function is called now for a listening socket, after the foreign address is set in the PCB, since the foreign address is used by the `tcp_mss` function: to get a route to the peer, and to check if the peer is "local" or "foreign" (with regard to the peer's network ID and subnet ID, used to select the MSS).

```

506         tp->t_template = tcp_template(tp);                                     tcp_input.c
507         if (tp->t_template == 0) {
508             tp = tcp_drop(tp, ENOBUFS);
509             dropsocket = 0; /* socket is already gone */
510             goto drop;
511         }
512         if (opttp)
513             tcp_dooptions(tp, opttp, optlen, ti,
514                           &ts_present, &ts_val, &ts_ecr);
515         if (iss)
516             tp->iss = iss;
517         else
518             tp->iss = tcp_iss;
519         tcp_iss += TCP_ISSINCR / 2;
520         tp->irs = ti->ti_seq;
521         tcp_sendseqinit(tp);
522         tcp_rcvseqinit(tp);
523         tp->t_flags |= TF_ACKNOW;
524         tp->t_state = TCFS_SYN_RECEIVED;
525         tp->t_timer[TCPT_KEEP] = TCPTV_KEEP_INIT;
526         dropsocket = 0; /* committed to socket */
527         tcpstat.tcps_accepts++;
528         goto trinthestep6;
529     }

```

Figure 28.17 `tcp_input` function: complete processing of SYN received in LISTEN state.

Initialize ISS

515-519 The initial send sequence number is normally copied from the global `tcp_iss`, which is then incremented by 64,000 (`TCP_ISSINCR` divided by 2). If the local variable `iss` is nonzero, however, its value is used instead of `tcp_iss` to initialize the send sequence number for the connection.

The local `iss` variable is used for the following scenario.

- A server is started on port 27 on the host with an IP address of 128.1.2.3.
- A client on host 192.3.4.5 establishes a connection with this server. The client's ephemeral port is 3000. The socket pair on the server is {128.1.2.3, 27, 192.3.4.5, 3000}.
- The server actively closes the connection, putting this socket pair into the `TIME_WAIT` state. While the connection is in this state, the last receive sequence number is remembered in the TCP control block. Assume its value is 100,000.
- Before this connection leaves the `TIME_WAIT` state, a new SYN is received from the same port on the same client host (192.3.4.5, port 3000), which locates the PCB corresponding to the connection in the `TIME_WAIT` state, not the PCB for the listening server. Assume the sequence number of this new SYN is 200,000.

- Since this connection does not correspond to a listening socket in the LISTEN state, the code we just looked at is not executed. Instead, the code in Figure 28.28 is executed, and we'll see that it contains the following logic: if the sequence number of the new SYN (200,000) is greater than the last sequence number received from this client (100,000), then (1) the local variable `iss` is set to 100,000 plus 128,000, (2) the connection in the TIME_WAIT state is completely closed (its PCB and TCP control block are deleted), and (3) a jump is made to `fIndpcb` (Figure 28.5).
- This time the server's listening PCB will be located (assuming the listening server is still running), causing the code in this section to be executed. The local variable `iss` (now 228,000) is used in Figure 28.17 to initialize `tcp_iss` for the new connection.

This logic, which is allowed by RFC 1122, lets the same client and server reuse the same socket pair as long as the server does the active close. This also explains why the global variable `tcp_iss` is incremented by 64,000 each time any process issues a `connect` (Figure 30.4): to ensure that if a single client reopens the same connection with the same server repeatedly, a larger ISS is used each time, even if no data was transferred on the previous connection, and even if the 500-ms timer (which increments `tcp_iss`) has not expired since the last connection.

Initialize sequence number variables in control block

520-522 In Figure 28.17, the initial receive sequence number (`irs`) is copied from the sequence number in the SYN segment. The following two macros initialize the appropriate variables in the TCP control block:

```
#define tcp_rcvseqinit(tp) \
    (tp)->rcv_adv = (tp)->rcv_next = (tp)->irs + 1

#define tcp_sndseqinit(tp) \
    (tp)->snd_una = (tp)->snd_next = (tp)->snd_max = (tp)->snd_up = \
    (tp)->iss
```

The addition of 1 in the first macro is because the SYN occupies a sequence number.

ACK the SYN and change state

523-525 The `TF_ACKNOW` flag is set since the ACK of a SYN is not delayed. The connection state becomes `SYN_RCVD`, and the connection-establishment timer is set to 75 seconds (`TCPTV_KEEP_INIT`). Since the `TF_ACKNOW` flag is set, at the bottom of this function `tcp_output` will be called. Looking at Figure 24.16 we see that `tcp_outflags` will cause a segment with the SYN and ACK flags to be sent.

526-528 TCP is now committed to the new socket created in Figure 28.7, so the `dropsocket` flag is cleared. The code at `trimthenstep6` is jumped to, to complete processing of the SYN segment. Remember that a SYN segment can contain data, although the data cannot be passed to the application until the connection enters the ESTABLISHED state.

Completion of Active Open

Figure 28.18 shows the first part of processing when the connection is in the SYN_SENT state. TCP is expecting to receive a SYN.

```

530      /* tcp_input.c
531      * If the state is SYN_SENT:
532      * if seg contains an ACK, but not for our SYN, drop the input.
533      * if seg contains an RST, then drop the connection.
534      * if seg does not contain SYN, then drop it.
535      * Otherwise this is an acceptable SYN segment
536      * initialize tp->rcv_nxt and tp->irs
537      * if seg contains ack then advance tp->snd_una
538      * if SYN has been acked change to ESTABLISHED else SYN_RCVD state
539      * arrange for segment to be acked (eventually)
540      * continue processing rest of data/controls, beginning with URG
541      */
542  case TCPS_SYN_SENT:
543      if ((tiflags & TH_ACK) &&
544          (SEQ_LEQ(ti->ti_ack, tp->iss) ||
545           SEQ_GT(ti->ti_ack, tp->snd_max)))
546          goto dropwithreset;
547      if (tiflags & TH_RST) {
548          if (tiflags & TH_ACK)
549              tp = tcp_drop(tp, ECONNREFUSED);
550          goto drop;
551      }
552      if ((tiflags & TH_SYN) == 0)
553          goto drop;

```

Figure 28.18 `tcp_input` function: check if SYN in response to active open.

Verify received ACK

When TCP sends a SYN in response to an active open by a process, we'll see in Figure 30.4 that the connection's `iss` is copied from the global `tcp_iss` and the macro `tcp_sendseqinit` (shown at the end of the previous section) is executed. Assuming the ISS is 365, Figure 28.19 shows the send sequence variables after the SYN is sent by `tcp_output`.

	SYN	366	367	...
	↑	↑		
snd_una = 365		snd_nxt = 366		
snd_up = 365		snd_max = 366		

Figure 28.19 Send variables after SYN is sent with sequence number 365.

`tcp_sendseqinit` sets all four of these variables to 365, then Figure 26.31 increments two of them to 366 when the SYN segment is output. Therefore, if the received segment in Figure 28.18 contains an ACK, and if the acknowledgment field is less than or equal to `iss` (365) or greater than `snd_max` (366), the ACK is invalid, causing the segment to be dropped and an RST sent in reply. Notice that the received segment for a connection in the SYN_SENT state need not contain an ACK. It can contain only a SYN, which is called a *simultaneous open* (Figure 24.15), and is described shortly.

Process and drop RST segment

547-551 If the received segment contains an RST, it is dropped. But the ACK flag was checked first because receipt of an acceptable ACK (which was just verified) *and* an RST in response to a SYN is how the other end tells TCP that its connection request was refused. Normally this is caused by the server process not being started on the other host. In this case `tcp_drop` sets the socket's `so_error` variable, causing an error to be returned to the process that called `connect`.

Verify SYN flag set

552-553 If the SYN flag is not set in the received segment, it is dropped.

The remainder of this case handles the receipt of a SYN (with an optional ACK) in response to TCP's SYN. The next part of `tcp_input`, shown in Figure 28.20, continues processing the SYN.

Process ACK

554-558 If the received segment contains an ACK, `snd_una` is set to the acknowledgment field. In Figure 28.19, `snd_una` becomes 366, since 366 is the only acceptable value for the acknowledgment field. If `snd_nxt` is less than `snd_una` (which shouldn't happen, given Figure 28.19), `snd_nxt` is set to `snd_una`.

Turn off connection-establishment timer

559 The connection-establishment timer is turned off.

This is a bug. This timer should be turned off only if the ACK flag is set, since the receipt of a SYN without an ACK is a simultaneous open, and doesn't mean the other end received TCP's SYN.

Initialize receive sequence numbers

560-562 The initial receive sequence number is copied from the sequence number of the received segment. The `tcp_rcvseqinit` macro (shown at the end of the previous section) initializes `rcv_adv` and `rcv_nxt` to the receive sequence number, plus 1. The `TF_ACKNOW` flag is set, causing `tcp_output` to be called at the bottom of this function. The segment it sends will contain `rcv_nxt` as the acknowledgment field (Figure 26.27), which acknowledges the SYN just received.

563-564 If the received segment contains an ACK, and if `snd_una` is greater than the `ISS` for the connection, the active open is complete, and the connection is established.

This second test appears superfluous. At the beginning of Figure 28.20 `snd_una` was set to the received acknowledgment field if the ACK flag was on. Also the `if` following the `case`

```

554         if (tiflags & TH_ACK) {                                     tcp_input.c
555             tp->snd_una = ti->ti_ack;
556             if (SEQ_LT(tp->snd_nxt, tp->snd_una))
557                 tp->snd_nxt = tp->snd_una;
558         }
559         tp->t_timer[TCPT_REXMT] = 0;
560         tp->irs = ti->ti_seq;
561         tcp_rcvseqinit(tp);
562         tp->t_flags |= TF_ACKNOW;
563         if (tiflags & TH_ACK && SEQ_GT(tp->snd_una, tp->iss)) {
564             tcpstat.tcps_connects++;
565             soisconnected(so);
566             tp->t_state = TCPS_ESTABLISHED;
567             /* Do window scaling on this connection? */
568             if ((tp->t_flags & (TF_RCVD_SCALE | TF_REQ_SCALE)) ==
569                 (TF_RCVD_SCALE | TF_REQ_SCALE)) {
570                 tp->snd_scale = tp->requested_s_scale;
571                 tp->rcv_scale = tp->request_r_scale;
572             }
573             (void) tcp_reass(tp, (struct tcpiphdr *) 0,
574                             (struct mbuf *) 0);
575             /*
576              * if we didn't have to retransmit the SYN,
577              * use its rtt as our initial srtt & rtt var.
578              */
579             if (tp->t_rtt)
580                 tcp_xmit_timer(tp, tp->t_rtt);
581         } else
582             tp->t_state = TCPS_SYN_RECEIVED;

```

Figure 28.20 `tcp_input` function: process received SYN in response to an active open.

statement in Figure 28.18 verified that the received acknowledgment field is greater than the ISS. So at this point in the code, if the ACK flag is set, we're already guaranteed that `snd_una` is greater than the ISS.

Connection is established

565-566 `soisconnected` sets the socket state to connected, and the state of the TCP connection is set to ESTABLISHED.

Check for window scale option

567-572 If TCP sent the window scale option in its SYN and the received SYN also contains the option, the option is enabled and the two variables `snd_scale` and `rcv_scale` are set. Since the TCP control block is initialized to 0 by `tcp_newtcpcb`, these two variables correctly default to 0 if the window scale option is not used.

Pass any queued data to process

573-574 Since data can arrive for a connection before the connection is established, any such data is now placed in the receive buffer by calling `tcp_reass` with a null pointer as the second argument.

This test is unnecessary. In this piece of code, TCP has just received the SYN with an ACK that moves it from the SYN_SENT state to the ESTABLISHED state. If data appears with this received SYN segment, it isn't processed until the label `do_data` near the end of the function. If TCP just received a SYN without an ACK (a simultaneous open) but with some data, that data is handled later (Figure 29.2) when the ACK is received that moves the connection from the SYN_RCVD state to the ESTABLISHED state.

Although it is valid for data to accompany a SYN, and Net/3 handles this type of received segment correctly, Net/3 never generates such a segment.

Update RTT estimators

575-580 If the SYN that is ACKed was being timed, `tcp_xmit_timer` initializes the RTT estimators based on the measured RTT for the SYN.

TCP ignores a received timestamp option here, and checks only the `t_rtt` counter. TCP sends a timestamp in a SYN generated by an active open (Figure 26.24) and if the other end agrees to the option, the other end should echo the received timestamp in its SYN. (Net/3 echoes the received timestamp in a SYN in Figure 28.10.) This would allow TCP to use the received timestamp here, instead of `t_rtt`, but since both have the same precision (500 ms) there's no advantage in using the timestamp value. The real advantage in using the timestamp option, instead of the `t_rtt` counter, is with large pipes, when lots of segments are in flight at once, providing more RTT timings and (it is hoped) better estimators.

Simultaneous open

581-582 When TCP receives a SYN without an ACK in the SYN_SENT state, it is a simultaneous open and the connection moves to the SYN_RCVD state.

The next piece of code, shown in Figure 28.21, handles any data received with the SYN. The label `trimthenstep6` is also jumped to at the end of Figure 28.17.

```

583         trimthenstep6:
584         /*
585          * Advance ti->ti_seq to correspond to first data byte.
586          * If data, trim to stay within window,
587          * dropping FIN if necessary.
588          */
589         ti->ti_seq++;
590         if (ti->ti_len > tp->rcv_wnd) {
591             todrop = ti->ti_len - tp->rcv_wnd;
592             m_adj(m, -todrop);
593             ti->ti_len = tp->rcv_wnd;
594             tiflags &= ~TH_FIN;
595             tcpstat.tcps_rcvpackafterwin++;
596             tcpstat.tcps_rcvbyteafterwin += todrop;
597         }
598         tp->snd_wll = ti->ti_seq - 1;
599         tp->rcv_up = ti->ti_seq;
600         goto step6;
601     }

```

Figure 28.21 `tcp_input` function: common processing for receipt of SYN.

584-589 The sequence number of the segment is incremented by 1 to account for the SYN. If there is any data in the segment, `ti_seq` now contains the starting sequence number of the first byte of data.

Drop any received data that follows receive window

590-597 `ti_len` is the number of data bytes in the segment. If it is greater than the receive window, the excess data (`ti_len` minus `rcv_wnd`) is dropped by `m_adj`. The negative argument to this function causes the data to be trimmed from the end of the mbuf chain (Figure 2.20). `ti_len` is updated to be the new amount of data in the mbuf chain and in case the FIN flag was set, it is cleared. This is because the FIN would follow the final data byte, which was just discarded because it was outside the receive window.

If too much data is received with a SYN, and if the SYN is in response to an active open, the other end received TCP's SYN, which contained a window advertisement. This means the other end ignored the advertised window and is exhibiting unsocial behavior. But if too much data accompanies a SYN performing an active open, the other end has not received a window advertisement, so it has to guess how much data can accompany its SYN.

Force update of window variables

598-599 `snd_wll` is set the received sequence number minus 1. We'll see in Figure 29.15 that this causes the three window update variables, `snd_wnd`, `snd_wll`, and `snd_wl2`, to be updated. The receive urgent pointer (`rcv_up`) is set to the received sequence number. A jump is made to `step6`, which refers to a step in RFC 793, and we cover this in Figure 29.15.

28.7 PAWS: Protection Against Wrapped Sequence Numbers

The next part of `tcp_input`, shown in Figure 28.22, provides protection against wrapped sequence numbers: the PAWS algorithm from RFC 1323. Also recall our discussion of the timestamp option in Section 26.6.

Basic PAWS test

602-613 `ts_present` was set by `tcp_dooptions` if a timestamp option was present. If the following three conditions are all true, the segment is dropped:

1. the RST flag is not set (Exercise 28.8),
2. TCP has received a valid timestamp from this peer (`ts_recent` is nonzero), and
3. the received timestamp in this segment (`ts_val`) is less than the previously received timestamp from this peer.

PAWS is built on the premise that the 32-bit timestamp values wrap around at a much lower frequency than the 32-bit sequence numbers, on a high-speed connection. Exercise 28.6 shows that even at the highest possible timestamp counter frequency (incrementing by 1 bit every millisecond), the sign bit of the timestamp wraps around only every 24 days. On a high-speed network such as a gigabit network, the sequence

```

602  /* tcp_input.c
603  * States other than LISTEN or SYN_SENT.
604  * First check timestamp, if present.
605  * Then check that at least some bytes of segment are within
606  * receive window. If segment begins before rcv_nxt,
607  * drop leading data (and SYN); if nothing left, just ack.
608  *
609  * RFC 1323 PAWS: If we have a timestamp reply on this segment
610  * and it's less than ts_recent, drop it.
611  */
612  if (ts_present && (tiflags & TH_RST) == 0 && tp->ts_recent &&
613      TSTMP_LT(ts_val, tp->ts_recent)) {
614
615      /* Check to see if ts_recent is over 24 days old. */
616      if ((int) (tcp_now - tp->ts_recent_age) > TCP_PAWS_IDLE) {
617          /*
618           * Invalidate ts_recent. If this segment updates
619           * ts_recent, the age will be reset later and ts_recent
620           * will get a valid value. If it does not, setting
621           * ts_recent to zero will at least satisfy the
622           * requirement that zero be placed in the timestamp
623           * echo reply when ts_recent isn't valid. The
624           * age isn't reset until we get a valid ts_recent
625           * because we don't want out-of-order segments to be
626           * dropped when ts_recent is old.
627           */
628          tp->ts_recent = 0;
629      } else {
630          tcpstat.tcps_rcvduppack++;
631          tcpstat.tcps_rcvdupbyte += ti->ti_len;
632          tcpstat.tcps_pawsdrop++;
633          goto dropafterack;
634      }

```

Figure 28.22 tcp_input function: process timestamp option.

number can wrap in 17 seconds (Section 24.3 of Volume 1). Therefore, if the received timestamp value is less than the most recent one from this peer, this segment is old and must be discarded (subject to the outdated timestamp test that follows). The packet might be discarded later in the input processing because the sequence number is "old," but PAWS is intended for high-speed connections where the sequence numbers can wrap quickly.

Notice that the PAWS algorithm is symmetric: it not only discards duplicate data segments but also discards duplicate ACKs. All received segments are subject to PAWS. Recall that the header prediction code also applied the PAWS test (Figure 28.11).

Check for outdated timestamp

614-627 There is a small possibility that the reason the PAWS test fails is because the connection has been idle for a long time. The received segment is not a duplicate; it is just that

because the connection has been idle for so long, the peer's timestamp value has wrapped around when compared to the most recent timestamp from that peer.

Whenever `ts_recent` is copied from the timestamp in a received segment, `ts_recent_age` records the current time (`tcp_now`). If the time at which `ts_recent` was saved is more than 24 days ago, it is set to 0 to invalidate it. The constant `TCP_PAWS_IDLE` is defined to be $(24 \times 24 \times 60 \times 60 \times 2)$, the final 2 being the number of ticks per second. The received segment is not dropped in this case, since the problem is not a duplicated segment, but an outdated timestamp. See also Exercises 28.6 and 28.7.

Figure 28.23 shows an example of an outdated timestamp. The system on the left is a non-Net/3 system that increments its timestamp clock at the highest frequency allowed by RFC 1323: once every millisecond. The system on the right is a Net/3 system.

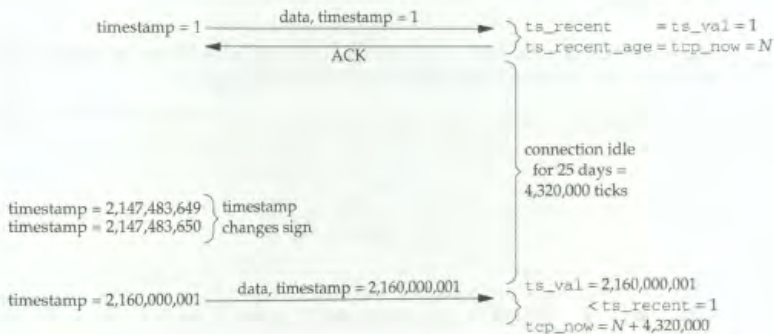


Figure 28.23 Example of outdated timestamp.

When the data segment arrives with a timestamp of 1, that value is saved in `ts_recent` and `ts_recent_age` is set to the current time (`tcp_now`), as shown in Figures 28.11 and 28.35. The connection is then idle for 25 days, during which time `tcp_now` will increase by 4,320,000 ($25 \times 24 \times 60 \times 60 \times 2$). During these 25 days the other end's timestamp clock will increase by 2,160,000,000 ($25 \times 24 \times 60 \times 60 \times 1000$). During this interval the timestamp "changes sign" with regard to the value 1, that is, 2,147,483,649 is greater than 1, but 2,147,483,650 is less than 1 (recall Figure 24.26). Therefore, when the data segment is received with a timestamp of 2,160,000,001, this value is less than `ts_recent` (1), when compared using the `TSTMP_LT` macro, so the PAWS test fails. But since `tcp_now` minus `ts_recent_age` is greater than 24 days, the reason for the failure is that the connection has been idle for more than 24 days, and the segment is accepted.

Drop duplicate segment

628-633

The segment is determined to be a duplicate based on the PAWS algorithm, and the timestamp is not outdated. It is dropped, after being acknowledged (since all duplicate segments are acknowledged).

Figure 24.5 shows a much smaller value for `tcps_pawdrop` (22) than for `tcps_rcvduppack` (46,953). This is probably because fewer systems support the timestamp option today, causing most duplicate packets to be discarded by later tests in TCP's input processing instead of by PAWS.

28.8 Trim Segment so Data is Within Window

This section trims the received segment so that it contains only data that is within the advertised window:

- duplicate data at the beginning of the received segment is discarded, and
- data that is beyond the end of the window is discarded from the end of the segment.

What remains is new data within the window. The code shown in Figure 28.24, checks if there is any duplicate data at the beginning of the segment.

```

635     todrop = tp->rcv_nxt - ti->ti_seq;
636     if (todrop > 0) {
637         if (tiflags & TH_SYN) {
638             tiflags &= ~TH_SYN;
639             ti->ti_seq++;
640             if (ti->ti_urp > 1)
641                 ti->ti_urp--;
642             else
643                 tiflags &= ~TH_URG;
644             todrop--;
645         }

```

tcp_input.c

tcp_input.c

Figure 28.24 `tcp_input` function: check for duplicate data at beginning of segment.

Check if any duplicate data at front of segment

635-636 If the starting sequence number of the received segment (`ti_seq`) is less than the next receive sequence number expected (`rcv_nxt`), data at the beginning of the segment is old and `todrop` will be greater than 0. These data bytes have already been acknowledged and passed to the application (Figure 24.18).

Remove duplicate SYN

637-645 If the SYN flag is set, it refers to the first sequence number in the segment, which is known to be old. The SYN flag is cleared and the starting sequence number of the segment is incremented by 1 to skip over the duplicate SYN. Furthermore, if the urgent offset in the received segment (`ti_urp`) is greater than 1, it must be decremented by 1, since the urgent offset is relative to the starting sequence number, which was just incremented. If the urgent offset is 0 or 1, it is left alone, but in case it was 1, the URG flag is cleared. Finally `todrop` is decremented by 1 (since the SYN occupies a sequence number).

The handling of duplicate data at the front of the segment continues in Figure 28.25.


```

646         if (todrop >= ti->ti_len) {
647             tcpstat.tcps_rcvdupack++;
648             tcpstat.tcps_rcvdupbyte += ti->ti_len;
649             /*
650              * If segment is just one to the left of the window,
651              * check two special cases:
652              * 1. Don't toss RST in response to 4.2-style keepalive.
653              * 2. If the only thing to drop is a FIN, we can drop
654              *    it, but check the ACK or we will get into FIN
655              *    wars if our FINs crossed (both CLOSING).
656              * In either case, send ACK to resynchronize,
657              * but keep on processing for RST or ACK.
658              */
659             if ((tiflags & TH_FIN && todrop == ti->ti_len + 1)
660                 ) {
661                 todrop = ti->ti_len;
662                 tiflags &= ~TH_FIN;
663                 tp->t_flags |= TF_ACKNOW;
664             } else {
665                 /*
666                  * Handle the case when a bound socket connects
667                  * to itself. Allow packets with a SYN and
668                  * an ACK to continue with the processing.
669                  */
670                 if (todrop != 0 || (tiflags & TH_ACK) == 0)
671                     goto dropafterack;
672             }
673         } else {
674             tcpstat.tcps_rcvpartdupack++;
675             tcpstat.tcps_rcvpartdupbyte += todrop;
676         }
677         m_adj(m, todrop);
678         ti->ti_seq += todrop;
679         ti->ti_len -= todrop;
680         if (ti->ti_urp > todrop)
681             ti->ti_urp -= todrop;
682         else {
683             tiflags &= ~TH_URG;
684             ti->ti_urp = 0;
685         }
686     }

```

Figure 28.25 tcp_input function: handle completely duplicate segment.

Check for entire duplicate packet

646-648 If the amount of duplicate data at the front of the segment is greater than or equal to the size of the segment, the entire segment is a duplicate.

Check for duplicate FIN

649-663 The next check is whether the FIN is duplicated. Figure 28.26 shows an example of this.

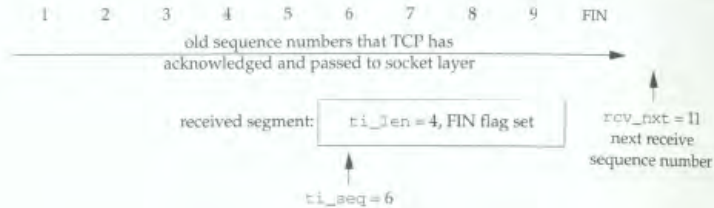


Figure 28.26 Example of duplicate packet with FIN flag set.

In this example `todrop` equals 5, which is greater than or equal to `ti_len` (4). Since the FIN flag is set and `todrop` equals `ti_len` plus 1, `todrop` is set to 4, the FIN flag is cleared, and the `TF_ACKNOW` flag is set, forcing an immediate ACK to be sent at the end of this function. This example also works for other segments if `ti_seq` plus `ti_len` equals 10.

The code contains the comment regarding 4.2BSD keepalives. This code (another test within the `if` statement) is omitted.

Generate duplicate ACK

664-672 If `todrop` is nonzero (the completely duplicate segment contains data) or the ACK flag is not set, the segment is dropped and an ACK is generated by `dropafterack`. This normally occurs when the other end did not receive our ACK, causing the other end to retransmit the segment. TCP generates another ACK.

Handle simultaneous open or self-connect

664-672 This code also handles either a simultaneous open or a socket that connects to itself. We go over both of these scenarios in the next section. If `todrop` equals 0 (there is no data in the completely duplicate segment) and the ACK flag is set, processing is allowed to continue.

This `if` statement is new with 4.4BSD. Earlier Berkeley-derived systems just had a jump to `dropafterack`. These systems could not handle either a simultaneous open or a socket connecting to itself.

Nevertheless, the piece of code in this figure still has bugs, which we describe at the end of this section.

Update statistics for partial duplicate segments

673-676 This `else` clause is executed when `todrop` is less than the segment length: only part of the segment contains duplicate bytes.

Remove duplicate data and update urgent offset

677-685 The duplicate bytes are removed from the front of the mbuf chain by `m_adj` and the starting sequence number and length adjusted appropriately. If the urgent offset points to data still in the mbuf, it is also adjusted. Otherwise the urgent offset is set to 0 and the URG flag is cleared.

The next part of input processing, shown in Figure 28.27, handles data that arrives after the process has terminated.

```

687  /*-----tcp_input.c
688  * If new data is received on a connection after the
689  * user processes are gone, then RST the other end.
690  */
691  if ((so->so_state & SS_NOFDREF) &&
692      tp->t_state > TCPS_CLOSE_WAIT && ti->ti_len) {
693      tp = tcp_close(tp);
694      tcpstat.tcps_rcvafterclose++;
695      goto dropwithreset;
696  }
-----tcp_input.c

```

Figure 28.27 tcp_input function: handle data that arrives after the process terminates.

687-696 If the socket has no descriptor referencing it, the process has closed the connection (the state is any one of the five with a value greater than CLOSE_WAIT in Figure 24.16), and there is data in the received segment, the connection is closed. The segment is then dropped and an RST is output.

Because of TCP's half-close, if a process terminates unexpectedly (perhaps it is terminated by a signal), when the kernel closes all open descriptors as part of process termination, a FIN is output by TCP. The connection moves into the FIN_WAIT_1 state. But the receipt of the FIN by the other end doesn't tell TCP whether this end performed a half-close or a full-close. If the other end assumes a half-close, and sends more data, it will receive an RST from the code in Figure 28.27.

The next piece of code, shown in Figure 28.28, removes any data from the end of the received segment that is beyond the right edge of the advertised window.

Calculate number of bytes beyond right edge of window

697-703 `todrop` contains the number of bytes of data beyond the right edge of the window. For example, in Figure 28.29, `todrop` would be $(6 + 5)$ minus $(4 + 6)$, or 1.

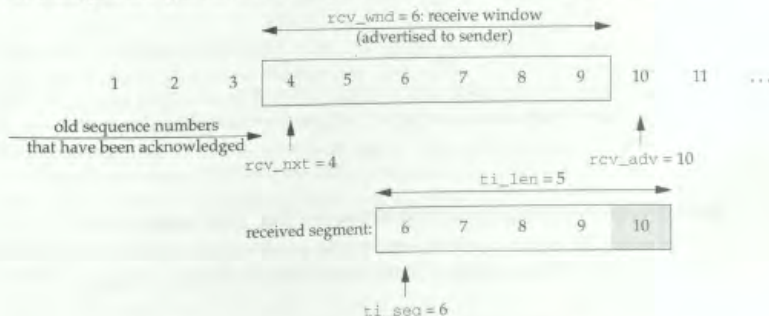


Figure 28.29 Example of received segment with data beyond right edge of window.

```

697  /*
698  * If segment ends after window, drop trailing data
699  * (and PUSH and FIN); if nothing left, just ACK.
700  */
701  todrop = (ti->ti_seq + ti->ti_len) - (tp->rcv_nxt + tp->rcv_wnd);
702  if (todrop > 0) {
703      tcpstat.tcps_rcvpackafterwin++;
704      if (todrop >= ti->ti_len) {
705          tcpstat.tcps_rcvbyteafterwin += ti->ti_len;
706          /*
707           * If a new connection request is received
708           * while in TIME_WAIT, drop the old connection
709           * and start over if the sequence numbers
710           * are above the previous ones.
711           */
712          if (tiflags & TH_SYN &&
713              tp->t_state == TCPS_TIME_WAIT &&
714              SEQ_GT(ti->ti_seq, tp->rcv_nxt)) {
715              iss = tp->rcv_nxt + TCP_ISSINCR;
716              tp = tcp_close(tp);
717              goto findpcb;
718          }
719          /*
720           * If window is closed can only take segments at
721           * window edge, and have to drop data and PUSH from
722           * incoming segments. Continue processing, but
723           * remember to ack. Otherwise, drop segment
724           * and ack.
725           */
726          if (tp->rcv_wnd == 0 && ti->ti_seq == tp->rcv_nxt) {
727              tp->t_flags |= TP_ACKNOW;
728              tcpstat.tcps_rcvwinprobe++;
729          } else
730              goto dropafterack;
731      } else
732          tcpstat.tcps_rcvbyteafterwin += todrop;
733      m_adj(m, -todrop);
734      ti->ti_len -= todrop;
735      tiflags &= ~(TH_PUSH | TH_FIN);
736  }

```

Figure 28.28 tcp_input function: remove data beyond right edge of window.

Check for new incarnation of a connection in the TIME_WAIT state

704-718 If `todrop` is greater than or equal to the length of the segment, the entire segment will be dropped. If the following three conditions are all true:

1. the SYN flag is set, and
2. the connection is in the TIME_WAIT state, and
3. the new starting sequence number is greater than the final sequence number for the connection,

this is a request for a new incarnation of a connection that was recently terminated and is currently in the TIME_WAIT state. This is allowed by RFC 1122, but the ISS for the new connection must be greater than the last sequence number used (`rcv_nxt`). TCP adds 128,000 (`TCP_ISSINCR`), which becomes the ISS when the code in Figure 28.17 is executed. The PCB and TCP control block for the connection in the TIME_WAIT state is discarded by `tcp_close`. A jump is made to `findpcb` (Figure 28.5) to locate the PCB for the listening server, assuming it is still running. The code in Figure 28.7 is then executed, creating a new socket for the new connection, and finally the code in Figures 28.16 and 28.17 will complete the new connection request.

Check for probe of closed window

719-728 If the receive window is closed (`rcv_wnd` equals 0) and the received segment starts at the left edge of the window (`rcv_nxt`), then the other end is probing TCP's closed window. An immediate ACK is sent as the reply, even though the ACK may still advertise a window of 0. Processing of the received segment also continues for this case.

Drop other segments that are completely outside window

729-730 The entire segment lies outside the window and it is not a window probe, so the segment is discarded and an ACK is sent as the reply. This ACK will contain the expected sequence number.

Handle segments that contain some valid data

731-735 The data to the right of the window is discarded from the mbuf chain by `m_adj` and `ti_len` is updated. In the case of a probe into a closed window, this discards all the data in the mbuf chain and sets `ti_len` to 0. Finally the FIN and PSH flags are cleared.

When to Drop an ACK

The code in Figure 28.25 has a bug that causes a jump to `dropafterack` in several cases when the code should fall through for further processing of the segment [Carlson 1993; Lanciani 1993]. In an actual scenario, when both ends of a connection had a hole in the data on the reassembly queue and both ends enter the persist state, the connection becomes deadlocked as both ends throw away perfectly good ACKs.

The fix is to simplify the code at the beginning of Figure 28.25. Instead of jumping to `dropafterack`, a completely duplicate segment causes the FIN flag to be turned off and an immediate ACK to be generated at the end of the function. Lines 646-676 in Figure 28.25 are replaced with the code shown in Figure 28.30. This code also corrects another bug present in the original code (Exercise 28.9).

```

if (todrop > ti->ti_len ||
    todrop == ti->ti_len && (tiflags & TH_FIN) == 0) {
    /*
     * Any valid FIN must be to the left of the window.
     * At this point the FIN must be a duplicate or
     * out of sequence; drop it.
     */
    tiflags &= ~TH_FIN;

    /*
     * Send an ACK to resynchronize and drop any data.
     * But keep on processing for RST or ACK.
     */
    tp->t_flags |= TF_ACKNOW;
    todrop = ti->ti_len;
    tcpstat.tcps_rcvdupbyte += todrop;
    tcpstat.tcps_rcvduppack++;
} else {
    tcpstat.tcps_rcvpartdupack++;
    tcpstat.tcps_rcvpartdupbyte += todrop;
}

```

Figure 28.30 Correction for lines 646–676 of Figure 28.25.

28.9 Self-Connects and Simultaneous Opens

It is instructive to look at the steps involved in a socket connecting to itself to see how the one-line fix to Figure 28.25 that was added to 4.4BSD allows this. This same fix allowed simultaneous opens to work, which wasn't handled correctly prior to 4.4BSD.

A process creates a socket and connects it to itself using the system calls: `socket`, `bind` a local port (say 3000), and then `connect` to this same port and some local IP address. If the `connect` succeeds, the socket is connected to itself: anything written to the socket can be read back from the socket. This is similar to a full-duplex pipe, but with a single descriptor instead of two descriptors. Although this is of limited use within a process, we'll see that the state transitions are the same as they are for a simultaneous open. If your system doesn't allow a socket to connect to itself, it probably doesn't handle simultaneous opens correctly either, and the latter are required by RFC 1122. Some people are surprised that a self-connect even works, given that a single Internet PCB and a single TCP control block are used. But TCP is a full-duplex, symmetric protocol and it maintains separate variables for each direction of data flow.

Figure 28.31 shows the send sequence space when the process calls `connect`. A SYN segment is sent and the state becomes `SYN_SENT`.

The SYN is received and processed in Figures 28.18 and 28.20, but since the SYN does not contain an ACK the resulting state is `SYN_RCVD`. According to the state transition diagram (Figure 24.15), this looks like a simultaneous open. Figure 28.32 shows the receive sequence space.

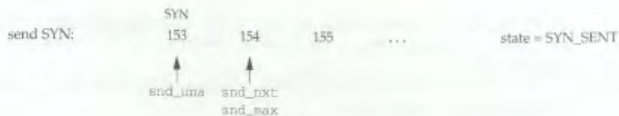


Figure 28.31 Send sequence space when SYN is sent for self-connect.



Figure 28.32 Receive sequence space after received SYN is processed.

Figure 28.20 sets the `TF_ACKNOW` flag and the segment generated by `tcp_output` will contain a SYN and an ACK (the `tcp_outFlags` value in Figure 24.16). The sequence number of the SYN is 153 and the acknowledgment number is 154.

Nothing changes in the send sequence space from Figure 28.20, except the state is now `SYN_SENT`. Figure 28.33 shows the receive sequence space when the segment with the SYN and ACK is received.

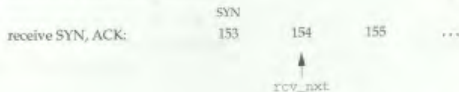


Figure 28.33 Receive sequence space when segment with SYN and ACK received.

Since the connection state is `SYN_RCVD`, the segment is not processed by the active open or passive open code that we saw earlier in this chapter. It must be processed by the `SYN_RCVD` code that we'll examine in Figure 29.2. But it is first processed by Figure 28.24, and it looks like a duplicate SYN:

```
todrop = rcv_nxt - ti_seq
        = 154 - 153
        = 1
```

Since the SYN flag is set, the flag is cleared, `ti_seq` becomes 154, and `todrop` becomes 0. But the test at the beginning of Figure 28.25 is true, because `todrop` equals the length of the segment (0). The segment is counted as a duplicate packet and the code with the comment "Handle the case when a bound socket connects to itself" is executed. Earlier releases jumped to `dropafterack`, which skipped the necessary code to handle the `SYN_RCVD` state, preventing the connection from ever being established. Instead, Net/3 continues processing the received segment if `todrop` equals 0 and the

ACK flag is set, both of which are true in this example. This allows the SYN_RCVD processing to happen later in the function, which moves the connection to the ESTABLISHED state.

It is also interesting to look at the sequence of function calls in this self-connect. This is shown in Figure 28.34.

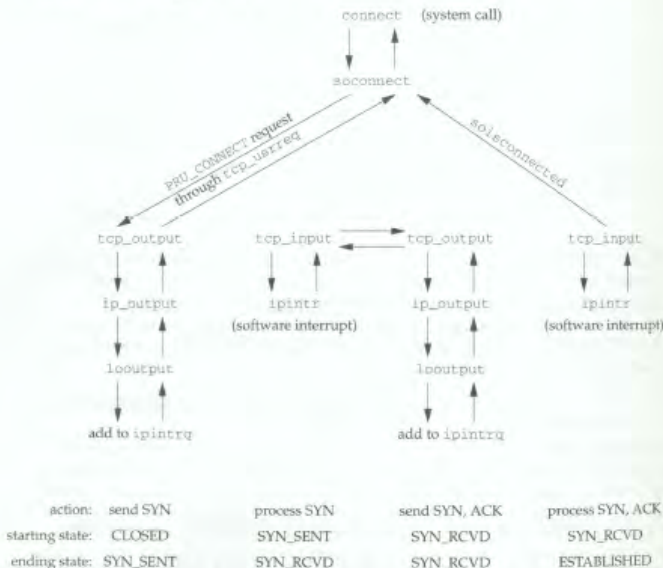


Figure 28.34 Sequence of function calls for self-connect.

The order of the operations goes from the left to the right. The steps that we show begin with the process calling `connect`. This issues the `PRU_CONNECT` request, which sends a SYN down the protocol stack. Since the segment is destined for the host's own IP address it is routed to the loopback interface, which adds the segment to `ipintrq` and generates a software interrupt.

The software interrupt causes `ipintr` to execute, which calls `tcp_input`. This function calls `tcp_output`, causing a SYN segment with an ACK to be sent down the protocol stack. It is again added to `ipintrq` by the loopback interface, and a software interrupt is generated. When this interrupt is processed by `ipintr`, the function `tcp_input` is called, and it moves the connection to the ESTABLISHED state.

28.10 Record Timestamp

The next part of `tcp_input`, shown in Figure 28.35, handles a received timestamp option.

```

737  /*
738  * If last ACK falls within this segment's sequence numbers,
739  * record its timestamp.
740  */
741  if (ts_present && SEQ_LEQ(ti->ti_seq, tp->last_ack_sent) &&
742      SEQ_LT(tp->last_ack_sent, ti->ti_seq + ti->ti_len +
743          ((ti_flags & (TH_SYN | TH_FIN)) != 0))) {
744      tp->ts_recent_age = tcp_now;
745      tp->ts_recent = ts_val;
746  }

```

tcp_input.c

Figure 28.35 `tcp_input` function: record timestamp.

137-746 If the received segment contains a timestamp, the timestamp value is saved in `ts_recent`. We discussed in Section 26.6 how this code used by Net/3 is flawed. The expression

```
((ti_flags & (TH_SYN|TH_FIN)) != 0)
```

is 0 if neither of the two flags is set, or 1 if either is set. This effectively adds 1 to `ti_len` if either flag is set.

28.11 RST Processing

Figure 28.36 shows the `switch` statement to handle the RST flag, which depends on the connection state.

SYN_RCVD state

759-761 The socket's error code is set to `ECONNREFUSED`, and a jump is made a few lines forward to close the socket.

This state can be entered from two directions. Normally it is entered from the `LISTEN` state, after a `SYN` has been received. TCP replied with a `SYN` and an `ACK` but received an `RST` in reply. Perhaps the other end sent its `SYN` and then terminated before the reply arrived, causing it to send an `RST`. In this case the socket referred to by `so` is the new socket created by `sonewconn` in Figure 28.7. Since `dropsocket` will still be true, the socket is discarded at the label `drop`. The listening descriptor isn't affected at all. This is why we show the state transition from `SYN_RCVD` back to `LISTEN` in Figure 24.15.

This state can also be entered by a simultaneous open, after a process has called `connect`. In this case the socket error is returned to the process.

```

747  /*
748  * If the RST bit is set examine the state:
749  *   SYN_RECEIVED state:
750  *   If passive open, return to LISTEN state.
751  *   If active open, inform user that connection was refused.
752  *   ESTABLISHED, FIN_WAIT_1, FIN_WAIT2, CLOSE_WAIT states:
753  *   Inform user that connection was reset, and close tcb.
754  *   CLOSING, LAST_ACK, TIME_WAIT states
755  *   Close the tcb.
756  */
757  if (tiflags & TH_RST)
758      switch (tp->t_state) {
759
760      case TCPS_SYN_RECEIVED:
761          so->so_error = ECONNREFUSED;
762          goto close;
763
764      case TCPS_ESTABLISHED:
765      case TCPS_FIN_WAIT_1:
766      case TCPS_FIN_WAIT_2:
767      case TCPS_CLOSE_WAIT:
768          so->so_error = ECONNRESET;
769          close;
770          tp->t_state = TCPS_CLOSED;
771          tcpstat.tcps_drops++;
772          tp = tcp_close(tp);
773          goto drop;
774
775      case TCPS_CLOSING:
776      case TCPS_LAST_ACK:
777      case TCPS_TIME_WAIT:
778          tp = tcp_close(tp);
779          goto drop;
780      }

```

Figure 28.36 tcp_input function: process RST flag.

Other states

762-777 The receipt of an RST in the ESTABLISHED, FIN_WAIT_1, FIN_WAIT_2, or CLOSE_WAIT states returns the error ECONNRESET. In the CLOSING, LAST_ACK, and TIME_WAIT state an error is not generated, since the process has closed the socket.

Allowing an RST to terminate a connection in the TIME_WAIT state circumvents the reason this state exists. RFC 1337 [Braden 1992] discusses this and other forms of "TIME_WAIT assassination hazards" and recommends not letting an RST prematurely terminate the TIME_WAIT state. See Exercise 28.10 for an example.

The next piece of code, shown in Figure 28.37, checks for erroneous SYN's and verifies that an ACK is present.

```

778      /*-----tcp_input.c
779      * If a SYN is in the window, then this is an
780      * error and we send an RST and drop the connection,
781      */
782      if (tiflags & TH_SYN) {
783          tp = tcp_drop(tp, ECONNRESET);
784          goto dropwithreset;
785      }
786      /*
787      * If the ACK bit is off we drop the segment and return.
788      */
789      if ((tiflags & TH_ACK) == 0)
790          goto drop;
-----tcp_input.c

```

Figure 28.37 tcp_input function: handle SYN-full and ACK-less segments.

778-785 If the SYN flag is still set, this is an error and the connection is dropped with the error ECONNRESET.

786-790 If the ACK flag is not set, the segment is dropped. The remainder of this function, which we continue in the next chapter, assumes the ACK flag is set.

28.12 Summary

This chapter has started our detailed look at TCP input. It continues in the next chapter.

The code in this chapter verifies the segment's checksum, processes any TCP options, handles SYNs that initiate or complete connection requests, trims excess data from the beginning or end of the segment, and processes the RST flag.

Header prediction is a successful attempt to handle common cases with the minimum amount of processing. Although the general processing steps that we've covered handle all possible cases (which they must), many segments are well behaved and the processing steps can be minimized.

Exercises

- 28.1 Given that the maximum size of a socket buffer is 262,144 in Net/3, what are the possible window scale shift factors calculated by Figure 28.7?
- 28.2 Given that the maximum size of a socket buffer is 262,144 in Net/3, what is the maximum throughput possible with a round-trip time of 60 ms? (*Hint:* See Figure 24.5 in Volume 1 and solve for the bandwidth.)
- 28.3 Why are the two timestamp values fetched using `bcopy` in Figure 28.10?
- 28.4 We mentioned in Section 26.6 that TCP correctly handles timestamp options in a format other than the one recommended in Appendix A of RFC 1323. While this is true, what is the penalty for not following the recommended format?

- 28.5 The `PRU_ATTACH` request allocates the PCB and the TCP control block, but doesn't call `tcp_template` to allocate the header template. Instead we saw in Figure 28.17 that the header template is allocated when the SYN arrives. Why doesn't the `PRU_ATTACH` request allocate this template?
- 28.6 Read RFC 1323 to determine why the limit of 24 days was chosen in Figure 28.22.
- 28.7 The comparison of `tcp_now` minus `ts_recent_age` to `TCP_PAWS_IDLE` in Figure 28.22 is also subject to sign bit wrap around, if the connection is idle for a period much longer than 24 days. With the 500-ms timestamp clock used by Net/3, when does this become a problem?
- 28.8 Read RFC 1323 to find out why RST segments are exempt from the PAWS test in Figure 28.22.
- 28.9 A client sends a SYN and the server responds with a SYN/ACK. The client moves to the ESTABLISHED state and responds with an ACK, but this ACK is lost. The server resends its SYN/ACK. Describe the processing steps when the client receives this duplicate SYN/ACK.
- 28.10 A client and server have an established connection and the server performs the active close. The connection terminates normally and the socket pair goes into the `TIME_WAIT` state on the server. Before this 2MSL wait expires on the server, the same client (i.e., the same socket pair on the client) sends a SYN to the server's socket pair but with a sequence number that is less than the ending sequence number from the previous incarnation of this connection. Describe what happens.

TCP Input (Continued)

29.1 Introduction

This chapter continues the discussion of TCP input processing, picking up where the previous chapter left off. Recall that the final test in Figure 28.37 was that either the ACK flag was set or, if not, the segment was dropped.

The ACK flag is handled, the window information is updated, the URG flag is processed, and any data in the segment is processed. Finally the FIN flag is processed and `tcp_output` is called, if required.

29.2 ACK Processing Overview

We begin this chapter with ACK processing, a summary of which is shown in Figure 29.1. The `SYN_RCVD` state is handled specially, followed by common processing for all remaining states. (Remember that a received ACK in either the `LISTEN` or `SYN_SENT` state was discussed in the previous chapter.) This is followed by special processing for the three states in which a received ACK causes a state transition, and for the `TIME_WAIT` state, in which the receipt of an ACK causes the 2MSL timer to be restarted.

29.3 Completion of Passive Opens and Simultaneous Opens

The first part of the ACK processing, shown in Figure 29.2, handles the `SYN_RCVD` state. As mentioned in the previous chapter, this handles the completion of a passive open (the common case) and also handles simultaneous opens and self-connects (the infrequent case).

```

switch (tp->t_state) {

case TCPS_SYN_RECEIVED:
    complete processing of passive open and process
    simultaneous open or self-connect;
    /* fall into ... */

case TCPS_ESTABLISHED:
case TCPS_FIN_WAIT_1:
case TCPS_FIN_WAIT_2:
case TCPS_CLOSE_WAIT:
case TCPS_CLOSING:
case TCPS_LAST_ACK:
case TCPS_TIME_WAIT:
    process duplicate ACK;
    update RTT estimators;
    if all outstanding data ACKed, turn off retransmission timer;
    remove ACKed data from socket send buffer;

    switch (tp->t_state) {

case TCPS_FIN_WAIT_1:
    if (FIN is ACKed) {
        move to FIN_WAIT_2 state;
        start FIN_WAIT_2 timer;
    }
    break;

case TCPS_CLOSING:
    if (FIN is ACKed) {
        move to TIME_WAIT state;
        start TIME_WAIT timer;
    }
    break;

case TCPS_LAST_ACK:
    if (FIN is ACKed)
        move to CLOSED state;
    break;

case TCPS_TIME_WAIT:
    restart TIME_WAIT timer;
    goto dropafterack;

}
}

```

Figure 29.1 Summary of ACK processing.

Verify received ACK

901-806

For the ACK to acknowledge the SYN that was sent, it must be greater than `snd_una` (which is set to the ISS for the connection, the sequence number of the SYN, by `tcp_sendseqinit`) and less than or equal to `snd_max`. If so, the socket is marked as connected and the state becomes ESTABLISHED.

```

791  /* tcp_input.c
792  * Ack processing.
793  */
794  switch (tp->t_state) {
795
796  /*
797  * In SYN_RECEIVED state if the ack ACKS our SYN then enter
798  * ESTABLISHED state and continue processing, otherwise
799  * send an RST.
800  */
801  case TCPS_SYN_RECEIVED:
802      if (SEQ_GT(tp->snd_una, ti->ti_ack) ||
803          SEQ_GT(ti->ti_ack, tp->snd_max))
804          goto dropwithreset;
805      tcpstat.tcps_connects++;
806      soisconnected(so);
807      tp->t_state = TCPS_ESTABLISHED;
808      /* Do window scaling? */
809      if ((tp->t_flags & (TF_RCVD_SCALE | TF_REQ_SCALE)) ==
810          (TF_RCVD_SCALE | TF_REQ_SCALE)) {
811          tp->snd_scale = tp->requested_s_scale;
812          tp->rcv_scale = tp->request_r_scale;
813      }
814      (void) tcp_reass(tp, (struct tcpiphdr *) 0, (struct mbuf *) 0);
815      tp->snd_wll = ti->ti_seq - 1;
816      /* fall into ... */

```

Figure 29.2 tcp_input function: received ACK in SYN_RCVD state.

Since `soisconnected` wakes up the process that performed the passive open (normally a server), we see that this doesn't occur until the last of the three segments in the three-way handshake has been received. If the server is blocked in a call to `accept`, that call now returns; if the server is blocked in a call to `select` waiting for the listening descriptor to become readable, it is now readable.

Check for window scale option

807-812 If TCP sent a window scale option and received one, the send and receive scale factors are saved in the TCP control block. Otherwise the default values of `snd_scale` and `rcv_scale` in the TCP control block are 0 (no scaling).

Pass queued data to process

813 Any data queued for the connection can now be passed to the process. This is done by `tcp_reass` with a null pointer as the second argument. This data would have arrived with the SYN that moved the connection into the SYN_RCVD state.

814 `snd_wll` is set to the received sequence number minus 1. We'll see in Figure 29.15 that this causes the three window update variables to be updated.

29.4 Fast Retransmit and Fast Recovery Algorithms

The next part of ACK processing, shown in Figure 29.3, handles duplicate ACKs and determines if TCP's fast retransmit and fast recovery algorithms [Jacobson 1990c] should come into play. The two algorithms are separate but are normally implemented together [Floyd 1994].

- The *fast retransmit* algorithm occurs when TCP deduces from a small number (normally 3) of consecutive duplicate ACKs that a segment has been lost and deduces the starting sequence number of the missing segment. The missing segment is retransmitted. The algorithm is mentioned in Section 4.2.2.21 of RFC 1122, which states that TCP may generate an immediate ACK when an out-of-order segment is received. We saw that Net/3 generates the immediate duplicate ACKs in Figure 27.15. This algorithm first appeared in the 4.3BSD Tahoe release and the subsequent Net/1 release. In these two implementations, after the missing segment was retransmitted, the slow start phase was entered.
- The *fast recovery* algorithm says that after the fast retransmit algorithm (that is, after the missing segment has been retransmitted), congestion avoidance but not slow start is performed. This is an improvement that allows higher throughput under moderate congestion, especially for large windows. This algorithm appeared in the 4.3BSD Reno release and the subsequent Net/2 release.

Net/3 implements both fast retransmit and fast recovery, as we describe shortly.

In the discussion of Figure 24.17 we noted that an acceptable ACK must be in the range

```
snd_una < acknowledgment field <= snd_max
```

This first test of the acknowledgment field compares it only to `snd_una`. The comparison against `snd_max` is in Figure 29.5. The reason for separating the tests is so that the following five tests can be applied to the received segment:

1. If the acknowledgment field is less than or equal to `snd_una`, and
2. the length of the received segment is 0, and
3. the advertised window (`tw`) has not changed, and
4. TCP has outstanding data that has not been acknowledged (the retransmission timer is nonzero), and
5. the received segment contains the biggest ACK TCP has seen (the acknowledgment field equals `snd_una`),

then this segment is a completely duplicate ACK. (Tests 1, 2, and 3 are in Figure 29.3; tests 4 and 5 are at the beginning of Figure 29.4.)

TCP counts the number of these duplicate ACKs that are received in a row (in the variable `t_dupacks`), and when the number reaches a threshold of 3 (`toprextthresh`), the lost segment is retransmitted. This is the *fast retransmit* algorithm described in Section 21.7 of Volume 1. It works in conjunction with the code we


```

816      /*
817      * In ESTABLISHED state: drop duplicate ACKs; ACK out-of-range
818      * ACKs. If the ack is in the range
819      * tp->snd_una < ti->ti_ack <= tp->snd_max
820      * then advance tp->snd_una to ti->ti_ack and drop
821      * data from the retransmission queue. If this ACK reflects
822      * more up-to-date window information we update our window information.
823      */
824     case TCPS_ESTABLISHED:
825     case TCPS_FIN_WAIT_1:
826     case TCPS_FIN_WAIT_2:
827     case TCPS_CLOSE_WAIT:
828     case TCPS_CLOSING:
829     case TCPS_LAST_ACK:
830     case TCPS_TIME_WAIT:
831
832     if (SEQ_LEQ(ti->ti_ack, tp->snd_una)) {
833         if (ti->ti_len == 0 && tiwin == tp->snd_wnd) {
834             tcpstat.tcps_rcvdupack++;
835             /*
836              * If we have outstanding data (other than
837              * a window probe), this is a completely
838              * duplicate ack (ie, window info didn't
839              * change), the ack is the biggest we've
840              * seen and we've seen exactly our rexmt
841              * threshold of them, assume a packet
842              * has been dropped and retransmit it.
843              * Kludge snd_next & the congestion
844              * window so we send only this one
845              * packet.
846              *
847              * We know we're losing at the current
848              * window size so do congestion avoidance
849              * (set ssthresh to half the current window
850              * and pull our congestion window back to
851              * the new ssthresh).
852              *
853              * Dup acks mean that packets have left the
854              * network (they're now cached at the receiver)
855              * so bump cwnd by the amount in the receiver
856              * to keep a constant cwnd packets in the
857              * network.
858              */
859         }
860     }
861 }

```

Figure 29.3 tcp_input function: check for completely duplicate ACK.

seen in Figure 27.15: when TCP receives an out-of-order segment, it is required to generate an immediate duplicate ACK, telling the other end that a segment might have been lost and telling it the value of the next expected sequence number. The goal of the fast retransmit algorithm is for TCP to retransmit immediately what appears to be the missing segment, instead of waiting for the retransmission timer to expire. Figure 21.7 of Volume 1 gives a detailed example of how this algorithm works.

The receipt of a duplicate ACK also tells TCP that a packet has “left the network,” because the other end had to receive an out-of-order segment to send the duplicate ACK. The *fast recovery* algorithm says that after some number of consecutive duplicate ACKs have been received, TCP should perform congestion avoidance (i.e., slow down) but need not wait for the pipe to empty between the two connection end points (slow start). The expression “a packet has left the network” means a packet has been received by the other end and has been added to the out-of-order queue for the connection. The packet is not still in transit somewhere between the two end points.

If only the first three tests shown earlier are true, the ACK is still a duplicate and is counted by the statistic `tcps_revdupack`, but the counter of the number of consecutive duplicate ACKs for this connection (`t_dupacks`) is reset to 0. If only the first test is true, the counter `t_dupacks` is reset to 0.

The remainder of the fast recovery algorithm is shown in Figure 29.4. When all five tests are true, the fast recovery algorithm processes the segment depending on the number of these consecutive duplicate ACKs that have been received.

1. `t_dupacks` equals 3 (`tcpexmtthresh`). Congestion avoidance is performed and the missing segment is retransmitted.
2. `t_dupacks` exceeds 3. Increase the congestion window and perform normal TCP output.
3. `t_dupacks` is less than 3. Do nothing.

Number of consecutive duplicate ACKs reaches threshold of 3

861-868 When `t_dupacks` reaches 3 (`tcpexmtthresh`), the value of `snd_nxt` is saved in `onxt` and the slow start threshold (`ssthresh`) is set to one-half the current congestion window, with a minimum value of two segments. This is what was done with the slow start threshold when the retransmission timer expired in Figure 25.27, but we’ll see later in this piece of code that the fast recovery algorithm does not set the congestion window to one segment, as was done with the timeout.

Turn off retransmission timer

869-870 The retransmission timer is turned off and, in case a segment is currently being timed, `t_rtt` is set to 0.

Retransmit missing segment

871-873 `snd_nxt` is set to the starting sequence number of the segment that appears to have been lost (the acknowledgment field of the duplicate ACK) and the congestion window is set to one segment. This causes `tcp_output` to send only the missing segment. (This is shown by segment 63 in Figure 21.7 of Volume 1.)

Set congestion window

874-875 The congestion window is set to the slow start threshold plus the number of segments that the other end has cached. By *cached* we mean the number of out-of-order segments that the other end has received and generated duplicate ACKs for. These cannot be passed to the process at the other end until the missing segment (which was just

```

858         if (tp->t_timer[TCPT_REXMT] == 0 ||                                     tcp_input.c
859             ti->ti_ack != tp->snd_una)
860             tp->t_dupacks = 0;
861         else if (++tp->t_dupacks == tcprexmtthresh) {
862             tcp_seq onxt = tp->snd_nxt;
863             u_int win =
864                 min(tp->snd_wnd, tp->snd_cwnd) / 2 /
865                 tp->t_maxseg;
866
867             if (win < 2)
868                 win = 2;
869             tp->snd_ssthresh = win * tp->t_maxseg;
870             tp->t_timer[TCPT_REXMT] = 0;
871             tp->t_rtt = 0;
872             tp->snd_nxt = ti->ti_ack;
873             tp->snd_cwnd = tp->t_maxseg;
874             (void) tcp_output(tp);
875             tp->snd_cwnd = tp->snd_ssthresh +
876                 tp->t_maxseg * tp->t_dupacks;
877             if (SEQ_GT(onxt, tp->snd_nxt))
878                 tp->snd_nxt = onxt;
879             goto drop;
880         } else if (tp->t_dupacks > tcprexmtthresh) {
881             tp->snd_cwnd += tp->t_maxseg;
882             (void) tcp_output(tp);
883             goto drop;
884         } else
885             tp->t_dupacks = 0;
886         break; /* beyond ACK processing (to step 6) */
887     }

```

Figure 29.4 tcp_input function: duplicate ACK processing.

sent) is received. Figures 21.10 and 21.11 in Volume 1 show what happens with the congestion window and slow start threshold when the fast recovery algorithm is in effect.

Set `snd_nxt`

876-878 The value of the next sequence number to send is set to the maximum of its previous value (`onxt`) and its current value. Its current value was modified by `tcp_output` when the segment was retransmitted. Normally this causes `snd_nxt` to be set back to its previous value, which means that only the missing segment is retransmitted, and that future calls to `tcp_output` carry on with the next segment in sequence.

Number of consecutive duplicate ACKs exceeds threshold of 3

879-883 The missing segment was retransmitted when `t_dupacks` equaled 3, so the receipt of each additional duplicate ACK means that another packet has left the network. The congestion window is incremented by one segment. `tcp_output` sends the next segment in sequence, and the duplicate ACK is dropped. (This is shown by segments 67, 69, and 71 in Figure 21.7 of Volume 1.)

884-885 This statement is executed when the received segment contains a duplicate ACK, but either the length is nonzero or the advertised window changed. Only the first of the five tests described earlier is true. The counter of consecutive duplicate ACKs is set to 0.

Skip remainder of ACK processing

886 This break is executed in three cases: (1) only the first of the five tests described earlier is true, or (2) only the first three of the five tests is true, or (3) the ACK is a duplicate, but the number of consecutive duplicates is less than the threshold of 3. For any of these cases the ACK is still a duplicate and the break goes to the end of the switch that started in Figure 29.2, which continues processing at the label `step6`.

To understand the purpose in this aggressive window manipulation, consider the following example. Assume the window is eight segments, and segments 1 through 8 are sent. Segment 1 is lost, but the remainder arrive OK and are acknowledged. After the ACKs for segments 2, 3, and 4 arrive, the missing segment (1) is retransmitted. TCP would like the subsequent ACKs for 5 through 8 to allow some of the segments starting with 9 to be sent, to keep the pipe full. But the window is 8, which prevents segments 9 and above from being sent. Therefore, the congestion window is temporarily inflated by one segment each time another duplicate ACK is received, since the receipt of the duplicate ACK tells TCP that another segment has left the pipe at the other end. When the acknowledgment of segment 1 is finally received, the next figure reduces the congestion window back to the slow start threshold. This increase in the congestion window as the duplicate ACKs arrive, and its subsequent decrease when the fresh ACK arrives, can be seen visually in Figure 21.10 of Volume 1.

29.5 ACK Processing

The ACK processing continues with Figure 29.5.

```

888      /*
889      * If the congestion window was inflated to account
890      * for the other side's cached packets, retract it.
891      */
892      if (tp->t_dupacks > tcprexmitthresh &&
893          tp->snd_cwnd > tp->snd_ssthresh)
894          tp->snd_cwnd = tp->snd_ssthresh;
895      tp->t_dupacks = 0;

896      if (SEQ_GT(ti->ti_ack, tp->snd_max)) {
897          tcpstat.tcps_rvacktoomuch++;
898          goto drops6terack;
899      }
900      acked = ti->ti_ack - tp->snd_una;
901      tcpstat.tcps_rvackpack++;
902      tcpstat.tcps_rvackbyte += acked;

```

Figure 29.5 tcp_input function: ACK processing continued.

Adjust congestion window

888-895 If the number of consecutive duplicate ACKs exceeds the threshold of 3, this is the first nonduplicate ACK after a string of four or more duplicate ACKs. The fast recovery algorithm is complete. Since the congestion window was incremented by one segment for every consecutive duplicate after the third, if it now exceeds the slow start threshold, it is set back to the slow start threshold. The counter of consecutive duplicate ACKs is set to 0.

Check for out-of-range ACK

896-899 Recall the definition of an acceptable ACK,
`snd_una < acknowledgment field <= snd_max`

If the acknowledgment field is greater than `snd_max`, the other end is acknowledging data that TCP hasn't even sent yet! This probably occurs on a high-speed connection when the sequence numbers wrap and a missing ACK reappears later. As we can see in Figure 24.5, this rarely happens (since today's networks aren't fast enough).

Calculate number of bytes acknowledged

900-902 At this point TCP knows that it has an acceptable ACK. `acked` is the number of bytes acknowledged.

The next part of ACK processing, shown in Figure 29.6, deals with RTT measurements and the retransmission timer.

Update RTT estimators

903-915 If either (1) a timestamp option was present, or (2) a segment was being timed and the acknowledgment number is greater than the starting sequence number of the segment being timed, `tcp_xmit_timer` updates the RTT estimators. Notice that the second argument to this function when timestamps are used is the current time (`tcp_now`) minus the timestamp echo reply (`ts_echo`) plus 1 (since the function subtracts 1).

Delayed ACKs are the reason for the greater-than test of the sequence numbers. For example, if TCP sends and times a segment with bytes 1-1024, followed by a segment with bytes 1025-2048, if an ACK of 2049 is returned, this test will consider whether 2049 is greater than 1 (the starting sequence number of the segment being timed), and since this is true, the RTT estimators are updated.

Check if all outstanding data has been acknowledged

916-924 If the acknowledgment field of the received segment (`ti_ack`) equals the maximum sequence number that TCP has sent (`snd_max`), all outstanding data has been acknowledged. The retransmission timer is turned off and the `needoutput` flag is set to 1. This flag forces a call to `tcp_output` at the end of this function. Since there is no more data waiting to be acknowledged, TCP may have more data to send that it has not been able to send earlier because the data was beyond the right edge of the window. Now that a new ACK has been received, the window will probably move to the right (`snd_una` is updated in Figure 29.8), which could allow more data to be sent.

```

903      /*
904      * If we have a timestamp reply, update smoothed
905      * round-trip time. If no timestamp is present but
906      * transmit timer is running and timed sequence
907      * number was acked, update smoothed round-trip time.
908      * Since we now have an rtt measurement, cancel the
909      * timer backoff (cf., Phil Karn's retransmit alg.).
910      * Recompute the initial retransmit timer,
911      */
912      if (ts_present)
913          tcp_xmit_timer(tp, tcp_now - ts_eck + 1);
914      else if (tp->t_rtt && SEQ_GT(ti->ti_ack, tp->t_rtsseq))
915          tcp_xmit_timer(tp, tp->t_rtt);
916      /*
917      * If all outstanding data is acked, stop retransmit
918      * timer and remember to restart (more output or persist).
919      * If there is more data to be acked, restart retransmit
920      * timer, using current (possibly backed-off) value
921      */
922      if (ti->ti_ack == tp->snd_max) {
923          tp->t_timer[TCPT_REXMT] = 0;
924          needoutput = 1;
925      } else if (tp->t_timer[TCPT_PERSIST] == 0)
926          tp->t_timer[TCPT_REXMT] = tp->t_rxtcur;

```

Figure 29.6 tcp_input function: RTT measurements and retransmission timer.

Unacknowledged data outstanding

925-926 Since there is additional data that has been sent but not acknowledged, if the persist timer is not on, the retransmission timer is restarted using the current value of `t_rxtcur`.

Karn's Algorithm and Timestamps

Notice that timestamps overrule the portion of Karn's algorithm (Section 21.3 of Volume 1) that says: when a timeout and retransmission occurs, the RTT estimators cannot be updated when the acknowledgment for the retransmitted data is received (the *retransmission ambiguity problem*). In Figure 25.26 we saw that `t_rtt` was set to 0 when a retransmission took place, because of Karn's algorithm. If timestamps are not present and it is a retransmission, the code in Figure 29.6 does not update the RTT estimators because `t_rtt` will be 0 from the retransmission. But if a timestamp is present, `t_rtt` isn't examined, allowing the RTT estimators to be updated using the received timestamp echo reply. With RFC 1323 timestamps the ambiguity is gone since the `ts_eck` value was copied by the other end from the segment being acknowledged. The other half of Karn's algorithm, specifying that an exponential backoff must be used with retransmissions, still holds, of course.

Figure 29.7 shows the next part of ACK processing, updating the congestion window.

```

927      /* tcp_input.c
928      * When new data is acked, open the congestion window.
929      * If the window gives us less than ssthresh packets
930      * in flight, open exponentially (maxseg per packet).
931      * Otherwise open linearly: maxseg per window
932      * (maxseg*2 / cwnd per packet), plus a constant
933      * fraction of a packet (maxseg/8) to help larger windows
934      * open quickly enough.
935      */
936      {
937          u_int   cw = tp->snd_cwnd;
938          u_int   incr = tp->t_maxseg;
939
940          if (cw > tp->snd_ssthresh)
941              incr = incr * incr / cw + incr / 8;
942          tp->snd_cwnd = min(cw + incr, TCP_MAXWIN << tp->snd_scale);
943      }
tcp_input.c

```

Figure 29.7 `tcp_input` function: open congestion window in response to ACKs.

Update congestion window

927-942 One of the rules of slow start and congestion avoidance is that a received ACK increases the congestion window. By default the congestion window is increased by one segment for each received ACK (slow start). But if the current congestion window is greater than the slow start threshold, it is increased by 1 divided by the congestion window, plus a constant fraction of a segment. The term

$$\text{incr} * \text{incr} / \text{cw}$$

is

$$t_maxseg * t_maxseg / \text{snd_cwnd}$$

which is 1 divided by the congestion window, taking into account that `snd_cwnd` is maintained in bytes, not segments. The constant fraction is the segment size divided by 8. The congestion window is then limited by the maximum value of the send window for this connection. Example calculations of this algorithm are in Section 21.8 of Volume 1.

Adding in the constant fraction (the segment size divided by 8) is wrong [Floyd 1994]. But it has been in the BSD sources since 4.3BSD Reno and is still in 4.4BSD and Net/3. It should be removed.

The next part of `tcp_input`, shown in Figure 29.8, removes the acknowledged data from the send buffer.

```

943     if (acked > so->so_snd.sb_cc) {                                     tcp_input.c
944         tp->snd_wnd -= so->so_snd.sb_cc;
945         sbdrop(&so->so_snd, (int) so->so_snd.sb_cc);
946         ourfinisacked = 1;
947     } else {
948         sbdrop(&so->so_snd, acked);
949         tp->snd_wnd -= acked;
950         ourfinisacked = 0;
951     }
952     if (so->so_snd.sb_flags & SB_NOTIFY)
953         sowakeup(so);
954     tp->snd_una = ti->tl_ack;
955     if (SEQ_LT(tp->snd_nxt, tp->snd_una))
956         tp->snd_nxt = tp->snd_una;

```

Figure 29.8 tcp_input function: remove acknowledged data from send buffer

Remove acknowledged bytes from the send buffer

943-946 If the number of bytes acknowledged exceeds the number of bytes on the send buffer, `snd_wnd` is decremented by the number of bytes in the send buffer and TCP knows that its FIN has been ACKed. That number of bytes is then removed from the send buffer by `sbdrop`. This method for detecting the ACK of a FIN works only because the FIN occupies 1 byte in the sequence number space.

947-951 Otherwise the number of bytes acknowledged is less than or equal to the number of bytes in the send buffer, so `ourfinisacked` is set to 0, and `acked` bytes of data are dropped from the send buffer.

Wakeup processes waiting on send buffer

952-956 `sowakeup` awakens any processes waiting on the send buffer. `snd_una` is updated to contain the oldest unacknowledged sequence number. If this new value of `snd_una` exceeds `snd_nxt`, the latter is updated, since the intervening bytes have been acknowledged.

Figure 29.9 shows how `snd_nxt` can end up with a sequence number that is less than `snd_una`. Assume two segments are transmitted, the first with bytes 1-512 and the second with bytes 513-1024.

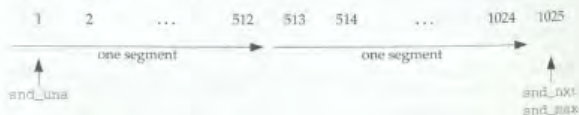


Figure 29.9 Two segments sent on a connection.

The retransmission timer then expires before an acknowledgment is returned. The code in Figure 25.26 sets `snd_nxt` back to `snd_una`, slow start is entered, `tcp_output` is called, and one segment containing bytes 1-512 is retransmitted. `tcp_output`

increases `snd_nxt` to 513, and we have the scenario shown in Figure 29.10.

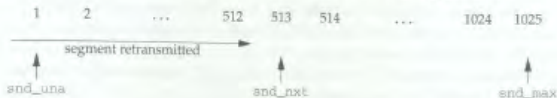


Figure 29.10 Continuation of Figure 29.9 after retransmission timer expires.

At this point an ACK of 1025 arrives (either the two original segments or the ACK was delayed somewhere in the network). The ACK is valid since it is less than or equal to `snd_max`, but `snd_nxt` will be less than the updated value of `snd_una`.

The general ACK processing is now complete, and the `switch` shown in Figure 29.11 handles four special cases.

```

957      switch (tp->t_state) {
958          /*
959           * In FIN_WAIT_1 state in addition to the processing
960           * for the ESTABLISHED state if our FIN is now acknowledged
961           * then enter FIN_WAIT_2.
962           */
963          case TCPS_FIN_WAIT_1:
964              if (ourfinisacked) {
965                  /*
966                   * If we can't receive any more
967                   * data, then closing user can proceed.
968                   * Starting the timer is contrary to the
969                   * specification, but if we don't get a FIN
970                   * we'll hang forever.
971                   */
972                  if (so->so_state & SS_CANTRCVMORE) {
973                      soisdisconnected(so);
974                      tp->t_timer[TCPT_2MSL] = tcp_maxidle;
975                  }
976                  tp->t_state = TCPS_FIN_WAIT_2;
977              }
978              break;

```

tcp_input.c

Figure 29.11 `tcp_input` function: receipt of ACK in `FIN_WAIT_1` state.

Receipt of ACK in `FIN_WAIT_1` state

958-977 In this state the process has closed the connection and TCP has sent the FIN. But other ACKs can be received for data segments sent before the FIN. Therefore the connection moves into the `FIN_WAIT_2` state only when the FIN has been acknowledged. The flag `ourfinisacked` is set in Figure 29.8; this depends on whether the number of bytes ACKed exceeds the amount of data in the send buffer or not.

Set FIN_WAIT_2 timer

972-975 We also described in Section 25.6 how Net/3 sets a FIN_WAIT_2 timer to prevent an infinite wait in the FIN_WAIT_2 state. This timer is set only if the process completely closed the connection (i.e., the `close` system call or its kernel equivalent if the process was terminated by a signal), and not if the process performed a half-close (i.e., the FIN was sent but the process can still receive data on the connection).

Figure 29.12 shows the receipt of an ACK in the CLOSING state.

```

979      /*
980      * In CLOSING state in addition to the processing for
981      * the ESTABLISHED state if the ACK acknowledges our FIN
982      * then enter the TIME-WAIT state, otherwise ignore
983      * the segment.
984      */
985      case TCPS_CLOSING:
986          if (ourfinisacked) {
987              tp->st_state = TCPS_TIME_WAIT;
988              tcp_canceltimers(tp);
989              tp->t_timer[TCPT_2MSL] = 2 * TCPTV_MSL;
990              goisdisconnected(so);
991          }
992          break;

```

Figure 29.12 `tcp_input` function: receipt of ACK in CLOSING state.

Receipt of ACK in CLOSING state

979-992 If the ACK is for the FIN (and not for some previous data segment), the connection moves into the TIME_WAIT state. Any pending timers are cleared (such as a pending retransmission timer), and the TIME_WAIT timer is started with a value of twice the MSL.

The processing of an ACK in the LAST_ACK state is shown in Figure 29.13.

```

993      /*
994      * In LAST_ACK, we may still be waiting for data to drain
995      * and/or to be acked, as well as for the ack of our FIN.
996      * If our FIN is now acknowledged, delete the TCB,
997      * enter the closed state, and return..
998      */
999      case TCPS_LAST_ACK:
1000         if (ourfinisacked) {
1001             tp = tcp_close(tp);
1002             goto drop;
1003         }
1004         break;

```

Figure 29.13 `tcp_input` function: receipt of ACK in LAST_ACK state.

Receipt of ACK in LAST_ACK state

893-1004 If the FIN is ACKed, the new state is CLOSED. This state transition is handled by `tcp_close`, which also releases the Internet PCB and TCP control block.

Figure 29.14 shows the processing of an ACK in the `TIME_WAIT` state.

```

1005          /*-----tcp_input.c
1006          * In TIME_WAIT state the only thing that should arrive
1007          * is a retransmission of the remote FIN. Acknowledge
1008          * it and restart the finack timer.
1009          */
1010          case TCPS_TIME_WAIT;
1011          tp->ti_timer[TCPT_3MSL] = 2 * TCPTV_MSL;
1012          goto dropafterack;
1013          }
1014      }
-----tcp_input.c

```

Figure 29.14 `tcp_input` function: receipt of ACK in `TIME_WAIT` state.

Receipt of ACK in TIME_WAIT state

7085-1014 In this state both ends have sent a FIN and both FINs have been acknowledged. If TCP's ACK of the remote FIN was lost, however, the other end will retransmit the FIN (with an ACK). TCP drops the segment and resends the ACK. Additionally, the `TIME_WAIT` timer must be restarted with a value of twice the MSL.

29.6 Update Window Information

There are two variables in the TCP control block that we haven't described yet: `snd_wl1` and `snd_wl2`.

- `snd_wl1` records the sequence number of the last segment used to update the send window (`snd_wnd`).
- `snd_wl2` records the acknowledgment number of the last segment used to update the send window.

Our only encounter with these variables so far was when a connection was established (active, passive, or simultaneous open) and `snd_wl1` was set to `ti_seq` minus 1. We said this was to guarantee a window update, which we'll see in the following code.

The send window (`snd_wnd`) is updated from the advertised window in the received segment (`tiwin`) if any one of the following three conditions is true:

1. The segment contains new data. Since `snd_wl1` contains the starting sequence number of the last segment that was used to update the send window, if

```
snd_wl1 < ti_seq
```

this condition is true.

- The segment does not contain new data (`snd_wll` equals `ti_seq`), but the segment acknowledges new data. The latter condition is true if

```
snd_wl2 < ti_ack
```

since `snd_wl2` records the acknowledgment number of the last segment that updated the send window.

- The segment does not contain new data, and the segment does not acknowledge new data, but the advertised window is larger than the current send window.

The purpose of these tests is to prevent an old segment from affecting the send window, since the send window is not an absolute sequence number, but is an offset from `snd_una`.

Figure 29.15 shows the code that implements the update of the send window.

```

1015 step6:                                     tcp_input.c
1016     /*
1017     * Update window information.
1018     * Don't look at window if an ACK: TAC's send garbage on first SYN.
1019     */
1020     if ((tiflags & TH_ACK) &&
1021         (SEQ_LT(tp->snd_wll, ti->ti_seq) || tp->snd_wll == ti->ti_seq &&
1022          (SEQ_LT(tp->snd_wl2, ti->ti_ack) ||
1023           tp->snd_wl2 == ti->ti_ack && tiwin > tp->snd_wnd))) {
1024         /* Keep track of pure window updates */
1025         if (ti->ti_len == 0 &&
1026             tp->snd_wl2 == ti->ti_ack && tiwin > tp->snd_wnd)
1027             tcpstat.tcps_rcvwinupd++;
1028
1029         tp->snd_wnd = tiwin;
1030         tp->snd_wll = ti->ti_seq;
1031         tp->snd_wl2 = ti->ti_ack;
1032         if (tp->snd_wnd > tp->max_sndwnd)
1033             tp->max_sndwnd = tp->snd_wnd;
1034         needoutput = 1;
1035     }

```

Figure 29.15 `tcp_input` function: update window information.

Check if send window should be updated

1015-1027 This if test verifies that the ACK flag is set along with any one of the three previously stated conditions. Recall that a jump was made to `step6` after the receipt of a SYN in either the LISTEN or SYN_SENT state, and in the LISTEN state the SYN does not contain an ACK.

The term TAC referred to in the comment is a "terminal access controller." These were Telet clients on the ARPANET.

1034-1037 If the received segment is a pure window update (the length is 0 and the ACK does not acknowledge new data, but the advertised window is larger), the statistic `tcps_rcvwinupd` is incremented.

Update variables

1038-1043 The send window is updated and new values of `snd_wll` and `snd_wl2` are recorded. Additionally, if this advertised window is the largest one TCP has received from this peer, the new value is recorded in `max_sndwnd`. This is an attempt to guess the size of the other end's receive buffer, and it is used in Figure 26.8. `needoutput` is set to 1 since the new value of `snd_wnd` might enable a segment to be sent.

29.7 Urgent Mode Processing

The next part of TCP input processing handles segments with the URG flag set.

```

1035      /*
1036      * Process segments with URG.
1037      */
1038      if ((tiflags & TB_URG) && ti->ti_urg &&
1039          TCPS_HAVERCVDFIN(tp->t_state) == 0) {
1040          /*
1041           * This is a kludge, but if we receive and accept
1042           * random urgent pointers, we'll crash in
1043           * soreceive. It's hard to imagine someone
1044           * actually wanting to send this much urgent data.
1045           */
1046          if (ti->ti_urg + so->so_rcv.sb_cc > sb_max) {
1047              ti->ti_urg = 0; /* XXX */
1048              tiflags &= ~TB_URG; /* XXX */
1049              goto dodata; /* XXX */
1050          }
1051      }

```

tcp_input.c

Figure 29.16 `tcp_input` function: urgent mode processing.

Check if URG flag should be processed

1035-1039 These segments must have the URG flag set, a nonzero urgent offset (`ti_urg`), and the connection must not have received a FIN. The macro `TCPS_HAVERCVDFIN` is true only for the `TIME_WAIT` state, so the URG is processed in any other state. This is contrary to a comment appearing later in the code stating that the URG flag is ignored in the `CLOSE_WAIT`, `CLOSING`, `LAST_ACK`, or `TIME_WAIT` states.

Ignore bogus urgent offsets

1046-1050 If the urgent offset plus the number of bytes already in the receive buffer exceeds the maximum size of a socket buffer, the urgent notification is ignored. The urgent offset is set to 0, the URG flag is cleared, and the rest of the urgent mode processing is skipped.

The next piece of code, shown in Figure 29.17, processes the urgent pointer.

```

1051      /*
1052      * If this segment advances the known urgent pointer,
1053      * then mark the data stream. This should not happen
1054      * in CLOSE_WAIT, CLOSING, LAST_ACK or TIME_WAIT states since
1055      * a FIN has been received from the remote side.
1056      * In these states we ignore the URG.
1057      *
1058      * According to RFC961 (Assigned Protocols),
1059      * the urgent pointer points to the last octet
1060      * of urgent data. We continue, however,
1061      * to consider it to indicate the first octet
1062      * of data past the urgent section as the original
1063      * spec states (in one of two places).
1064      */
1065      if (SEQ_GT(ti->ti_seq + ti->ti_urg, tp->rcv_up)) {
1066          tp->rcv_up = ti->ti_seq + ti->ti_urg;
1067          so->so_oobmark = so->so_rcv.sb_cc +
1068              (tp->rcv_up - tp->rcv_next) - 1;
1069          if (so->so_oobmark == 0)
1070              so->so_state |= SS_RCVATMARK;
1071          sohasoutofband(so);
1072          tp->t_oobflags &= ~(TCPOOB_HAVEDATA | TCPOOB_HADDATA);
1073      }
1074      /*
1075      * Remove out-of-band data so doesn't get presented to user.
1076      * This can happen independent of advancing the URG pointer,
1077      * but if two URG's are pending at once, some out-of-band
1078      * data may creep in... ick.
1079      */
1080      if (ti->ti_urg <= ti->ti_len
1081      #ifdef SO_OOBINLINE
1082          && (so->so_options & SO_OOBINLINE) == 0
1083      #endif
1084          )
1085          tcp_pulloutofband(so, ti, m);
1086      } else {
1087          /*
1088          * If no out-of-band data is expected, pull receive
1089          * urgent pointer along with the receive window.
1090          */
1091          if (SEQ_GT(tp->rcv_next, tp->rcv_up))
1092              tp->rcv_up = tp->rcv_next;
1093      }

```

Figure 29.17 tcp_input function: processing of received urgent pointer.

1051-1065 If the starting sequence number of the received segment plus the urgent offset exceeds the current receive urgent pointer, a new urgent pointer has been received. For example, when the 3-byte segment that was sent in Figure 26.30 arrives at the receiver, we have the scenario shown in Figure 29.18.

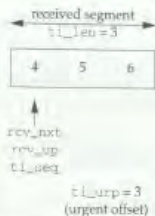


Figure 29.18 Receiver side when segment from Figure 26.30 arrives.

Normally the receive urgent pointer (`rcv_up`) equals `rcv_nxt`. In this example, since the `if` test is true (4 plus 3 is greater than 4), the new value of `rcv_up` is calculated as 7.

Calculate receive urgent pointer

1056-1070 The out-of-band mark in the socket's receive buffer is calculated, taking into account any data bytes already in the receive buffer (`so_rcv.sb_off`). In our example, assuming there is no data already in the receive buffer, `so_oobmark` is set to 2: that is, the byte with the sequence number 6 is considered the out-of-band byte. If this out-of-band mark is 0, the socket is currently at the out-of-band mark. This happens if the `send` system call that sends the out-of-band byte specifies a length of 1, and if the receive buffer is empty when this segment arrives at the other end. This reiterates that Berkeley-derived systems consider the urgent pointer to point to the first byte of data after the out-of-band byte.

Notify process of TCP's urgent mode

1071-1072 `sohasoutofband` notifies the process that out-of-band data has arrived for the socket. The two flags `TCPOOB_HAVEDATA` and `TCPOOB_HADDATA` are cleared. These two flags are used with the `PRU_RCVOOB` request in Figure 30.8.

Pull out-of-band byte out of normal data stream

1074-1085 If the urgent offset is less than or equal to the number of bytes in the received segment, the out-of-band byte is contained in the segment. With TCP's urgent mode it is possible for the urgent offset to point to a data byte that has not yet been received. If the `SO_OOBINLINE` constant is defined (which it always is for Net/3), and if the corresponding socket option is not enabled, the receiving process wants the out-of-band byte pulled out of the normal stream of data and placed into the variable `t_ioobc`. This is done by `tcp_pulloutofband`, which we cover in the next section.

Notice that the receiving process is notified that the sender has entered urgent mode, regardless of whether the byte pointed to by the urgent pointer is readable or not. This is a feature of TCP's urgent mode.

Adjust receive urgent pointer if not urgent mode

1086-1093 When the receiver is not processing an urgent pointer, if `rcv_nxt` is greater than the receive urgent pointer, `rcv_up` is moved to the right and set equal to `rcv_nxt`. This keeps the receive urgent pointer at the left edge of the receive window so that the

comparison using `SEQ_GT` at the beginning of Figure 29.17 will work correctly when an URG flag is received.

If the solution to Exercise 26.6 is implemented, corresponding changes will have to go into Figures 29.16 and 29.17 also.

29.8 `tcp_pulloutofband` Function

This function is called from Figure 29.17 when

1. urgent mode notification arrives in a received segment, and
2. the out-of-band byte is contained within the segment (i.e., the urgent pointer points into the received segment), and
3. the `SO_OOBINLINE` socket option is not enabled for this socket.

This function removes the out-of-band byte from the normal stream of data (i.e., the mbuf chain containing the received segment) and places it into the `t_iobc` variable in the TCP control block for the connection. The process reads this variable using the `MSG_OOB` flag with the `recv` system call: the `PRU_RCVOOB` request in Figure 30.8. Figure 29.19 shows the function.

```

1282 void
1283 tcp_pulloutofband(so, ti, m)
1284 struct socket *so;
1285 struct tcpiphdr *ti;
1286 struct mbuf *m;
1287 {
1288     int     cnt = ti->ti_urp - 1;
1289     while (cnt >= 0) {
1290         if (m->m_len > cnt) {
1291             char *cp = mtod(m, caddr_t) + cnt;
1292             struct tcpcb *tp = sototcpcb(so);
1293             tp->t_iobc = *cp;
1294             tp->t_oobflags |= TCPOOB_HAVEDATA;
1295             bcopy(cp + 1, cp, (unsigned) (m->m_len - cnt - 1));
1296             m->m_len--;
1297             return;
1298         }
1299         cnt -= m->m_len;
1300         m = m->m_next;
1301         if (m == 0)
1302             break;
1303     }
1304     panic("tcp_pulloutofband");
1305 }

```

Figure 29.19 `tcp_pulloutofband` function: place out-of-band byte into `t_iobc`.

1282-1289

Consider the example in Figure 29.20. The urgent offset is 3, therefore the urgent pointer is 7, and the sequence number of the out-of-band byte is 6. There are 5 bytes in the received segment, all contained in a single mbuf.

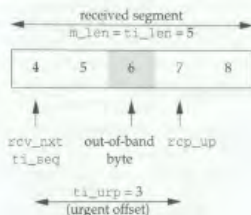


Figure 29.20 Received segment with an out-of-band byte.

The variable `cnt` is 2 and since `m_len` (which is 5) is greater than 2, the true portion of the `if` statement is executed.

1290-1298

`cp` points to the shaded byte with a sequence number of 6. This is placed into the variable `t_iobc`, which contains the out-of-band byte. The `TCPOOB_HAVEDATA` flag is set and `bcopy` moves the next 2 bytes (with sequence numbers 7 and 8) left 1 byte, giving the arrangement shown in Figure 29.21.

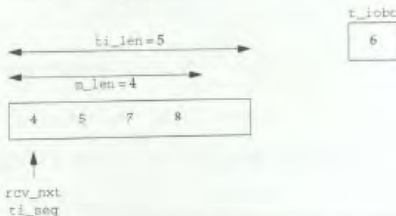


Figure 29.21 Result from Figure 29.20 after removal of out-of-band byte.

Remember that the numbers 7 and 8 specify the sequence numbers of the data bytes, not the contents of the data bytes. The length of the mbuf is decremented from 5 to 4 but `ti_len` is left as 5, for sequencing of the segment into the socket's receive buffer. Both the `TCP_REASS` macro and the `tcp_reass` function (which are called in the next section) increment `rcv_nxt` by `ti_len`, which in this example must be 5, because the next expected receive sequence number is 9. Also notice in this function that the length field in the packet header (`m_pkthdr.len`) in the first mbuf is not decremented by 1. This is because that length field is not used by `sbandend`, which appends the data to the socket's receive buffer.

Skip to next mbuf in chain

2299-2303 The out-of-band byte is not contained in this mbuf, so `cnt` is decremented by the number of bytes in the mbuf and the next mbuf in the chain is processed. Since this function is called only when the urgent offset points into the received segment, if there is not another mbuf on the chain, the `break` causes the call to `panic`.

29.9 Processing of Received Data

`tcp_input` continues by taking the received data (if any) and either appending it to the socket's receive buffer (if it is the next expected segment) or placing it onto the socket's out-of-order queue. Figure 29.22 shows the code that performs this task.

```

1094 sockData; /* XXX */
1095 /*
1096  * Process the segment text, merging it into the TCP sequencing queue,
1097  * and arranging for acknowledgment of receipt if necessary.
1098  * This process logically involves adjusting tp->rcv_wnd as data
1099  * is presented to the user (this happens in tcp_usrreq.c,
1100  * case PRU_RCVD). If a FIN has already been received on this
1101  * connection then we just ignore the text.
1102  */
1103 if ((ti->ti_len || (tiflags & TH_FIN)) &&
1104     TCPS_HAVERCVDFIN(tp->t_state) == 0) {
1105     TCP_REASS(tp, ti, m, so, tiflags);
1106     /*
1107      * Note the amount of data that peer has sent into
1108      * our window, in order to estimate the sender's
1109      * buffer size.
1110      */
1111     len = so->so_rcv.sb_hiwat - (tp->rcv_adv - tp->rcv_next);
1112 } else {
1113     m_freem(m);
1114     tiflags &= ~TH_FIN;
1115 }

```

tcp_input.c

Figure 29.22 `tcp_input` function: merge received data into sequencing queue for socket

2304-2305 Segment data is processed if

1. the length of the received data is greater than 0 or the FIN flag is set, and
2. a FIN has not yet been received for the connection.

The macro `TCP_REASS` processes the data. If the data is in sequence (i.e., the next expected data for this connection), the delayed-ACK flag is set, `rcv_next` is incremented, and the data is appended to the socket's receive buffer. If the data is out of order, the macro calls `tcp_reass` to add the data to the connection's reassembly queue (which might fill a hole and cause already-queued data to be appended to the socket's receive buffer).

Recall that the final argument to the macro (`tiflags`) can be modified. Specifically, if the data is out of order, `tcp_reass` sets `tiflags` to 0, clearing the FIN flag (if it was set). That's why the `if` statement is true if the FIN flag is set even if there is no data in the segment.

Consider the following example. A connection is established and the sender immediately transmits three segments: one with bytes 1–1024, another with bytes 1025–2048, and another with the FIN flag but no data. The first segment is lost, so when the second arrives (bytes 1025–2048) the receiver places it onto the out-of-order list and generates an immediate ACK. When the third segment with the FIN flag is received, the code in Figure 29.22 is executed. Even though the data length is 0, since the FIN flag is set, `TCP_REASS` is invoked, which calls `tcp_reass`. Since `ti_seq` (2049), the sequence number of the FIN) does not equal `rcv_nxt` (1), `tcp_reass` returns 0 (Figure 27.23), which in the `TCP_REASS` macro sets `tiflags` to 0. This clears the FIN flag, preventing the code that follows (Section 29.10) from processing the FIN flag.

Guess size of other end's send buffer

1106–1111 The calculation of `len` is attempt to guess the size of the other end's send buffer. Consider the following example. A socket has a receive buffer size of 8192 (the Net/3 default), so TCP advertises a window of 8192 in its SYN. The first segment with bytes 1–1024 is then received. Figure 29.23 shows the state of the receive space after `TCP_REASS` has incremented `rcv_nxt` to account for the received segment.

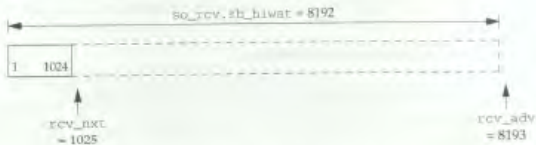


Figure 29.23 Receipt of bytes 1–1024 into a 8192-byte receive window.

The calculation of `len` yields 1024. The value of `len` will increase as the other end sends more data into the receive window, but it will never exceed the size of the other end's send buffer. Recall that the variable `max_sndwnd`, calculated in Figure 29.15, is an attempt to guess the size of the other end's receive buffer.

This variable `len` is never used! It is left over code from Net/1 when the variable `max_rcvwd` was stored in the TCP control block after the calculation of `len`:

```
if (len > tp->max_rcvwd)
    tp->max_rcvwd = len;
```

But even in Net/1 the variable `max_rcvwd` was never used.

1112–1115 If the length is 0 and the FIN flag is not set, or if a FIN has already been received for the connection, the received mbuf chain is discarded and the FIN flag is cleared.

29.10 FIN Processing

The next step in `tcp_input`, shown in Figure 29.24, handles the FIN flag.

```

1116  /*
1117  * If FIN is received ACK the FIN and let the user know
1118  * that the connection is closing.
1119  */
1120  if (tflags & TF_FIN) {
1121      if (TCPS_HAVERCVDPIN(tp->t_state) == 0) {
1122          socantrcvmore(s0);
1123          tp->t_flags |= TF_ACKNOW;
1124          tp->rcv_nxt++;
1125      }
1126      switch (tp->t_state) {
1127          /*
1128           * In SYN_RECEIVED and ESTABLISHED states
1129           * enter the CLOSE_WAIT state.
1130           */
1131          case TCPS_SYN_RECEIVED:
1132          case TCPS_ESTABLISHED:
1133              tp->t_state = TCPS_CLOSE_WAIT;
1134              break;

```

Figure 29.24 `tcp_input` function: FIN processing, first half.

Process first FIN received on connection

If the FIN flag is set and this is the first FIN received for this connection, `socantrcvmore` marks the socket as write-only, `TF_ACKNOW` is set to acknowledge the FIN immediately (i.e., it is not delayed), and `rcv_nxt` steps over the FIN in the sequence space.

The remainder of FIN processing is handled by a `switch` that depends on the connection state. Notice that the FIN is not processed in the `CLOSED`, `LISTEN`, or `SYN_SENT` states, since in these three states a SYN has not been received to synchronize the received sequence number, making it impossible to validate the sequence number of the FIN. A FIN is also ignored in the `CLOSING`, `CLOSE_WAIT`, and `LAST_ACK` states, because in these three states the FIN is a duplicate.

SYN_RCVD or ESTABLISHED states

From either the `ESTABLISHED` or `SYN_RCVD` states, the `CLOSE_WAIT` state is entered.

The receipt of a FIN in the `SYN_RCVD` state is unusual, but legal. It is not shown in Figure 24.15. It means a socket is in the `LISTEN` state when a segment containing a SYN and a FIN is received. Alternatively, a SYN is received for a listening socket, moving the connection to the `SYN_RCVD` state but before the ACK is received. (We know the segment does not contain a valid ACK, because if it did the code in Figure 29.2 would have moved the connection to the `ESTABLISHED` state.)

The next part of FIN processing is shown in Figure 29.25

```

1135      /*
1136      * If still in FIN_WAIT_1 state FIN has not been acked so
1137      * enter the CLOSING state.
1138      */
1139      case TCPS_FIN_WAIT_1:
1140          tp->t_state = TCPS_CLOSING;
1141          break;
1142      /*
1143      * In FIN_WAIT_2 state enter the TIME_WAIT state,
1144      * starting the time-wait timer, turning off the other
1145      * standard timers.
1146      */
1147      case TCPS_FIN_WAIT_2:
1148          tp->t_state = TCPS_TIME_WAIT;
1149          tcp_canceltimers(tp);
1150          tp->t_timer[TCPT_2MSL] = 2 * TCPTV_MSL;
1151          so_ladisonconnected(so);
1152          break;
1153      /*
1154      * In TIME_WAIT state restart the 2 MSL time_wait timer.
1155      */
1156      case TCPS_TIME_WAIT:
1157          tp->t_timer[TCPT_2MSL] = 2 * TCPTV_MSL;
1158          break;
1159      }
1160 }

```

tcp_input.c

Figure 29.25 tcp_input function: FIN processing, second half.

FIN_WAIT_1 state

1135-1141 Since ACK processing is already complete for this segment, if the connection is in the FIN_WAIT_1 state when the FIN is processed, it means a simultaneous close is taking place—the two FINs from each end have passed in the network. The connection enters the CLOSING state.

FIN_WAIT_2 state

1142-1148 The receipt of the FIN moves the connection into the TIME_WAIT state. When a segment containing a FIN and an ACK is received in the FIN_WAIT_1 state (the typical scenario), although Figure 24.15 shows the transition directly from the FIN_WAIT_1 state to the TIME_WAIT state, the ACK is processed in Figure 29.11, moving the connection to the FIN_WAIT_2 state. The FIN processing here moves the connection into the TIME_WAIT state. Because the ACK is processed before the FIN, the FIN_WAIT_2 state is always passed through, albeit momentarily.

Start TIME_WAIT timer

1149-1152 Any pending TCP timer is turned off and the TIME_WAIT timer is started with a value of twice the MSL. (If the received segment contained a FIN and an ACK, Figure 29.11 started the FIN_WAIT_2 timer.) The socket is disconnected.

TIME_WAIT state

1153-1159 If a FIN arrives in the TIME_WAIT state, it is a duplicate, and similar to Figure 29.14, the TIME_WAIT timer is restarted with a value of twice the MSL.

29.11 Final Processing

The final part of the slow path through `tcp_input` along with the label `dropafterack` is shown in Figure 29.26.

```

1161     if (so->so_options & SO_DEBUG)                                     tcp_input.c
1162         tcp_trace(TA_INPUT, ostate, tp, &tcp_savet, 0);
1163     /*
1164      * Return any desired output.
1165      */
1166     if (needoutput || (tp->t_flags & TF_ACKNOW))
1167         (void) tcp_output(tp);
1168     return;
1169 dropafterack:
1170     /*
1171      * Generate an ACK dropping incoming segment if it occupies
1172      * sequence space, where the ACK reflects our state.
1173      */
1174     if (tiflags & TH_RST)
1175         goto drop;
1176     m_freem(m);
1177     tp->t_flags |= TF_ACKNOW;
1178     (void) tcp_output(tp);
1179     return;

```

Figure 29.26 `tcp_input` function: final processing.

SO_DEBUG socket option

1161-1162 If the SO_DEBUG socket option is enabled, `tcp_trace` appends the trace record to the kernel's circular buffer. Remember that the code in Figure 28.7 saved both the original connection state and the IP and TCP headers, since these values may have changed in this function.

Call tcp_output

1163-1168 If either the `needoutput` flag was set (Figures 29.6 and 29.15) or if an immediate ACK is required, `tcp_output` is called.

dropafterack

1169-1179 An ACK is generated only if the RST flag was not set. (A segment with an RST is never ACKed.) The mbuf chain containing the received segment is released, and `tcp_output` generates an immediate ACK.

Figure 29.27 completes the `tcp_input` function.

```

1180 dropwithreset;                                     tcp_input.c
1181 /*
1182  * Generate an RST, dropping incoming segment.
1183  * Make ACK acceptable to originator of segment.
1184  * Don't bother to respond if destination was broadcast/multicast.
1185  */
1186 if ((tiflags & TH_RST) || m->m_flags & (M_BCAST | M_MCAST) ||
1187     IN_MULTICAST(ti->ti_dst.s_addr))
1188     goto drop;
1189 if (tiflags & TH_ACK)
1190     tcp_respond(tp, ti, m, (tcp_seq) 0, ti->ti_ack, TH_RST);
1191 else {
1192     if (tiflags & TH_SYN)
1193         ti->ti_len++;
1194     tcp_respond(tp, ti, m, ti->ti_seq + ti->ti_len, (tcp_seq) 0,
1195                 TH_RST | TH_ACK);
1196 }
1197 /* destroy temporarily created socket */
1198 if (dropsocket)
1199     (void) sobort(so);
1200 return;
1201 drop:
1202 /*
1203  * Drop space held by incoming segment and return.
1204  */
1205 if (tp && (tp->t_inpcb->inp_socket->so_options & SO_DEBUG){
1206     tcp_trace(TA_DROP, ostate, tp, &tcp_saveti, 0);
1207     m_freem(m);
1208     /* destroy temporarily created socket */
1209     if (dropsocket)
1210         (void) sobort(so);
1211     return;
1212 }

```

Figure 29.27 tcp_input function: final processing.

dropwithreset

1180-1188 An RST is generated unless the received segment also contained an RST, or the received segment was sent as a broadcast or multicast. An RST is never generated in response to an RST, since this could lead to RST storms (a continual exchange of RST segments between two end points).

This code contains the same error that we noted in Figure 28.16: it does not check whether the destination address of the received segment was a broadcast address.

Similarly, the destination address argument to `IN_MULTICAST` needs to be converted to host byte order.

Sequence number and acknowledgment number of RST segment

1189-1196 The values of the sequence number field, the acknowledgment field, and the ACK flag of the RST segment depend on whether the received segment contained an ACK.

Figure 29.28 summarizes these fields in the RST segment that is generated.

received segment	RST segment generated		
	seq#	ack. field	flags
contains ACK	received ack. field	0	TH_RST
ACK-less	0	received seq# field	TH_RST TH_ACK

Figure 29.28 Values of fields in RST segment generated.

Realize that the ACK flag is normally set in all segments except when an initial SYN is sent (Figure 24.16). The fourth argument to `tcp_respond` is the acknowledgment field, and the fifth argument is the sequence number.

Rejecting connections

If the SYN flag is set, `ti_len` must be incremented by 1, causing the acknowledgment field of the RST to be 1 greater than the received sequence number of the SYN. This code is executed when a SYN arrives for a nonexistent server. When the Internet PCB is not found in Figure 28.6, a jump is made to `dropwithreset`. But for the received RST to be acceptable to the other end, the acknowledgment field must ACK the SYN (Figure 28.18). Figure 18.14 of Volume 1 contains an example of this type of RST segment.

Finally note that `tcp_respond` builds the RST in the first mbuf of the received chain and releases any remaining mbufs in the chain. When that mbuf finally makes its way to the device driver, it will be discarded.

Destroy temporarily created socket

If a temporary socket was created in Figure 28.7 for a listening server, but the code in Figure 28.16 found the received segment to contain an error, `dropsocket` will be 1. If so, that socket is now destroyed.

Drop (without ACK or RST)

`tcp_trace` is called when a segment is dropped without generating an ACK or an RST. If the `SO_DEBUG` flag is set and an ACK is generated, `tcp_output` generates a trace record. If the `SO_DEBUG` flag is set and an RST is generated, a trace record is not generated for the RST.

The mbuf chain containing the received segment is released and the temporary socket is destroyed if `dropsocket` is nonzero.

29.12 Implementation Refinements

The refinements to speed up TCP processing are similar to the ones described for UDP (Section 23.12). Multiple passes over the data should be avoided and the checksum computation should be combined with a copy. [Dalton et al. 1993] describe these modifications.

The linear search of the TCP PCBs is also a bottleneck when the number of connections increases. [McKenney and Dove 1992] address this problem by replacing the linear search with hash tables.

[Partridge 1993] describes a research implementation being developed by Van Jacobson that greatly reduces the TCP input processing. The received packet is processed by IP (about 25 instructions on a RISC system), then by a demultiplexer to locate the PCB (about 10 instructions), and then by TCP (about 30 instructions). These 30 instructions perform header prediction and calculate the pseudo-header checksum. If the segment passes the header prediction test, contains data, and the process is waiting for the data, the data is copied into the process buffer and the remainder of the TCP checksum is calculated and verified (a one-pass copy and checksum). If the TCP header prediction fails, the slow path through the TCP input processing occurs.

29.13 Header Compression

We now describe TCP *header compression*. Although header compression is not part of TCP input, we needed to cover TCP thoroughly before describing header compression. Header compression is described in detail in RFC 1144 [Jacobson 1990a]. It was designed by Van Jacobson and is sometimes called *VJ header compression*. Our purpose in this section is not to go through the header compression source code (a well-commented version of which is presented in RFC 1144, and which is approximately the same size as `tcp_output`), but to provide an overview of the algorithm. Be sure to distinguish between header prediction (Section 28.4) and header compression.

Introduction

Most implementations of SLIP and PPP support header compression. Although header compression could, in theory, be used with any data link, it is intended for slow-speed serial links. Header compression works with TCP segments only—it does nothing with other IP datagrams (e.g., ICMP, IGMP, UDP, etc.). Header compression reduces the size of the combined IP/TCP header from its normal 40 bytes to as few as 3 bytes. This reduces the size of a typical TCP segment from an interactive application such as Rlogin or Telnet from 41 bytes to 4 bytes—a big saving on a slow-speed serial link.

Each end of the serial link maintains two connection state tables, one for datagrams sent and one for datagrams received. Each table allows a maximum of 256 entries, but typically there are 16 entries in this table, allowing up to 16 different TCP connections to be compressed at any time. Each entry contains an 8-bit connection ID (hence the limit of 256), some flags, and the complete uncompressed IP/TCP header from the most recent datagram. The 96-bit socket pair that uniquely identifies each connection—the source and destination IP addresses and source and destination TCP ports—are contained in this uncompressed header. Figure 29.29 shows an example of these tables.

Since a TCP connection is full duplex, header compression can be applied in both directions. Each end must implement both compression and decompression. A connection appears in both tables, as shown in Figure 29.29. In this example, the entry with a connection ID of 1 in the top two tables has a source IP address of 128.1.2.3, source TCP port of 1500, destination IP address of 192.3.4.5, and a destination TCP port of 25. The entry with a connection ID of 2 in the bottom two tables is for the other direction of the same connection.

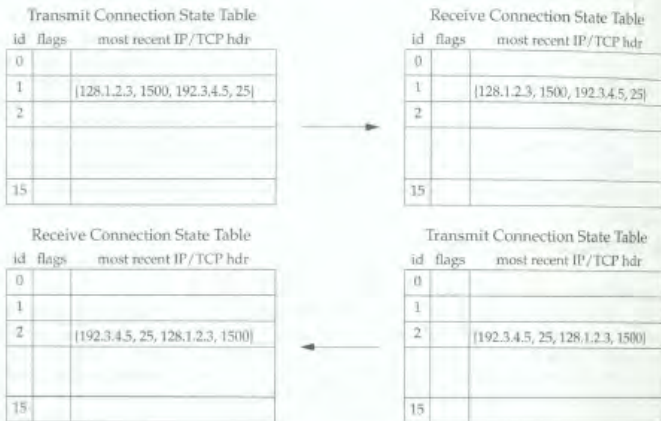


Figure 29.29 A pair of connection state tables at each end of a link (e.g., SLIP link).

We show these tables as arrays, but the source code defines each entry as a structure, and a connection table is a circular linked list of these structures. The most recently used structure is stored at the head of the list.

By saving the most recent uncompressed header at each end, only the *differences* in various header fields from the previous datagram to the current datagram are transmitted across the link (along with a special first byte indicating which fields follow). Since some header fields don't change at all from one datagram to the next, and other header fields change by small amounts, this differential coding provides the savings. Header compression works with the IP and TCP headers only—the data contents of the TCP segment are not modified.

Figure 29.30 shows the steps involved at the sending side when it has an IP datagram to send across a link using header compression.

Three different types of datagrams are sent and must be recognized at the receiver:

1. Type `IP` is specified with the high-order 4 bits of the first byte equal to 4. This is the normal IP version number in the IP header (Figure 8.8). The normal, uncompressed datagram is transmitted across the link.
2. Type `COMPRESSED_TCP` is specified by setting the high-order bit of the first byte. This looks like an IP version between 8 and 15 (i.e., the remaining 7 bits of this byte are used by the compression algorithm). The compressed header and uncompressed data are transmitted across the link, as we describe later in this section.

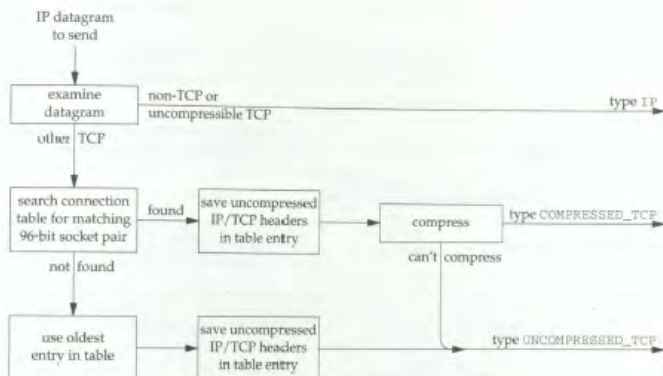


Figure 29.30 Steps involved in header compression at sender side.

3. Type `UNCOMPRESSED_TCP` is specified with the high-order 4 bits of the first byte equal to 7. The normal, uncompressed datagram is transmitted across the link, but the IP protocol field (which equals 6 for TCP), is replaced with the connection ID. This identifies the connection state table entry for the receiver.

The receiver can identify the datagram type by examining its first byte. The code that does this was shown in Figure 5.13. In Figure 5.16 the sender calls `sl_compress_tcp` to check if a TCP segment is compressible, and the return value of this function is logically ORed into the first byte of the datagram.

Figure 29.31 shows an illustration of the first byte that is sent across the link.

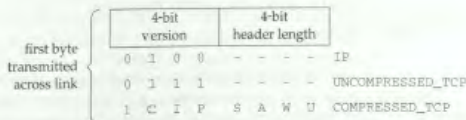


Figure 29.31 First byte transmitted across link.

The 4 bits shown as “-” comprise the normal IP header length field. The 7 bits shown as C, I, P, S, A, W, and U indicate which optional fields follow. We describe these fields shortly.

Figure 29.32 shows the complete IP datagram for the various datagrams that are sent.

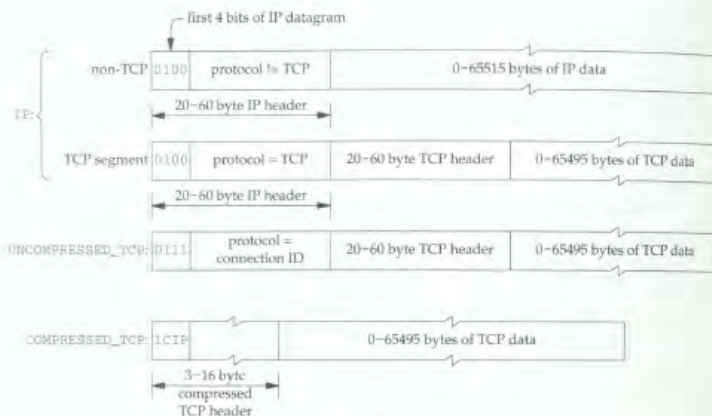


Figure 29.32 Different types of IP datagrams possible with header compression.

We show two datagrams with a type of IP: one that is not a TCP segment (e.g., a protocol of UDP, ICMP, or IGMP), and one that is a TCP segment. This is to illustrate the differences between the TCP segment sent as type IP and the TCP segment sent as type UNCOMPRESSED_TCP: the first 4 bits are different as is the protocol field in the IP header.

Datagrams are not candidates for header compression if the protocol is not TCP, or if the protocol is TCP but any one of the following conditions is true.

- The datagram is an IP fragment: either the fragment offset is nonzero or the more-fragments bit is set.
- Any one of the SYN, FIN, or RST flags is set.
- The ACK flag is not set.

If any one of these three conditions is true, the datagram is sent as type IP.

Furthermore, even if the datagram is a TCP segment that looks compressible, it is possible to abort the compression and send the datagram as type UNCOMPRESSED_TCP if certain fields have changed between the current datagram and the last datagram sent for this connection. These are fields that normally do not change for a given connection, so the compression scheme was not designed to encode their differences from one datagram to the next. The TOS field and the don't fragment bit are examples. Also, when the differences in some fields are greater than 65535, the compression algorithm fails and the datagram is sent uncompressed.

Compression of Header Fields

We now describe how the fields in the IP and TCP headers, shown in Figure 29.33, are compressed. The shaded fields normally don't change during a connection.

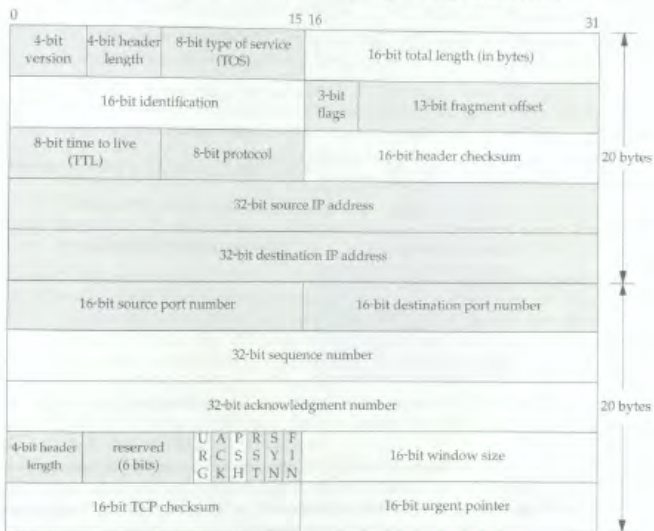


Figure 29.33 Combined IP and TCP headers: shaded fields normally don't change.

If any of the shaded fields have changed from the previous segment on this connection to the current segment, the segment is sent uncompressed. We don't show IP options or TCP options in this figure, but if either are present and have changed from the previous segment, the segment is sent uncompressed (Exercise 29.7).

If the algorithm transmitted only the nonshaded fields when the shaded fields do not change from the previous segment, about a 50% savings would result. VJ header compression does even better than this, by knowing which fields in the IP and TCP headers *normally* don't change. Figure 29.34 shows the format of the compressed IP/TCP header.

The smallest compressed header consists of 3 bytes: the first byte (the flag bits) followed by the 16-bit TCP checksum. For protection against possible link errors, the TCP checksum is always transmitted without any change. (SLIP provides no link-layer checksum, although PPP does provide one.)

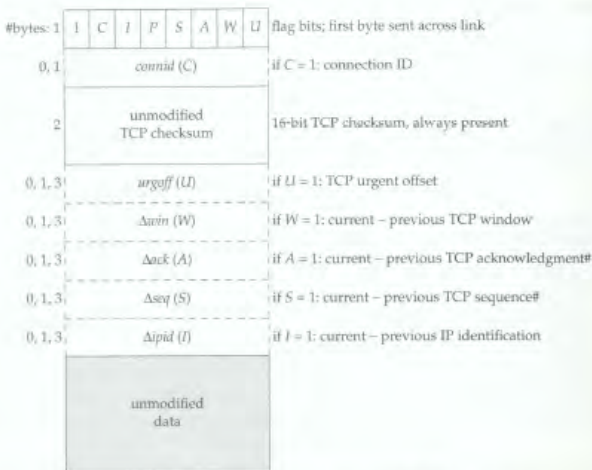


Figure 29.34 Format of compressed IP/TCP header.

The other six fields, *connid*, *urgoff*, *Δwin* , *Δack* , *Δseq* , and *Δpid* , are optional. We show the number of bytes used to encode all the fields to the left of the field in Figure 29.34. The largest compressed header appears to be 19 bytes, but we'll see shortly that the 4 bits *SAWU* can never be set at the same time in a compressed header, so the largest size is actually 16 bytes.

Six of the 7 bits in the first byte specify which of the six optional fields are present. The high-order bit of the first byte is always set to 1. This identifies the datagram type as *COMPRESSED_TCP*. Figure 29.35 summarizes the 7 bits, which we now describe.

Flag bit	Description	Structure member	Meaning if flag = 0	Meaning if flag = 1
C	connection ID		same connection ID as last	<i>connid</i> = connection ID
I	IP identification	<i>ip_id</i>	<i>ip_id</i> has increased by 1	<i>Δpid</i> = current - previous
P	TCP push flag		PSH flag off	PSH flag on
S	TCP sequence#	<i>th_seq</i>	same <i>th_seq</i> as last	<i>Δseq</i> = current - previous
A	TCP acknowledgment#	<i>th_ack</i>	same <i>th_ack</i> as last	<i>Δack</i> = current - previous
W	TCP window	<i>th_win</i>	same <i>th_win</i> as last	<i>Δwin</i> = current - previous
U	TCP urgent offset	<i>th_urg</i>	URG flag not set	<i>urgoff</i> = urgent offset

Figure 29.35 The 7 bits in the compressed header.

- C If this bit is 0, this segment has the same connection ID as the previous compressed or uncompressed segment. If this flag is 1, *connid* is the connection ID, a value between 0 and 255.
- I If this bit is 0, the IP identification field has increased by 1 (the typical case). If this bit is 1, *Δipid* is the current value of *ip_id* minus its previous value.
- P This bit is a copy of the PSH flag from the TCP segment. Since the PSH flag doesn't follow any established pattern, it must be explicitly specified for each segment.
- S If this bit is 0, the TCP sequence number has not changed. If this bit is 1, *Δseq* is the current value of *th_seq* minus its previous value.
- A If this bit is 0, the TCP acknowledgment number has not changed (the typical case). If this bit is 1, *Δack* is the current value of *th_ack* minus its previous value.
- W If this bit is 0, the TCP window has not changed (the typical case). If this bit is 1, *Δwin* is the current value of *th_win* minus its previous value.
- U If this bit is 0, the URG flag in the segment is not set and the urgent offset has not changed from its previous value (the typical case). If this bit is 1, *urgoff* is the current value of *th_urg* and the URG flag is set. If the urgent offset changes without the URG flag being set, the segment is sent uncompressed. (This often occurs in the first segment following urgent data.)

The differences are encoded as the current value minus the previous value, because most of these differences will be small positive numbers (with *Δwin* being an exception) given the way these fields normally change.

We note that five of the optional fields in Figure 29.34 are encoded in 0, 1, or 3 bytes.

- 0 bytes: If the corresponding flag is not set, nothing is encoded for the field.
- 1 byte: If the value to send is between 1 and 255, a single byte encodes the value.
- 3 bytes: If the value to send is either 0 or between 256 and 65535, 3 bytes encode the value: the first byte is 0, followed by the 2-byte value. This always works for the three 16-bit values, *urgoff*, *Δwin*, and *Δipid*; but if the difference to encode for the two 32-bit values, *Δack* and *Δseq*, is less than 0 or greater than 65535, the segment is sent uncompressed.

If we compare the nonshaded fields in Figure 29.33 with the possible fields in Figure 29.34 we notice that some fields are never transmitted.

- The IP total length field is not transmitted since most link layers provide the length of a received message to the receiver.
- Since the only field in the IP header that is being transmitted is the identification field, the IP checksum is also omitted. This is a hop-by-hop checksum that protects only the IP header across any given link.

Special Cases

Two common cases are detected and transmitted as special combinations of the 4 low-order bits: *SAWLI*. Since urgent data is rare, if the URG flag in the segment is set and both the sequence number and window also change (implying that the 4 low-order bits would be 1011 or 1111), the segment is sent uncompressed. Therefore if the 4 low-order bits are sent as 1011 (called *5A) or 1111 (called *5), the following two special cases apply:

- *5A The sequence number and acknowledgment number both increase by the amount of data in the last segment, the window and urgent offset don't change, and the URG flag is not set. This special case avoids encoding both Δseq and Δack .

This case occurs frequently for both directions of echoed terminal traffic. Figures 19.3 and 19.4 of Volume 1 give examples of this type of data flow across an Rlogin connection.

- *5 The sequence number changes by the amount of data in the last segment, the acknowledgment number, window, and urgent offset don't change, and the URG flag is not set. This special case avoids encoding Δseq .

This case occurs frequently for the sending side of a unidirectional data transfer (e.g., FTP). Figures 20.1, 20.2, and 20.3 of Volume 1 give examples of this type of data transfer. This case also occurs for the sender of nonechoed terminal traffic (e.g., commands that are not echoed by a full-screen editor).

Examples

Two simple examples were run across the SLIP link between the systems *bsd1* and *slip* in Figure 1.17. This SLIP link uses header compression in both directions. The *tcpdump* program described in Appendix A of Volume 1 was also run on the host *bsd1* to save a copy of all the frames. This program has an option that outputs the compressed header, showing all the fields in Figure 29.34.

Two traces were obtained: a short portion of an Rlogin connection and a file transfer from *bsd1* to *slip* using FTP. Figure 29.36 shows a summary of the different frame types for both connections.

The two entries of 75 verify our claim that this special case often occurs for both directions of echoed terminal traffic. The entry of 325 verifies our claim that this special case occurs frequently for the sending side of a unidirectional data transfer.

The 10 frames of type *IP* for the FTP example correspond to four segments with the SYN flag set and six segments with the FIN flag set. FTP uses two connections: one for the interactive commands and one for the file transfer.

The *UNCOMPRESSED_TCP* frame types normally correspond to the first segment following connection establishment, the one that establishes the connection ID. An additional few are seen in these examples when the type of service is set (the Net/3 Rlogin and FTP clients and servers all set the TOS field *after* the connection is established).

frame type	Rlogin		FTP	
	input	output	input	output
IP	1	1	5	5
UNCOMPRESSED_TCP	3	2	2	3
COMPRESSED_TCP				
*SA special case	75	75	0	0
*S special case	25	1	1	325
nonspecial	9	93	337	13
Total	113	172	345	346

Figure 29.36 Counts of different frame types for Rlogin and FTP connections.

#bytes	Rlogin		FTP	
	input	output	input	output
3	102	44	2	250
4		94		78
5	7	12	5	2
6		6	325	5
7		13	2	1
8				1
9			4	1
Total	109	169	338	338

Figure 29.37 Distribution of compressed-header sizes.

Figure 29.37 shows the distribution of the compressed-header sizes. The average size of the compressed header for the final four columns in Figure 29.37 is 3.1, 4.1, 6.0, and 3.3 bytes, a significant savings compared to the uncompressed 40-byte headers, especially for the interactive connection.

Almost all of the 325 6-byte headers in the FTP input column contain only a *ack* of 256, which being greater than 255 is encoded in 3 bytes. The SLIP MTU is 296, so TCP uses an MSS of 256. Almost all of the 250 3-byte headers in the FTP output column contain the *S special case (sequence number change only) with a change of 256 bytes. But since this change refers to the amount of data in the previous segment, nothing is transmitted other than the flag byte and the TCP checksum. The 78 4-byte headers in the FTP output column are this same special case, but with a change in the IP identification field also (Exercise 29.8).

Configuration

Header compression must be enabled on a given SLIP or PPP link. With a SLIP link there are normally two flags that can be set when the interface is configured: enable header compression and autoenable header compression. These two flags are set using

the `link0` and `link2` flags to the `ifconfig` command, respectively. Normally a client (the dialin host) decides whether to use header compression or not. The server (the host or terminal server to which the client dials in) specifies the autoenable flag only: if header compression is enabled by the client, its TCP will send a datagram of type `UNCOMPRESSED_TCP` to specify the connection ID. When the server sees this packet it enables header compression (since it was in the autoenable mode). If the server never sees this type of packet, it never enables header compression for this line.

PPP allows the negotiation of options between the two ends of the link when the link is established. One of the options that can be negotiated is whether to use header compression or not.

29.14 Summary

This chapter completes our detailed look at TCP input processing. We started with the processing of an ACK in the `SYN_RCVD` state, which completes a passive open, a simultaneous open, or a self connect.

The fast retransmit algorithm lets TCP detect a dropped segment after receiving a specified number of consecutive duplicate ACKs and retransmit the segment before the retransmission timer expires. Net/3 combines the fast retransmit algorithm with the fast recovery algorithm, which tries to keep the data flowing from the sender to the receiver, albeit at a slower rate, using congestion avoidance but not slow start.

ACK processing then discards the acknowledged data from the socket's send buffer and handles a few TCP states specially, when the receipt of an ACK changes the connection state.

The URG flag is processed, if set, and TCP's urgent mode is mapped into the socket abstraction of out-of-band data. This is complicated because the process can receive the out-of-band byte inline or in a special out-of-band buffer, and TCP can receive urgent notification before the data byte referenced by the urgent pointer has been received.

TCP input processing completes by calling `TCP_REASS` to merge the received data into either the socket's receive buffer or the socket's out-of-order queue, processing the FIN flag, and calling `tcp_output` if a segment must be generated in response to the received segment.

TCP header compression is a technique used on SLIP and PPP links to reduce the size of the IP and TCP headers from the normal 40 bytes to around 3–6 bytes (typically). This is done by recognizing that most fields in these headers don't change from one segment to the next on a given connection, and the fields that do change often change by a small amount. This allows a flag byte to be sent indicating which fields have changed, and the changes are encoded as differences from the previous segment.

Exercises

- 29.1 A client connects to a server and no segments are lost. Which process, the client or server, completes its open of the connection first?
- 29.2 A Net/3 system receives a SYN for a listening socket and the SYN segment also contains 50 bytes of data. What happens?
- 29.3 Continue the previous exercise assuming that the client does not retransmit the 50 bytes of data; instead the client responds with a segment that acknowledges the server's SYN/ACK and contains a FIN. What happens?
- 29.4 A Net/3 client performs a passive open to a listening server. The server's response to the client's SYN is a segment with the expected SYN/ACK, but the segment also contains 50 bytes of data and the FIN flag. List the processing steps for the client's TCP.
- 29.5 Figure 18.19 in Volume 1 and Figure 14 in RFC 793 both show four segments exchanged during a simultaneous close. But if we trace a simultaneous close between two Net/3 systems, or if we watch the close sequence following a self-connect on a Net/3 system, we see six segments, not four. The extra two segments are a retransmission of the FIN by each end when the other's FIN is received. Where is the bug and what is the fix?
- 29.6 Page 72 of RFC 793 says that when data in the send buffer is acknowledged by the other end "Users should receive positive acknowledgments for buffers which have been sent and fully acknowledged (i.e., send buffer should be returned with 'ok' response)." Does Net/3 provide this notification?
- 29.7 What effect do the options defined in RFC 1323 have on TCP header compression?
- 29.8 What effect does the Net/3 assignment of the IP identification field have on TCP header compression?

TCP User Requests

30.1 Introduction

This chapter looks at the TCP user-request function `tcp_usrreq`, which is called as the protocol's `pr_usrreq` function to handle many of the system calls that reference a TCP socket. We also look at `tcp_ctloutput`, which is called when the process calls `setsockopt` for a TCP socket. Finally, since this is the last of the TCP chapters, we also look at how the Net/3 implementation of TCP handles the RFC 1122 requirements.

30.2 `tcp_usrreq` Function

TCP's user-request function is called for a variety of operations. Figure 30.1 shows the beginning and end of `tcp_usrreq`. The body of the `switch` is shown in following figures. The function arguments, some of which differ depending on the request, are described in Figure 15.17.

`in_control` processes `ioctl` requests

45-58 The `PRU_CONTROL` request is from the `ioctl` system call. The function `in_control` processes the request completely.

Control information is invalid

59-64 A call to `sendmsg` specifying control information is invalid for a TCP socket. If this happens, the `mbufs` are released and `EINVAL` is returned.

65-66 This remainder of the function executes at `splnet`. This is overly conservative locking to avoid sprinkling the individual `case` statements with calls to `splnet` when the calls are really necessary. As we mentioned with Figure 23.15, setting the processor priority to `splnet` only stops a software interrupt from causing the IP input routine to

```

45 int
46 tcp_usrreq(so, req, m, nam, control)
47 struct socket *so;
48 int req;
49 struct mbuf *m, *nam, *control;
50 {
51     struct inpcb *inp;
52     struct tcpcb *tp;
53     int s;
54     int error = 0;
55     int ostate;
56     if (req == PRU_CONTROL)
57         return (in_control(so, (int) m, (caddr_t) nam,
58                             (struct ifnet *) control));
59     if (control && control->m_len) {
60         m_freem(control);
61         if (m)
62             m_freem(m);
63         return (EINVAL);
64     }
65     s = splnet();
66     inp = sotoinpcb(so);
67     /*
68      * When a TCP is attached to a socket, then there will be
69      * a (struct inpcb) pointed at by the socket, and this
70      * structure will point at a subsidiary (struct tcpcb).
71      */
72     if (inp == 0 && req != PRU_ATTACH) {
73         splx(s);
74         return (EINVAL); /* XXX */
75     }
76     if (inp) {
77         tp = intotcpcb(inp);
78         /* WHAT IF TP IS 0? */
79         ostate = tp->t_state;
80     } else
81         ostate = 0;
82     switch (req) {
83
84         /* switch cases */
85
86     default:
87         panic("tcp_usrreq");
88     }
89     if (tp && (so->so_options & SO_DEBUG))
90         tcp_trace(TA_USER, ostate, tp, (struct tcpiphdr *) 0, req);
91     splx(s);
92     return (error);
93 }

```

Figure 30.1 Body of tcp_usrreq function.

be executed (which could call `tcp_input`). It does not prevent the interface layer from accepting incoming packets and placing them onto IP's input queue.

The pointer to the Internet PCB is obtained from the `socket` structure pointer. The only time the resulting PCB pointer is allowed to be a null pointer is when the `PRU_ATTACH` request is issued, which occurs in response to the `socket` system call.

67-81 If `inp` is nonnull, the current connection state is saved in `ostate` for the call to `tcp_trace` at the end of the function.

We now discuss the individual case statements. The `PRU_ATTACH` request, shown in Figure 30.2, is issued by the `socket` system call and by `sonewconn` when a connection request arrives for a listening socket (Figure 28.7).

```

83      /*
84      * TCP attaches to socket via PRU_ATTACH, reserving space,
85      * and an internet control block.
86      */
87      case PRU_ATTACH:
88          if (inp) {
89              error = EISCONN;
90              break;
91          }
92          error = tcp_attach(so);
93          if (error)
94              break;
95          if ((so->so_options & SO_LINGER) && so->so_linger == 0)
96              so->so_linger = TCP_LINGERTIME;
97          tp = sototcpch(so);
98          break;
99      /*
100     * PRU_DETACH detaches the TCP protocol from the socket.
101     * If the protocol state is non-embryonic, then can't
102     * do this directly; have to initiate a PRU_DISCONNECT,
103     * which may finish later; embryonic TCB's can just
104     * be discarded here.
105     */
106     case PRU_DETACH:
107         if (tp->t_state > TCPS_LISTEN)
108             tp = tcp_disconnect(tp);
109         else
110             tp = tcp_close(tp);
111         break;

```

Figure 30.2 `tcp_usrreq` function: `PRU_ATTACH` and `PRU_DETACH` requests.

PRU_ATTACH request

83-94 If the socket structure already points to a PCB, `EISCONN` is returned. `tcp_attach` completes the processing; it allocates and initializes the Internet PCB and the TCP control block.

95-96 If the `SO_LINGER` socket option is set, and the linger time is 0, it is set to 120 (`TCP_LINGERTIME`).

How can a socket option be set before the `PRU_ATTACH` request is issued? It is impossible to set a socket option before calling `socket`, but `socketconn` also issues the `PRU_ATTACH` request. The `PRU_ATTACH` request is issued after `socketconn` copies the `so_options` from the listening socket to the newly created socket. This code prevents a newly accepted connection from inheriting a linger time of 0 from the listening socket.

There is a bug here. The constant `TCP_LINGERTIME` is initialized to 120 in the header `tcp_timer.h` with the comment "linger at most 2 minutes." But the `so_linger` value becomes the final argument to the kernel's `msleep` function (called from `socket`), which becomes the final argument to the kernel's `timeout` function and is in clock ticks, not seconds. If the system's clock-tick frequency (`hz`) is 100, this value for the linger time is 1.2 seconds, not 2 minutes.

97 `tp` is now set to the pointer to the socket's TCP control block. This is required at the end, in case the `SO_DEBUG` socket option is set.

PRU_DETACH request

98-111 The `close` system call issues the `PRU_DETACH` request if the `PRU_DISCONNECT` request fails. If the connection has not been completed (the connection state is less than `ESTABLISHED`), nothing needs to be sent to the other end. But if the connection has been established, `tcp_disconnect` initiates TCP's connection-close sequence (e.g., any pending data is sent, followed by a `FIN`).

The test for the state being greater than `LISTEN` is incorrect, because if the state is `SYN_SENT` or `SYN_RCVD`, both of which are greater than `LISTEN`, `tcp_disconnect` just calls `tcp_close`. This case could be simplified by just calling `tcp_disconnect`.

Figure 30.3 shows the processing for the `bind` and `listen` system calls.

```

112      /*
113      * Give the socket an address.
114      */
115      case PRU_BIND:
116          error = in_pcbbind(inp, nam);
117          if (error)
118              break;
119          break;
120      /*
121      * Prepare to accept connections.
122      */
123      case PRU_LISTEN:
124          if (inp->inp_lport == 0)
125              error = in_pcbbind(inp, (struct mbuf *) 0);
126          if (error == 0)
127              tcp->t_state = TCPS_LISTEN;
128          break;

```

Figure 30.3 `tcp_usrreq` function: `PRU_BIND` and `PRU_LISTEN` requests.

112-119 All the work for a `PRU_BIND` request is done by `in_pcbbind`.

120-128 For the PRU_LISTEN request, if the socket has not been bound with a local port, `in_pcbbind` assigns one automatically. This is rare, since most servers explicitly bind their well-known port, although RPC (remote procedure call) servers typically bind an ephemeral port and then register the port with the *Port Mapper*. (Section 29.4 of Volume 1 describes the Port Mapper.) The connection state is set to LISTEN. This is the main purpose of `listen`: to set the socket's state so that incoming connections are accepted (i.e., a passive open).

Figure 30.4 shows the processing for the `connect` system call: an active open normally initiated by a client.

```

129      /* tcp_usrreq.c
130      * Initiate connection to peer.
131      * Create a template for use in transmissions on this connection.
132      * Enter SYN_SENT state, and mark socket as connecting.
133      * Start keepalive timer, and seed output sequence space.
134      * Send initial segment on connection.
135      */
136      case PRU_CONNECT:
137          if (inp->inp_lport == 0) {
138              error = in_pcbbind(inp, (struct mbuf *) 0);
139              if (error)
140                  break;
141          }
142          error = in_pcbconnect(inp, nam);
143          if (error)
144              break;
145          tp->t_template = tcp_template(tp);
146          if (tp->t_template == 0) {
147              in_pcbdisconnect(inp);
148              error = ENOBUFS;
149              break;
150          }
151          /* Compute window scaling to request. */
152          while (tp->request_r_scale < TCP_MAX_WINSHIFT &&
153                (TCP_MAXWIN << tp->request_r_scale) < so->so_rcv.sb_hiwat)
154              tp->request_r_scale++;
155          soisconnecting(so);
156          tcpstat.tcps_connattempt++;
157          tp->t_state = TCPS_SYN_SENT;
158          tp->t_timer[TCPP_KEEP] = TCPTV_KEEP_INIT;
159          tp->iss = tcp_iss;
160          tcp_iss += TCP_ISSINCR / 2;
161          tcp_sendseqinit(tp);
162          error = tcp_output(tp);
163          break;
tcp_usrreq.c

```

Figure 30.4 `tcp_usrreq` function: PRU_CONNECT request.

Assign ephemeral port

128-141 If the socket has not been bound with a local port, `in_pcbbind` assigns one automatically. This is typical for clients, which normally don't care about the value of the local port.

Connect PCB

163-166 `in_pcbconnect` acquires a route to the destination, determines the outgoing interface, and verifies that the socket pair is unique.

Initialize IP and TCP headers

167-150 `tcp_template` allocates an mbuf for a copy of the IP and TCP headers, and it initializes both headers with as much information as possible. The only way for this function to fail is for the kernel to run out of mbufs.

Calculate window scale factor

151-154 The window scale value for the receive buffer is calculated: 65535 (`TCP_MAXWIN`) is left shifted until the value is greater than or equal to the size of the receive buffer (`so_rcv.sb_hiwat`). The resulting shift count (between 0 and 14) is the scale factor that will be sent in the SYN. (We saw identical code in Figure 28.7 that was executed for a passive open.) Since the window scale option is sent in the SYN resulting from a connect, the process must set the `SO_RCVBUF` socket option before calling `connect`, or the default buffer size is used (`tcp_recvspace` from Figure 24.3).

Set socket and connection state

155-158 `soisconnecting` sets the appropriate bits in the socket's state variable, and the state of the TCP connection is set to `SYN_SENT`. This causes the call to `tcp_output` that follows to send the SYN (see the `tcp_outflags` value in Figure 24.16). The connection-establishment timer is initialized to 75 seconds. `tcp_output` will also set the retransmission timer for the SYN, as shown in Figure 25.16.

Initialize sequence numbers

161-161 The initial send sequence number is copied from the global `tcp_iss`. This global is then incremented by 64,000 (`TCP_ISSINCR` divided by 2). We saw this same handling of `tcp_iss` when the ISS was initialized after a listening server received a SYN (Figure 28.17). The send sequence numbers are then initialized by `tcp_sendseqinit`.

Send initial SYN

162 `tcp_output` sends the initial SYN to initiate the connection. A local error (for example, out of mbufs or no route to destination) is returned by `tcp_output`, which becomes the return value from `tcp_usrreq`, which is returned to the process.

Figure 30.5 shows the processing for the `PRU_CONNECT2`, `PRU_DISCONNECT`, and `PRU_ACCEPT` requests.

164-169 The `PRU_CONNECT2` request, a result of the `socketpair` system call, is invalid for the TCP protocol.

170-183 The `close` system call issues the `PRU_DISCONNECT` request. If the connection has been established, a FIN must be sent and the normal TCP close sequence followed. This is done by `tcp_disconnect`.

```

164      /*
165       * Create a TCP connection between two sockets.
166       */
167      case PRU_CONNECT2:
168          error = EOPNDTSUPP;
169          break;
170
171      /*
172       * Initiate disconnect from peer.
173       * If connection never passed embryonic stage, just drop;
174       * else if don't need to let data drain, then can just drop anyway.
175       * else have to begin TCP shutdown process: mark socket disconnecting,
176       * drain unread data, state switch to reflect user close, and
177       * send segment (e.g. FIN) to peer. Socket will be really disconnected
178       * when peer sends FIN and acks ours.
179       *
180       * SHOULD IMPLEMENT LATER PRU_CONNECT VIA REALLOC TCPCB.
181       */
182      case PRU_DISCONNECT:
183          tp = tcp_disconnect(tp);
184          break;
185
186      /*
187       * Accept a connection. Essentially all the work is
188       * done at higher levels; just return the address
189       * of the peer, starting through addr.
190       */
191      case PRU_ACCEPT:
192          in_setpeeraddr(inp, nam);
193          break;

```

Figure 30.5 tcp_usrreq function: PRU_CONNECT2, PRU_DISCONNECT, and PRU_ACCEPT requests.

The comment beginning with "SHOULD IMPLEMENT" refers to the fact that a socket that encounters an error cannot be reused. For example, if a client issues a `connect` and receives an error, it cannot issue another `connect` on the same socket. Instead, the socket with the error must be closed, a new socket created with `socket`, and the `connect` issued on the new socket.

186-191 All the work associated with the `accept` system call is done by the socket layer and the protocol layer. The `PRU_ACCEPT` request just returns the IP address and port number of the peer to the process.

The `PRU_SHUTDOWN`, `PRU_RCVD`, and `PRU_SEND` requests are processed in Figure 30.6.

PRU_SHUTDOWN request

192-200 This request is issued by `soshutdown` when the process calls `shutdown` to prevent any further output. `socant_sendmore` sets the socket's flags to prevent any future output. `tcp_usrclosed` sets the connection state according to Figure 24.15. `tcp_output` attempts to send the FIN, but if there is still pending data to send to the other end, that data is sent before the FIN is sent.

```

192      /*
193      * Mark the connection as being incapable of further output.
194      */
195      case PRU_SHUTDOWN:
196          socant_sendmore(sc);
197          tp = tcp_usrclosed(tp);
198          if (tp)
199              error = tcp_output(tp);
200          break;

201      /*
202      * After a receive, possibly send window update to peer.
203      */
204      case PRU_RCVD:
205          (void) tcp_output(tp);
206          break;

207      /*
208      * Do a send by putting data in output queue and updating urgent
209      * marker (if URG set). Possibly send more data.
210      */
211      case PRU_SEND:
212          sbappend(&so->so_snd, m);
213          error = tcp_output(tp);
214          break;

```

tcp_usrreq.c

Figure 30.6 tcp_usrreq function: PRU_SHUTDOWN, PRU_RCVD, and PRU_SEND requests.

PRU_RCVD request

201-206 This request is issued by `soreceive` after the process has read data from the socket's receive buffer. TCP needs to know about this since the receive buffer may now have enough room to allow the advertised window to increase. `tcp_output` will determine whether a window update segment should be sent.

PRU_SEND request

207-214 In Figure 23.14 we showed how the five write functions ended up issuing this request. `sbappend` adds the data to the socket's send buffer (where it must wait until acknowledged by the other end), and `tcp_output` sends a segment, if possible.

Figure 30.7 shows the processing of the PRU_ABORT and PRU_SENSE requests.

PRU_ABORT request

218-220 A PRU_ABORT request is issued for a TCP socket by `socklose` if the socket is a listening socket (e.g., a server) and if there are pending connections for the server that have already initiated or completed the three-way handshake, but have not been accepted by the server yet. `tcp_drop` sends an RST if the connection is synchronized.


```

215      /*
216       * Abort the TCP.
217       */
218      case PRU_ABORT:
219          tp = tcp_drop(tp, ECONNABORTED);
220          break;
221
222      case PRU_SENSE:
223          ((struct stat *) m)->st_blksize = so->so_snd.sb_hiwat;
224          (void) splx(s);
225          return (0);

```

Figure 30.7 tcp_usrreq function: PRU_ABORT and PRU_SENSE requests.

PRU_SENSE request

221-224 The `fstat` system call generates the PRU_SENSE request. TCP returns the size of the send buffer as the `st_blksize` element of the `stat` structure.

Figure 30.8 shows the PRU_RCVOOB request, issued by `soreceive` when the process issues a read system call specifying the MSG_OOB flag to read out-of-band data.

```

225      case PRU_RCVOOB:
226          if ((so->so_oobmark == 0 &&
227              (so->so_state & SS_RCVATMARK) == 0) ||
228              (so->so_options & SO_COBINLINE)) {
229              tp->t_oobflags & TCPOOB_HADDATA) {
230                  error = EINVAL;
231                  break;
232              }
233              if ((tp->t_oobflags & TCPOOB_HAVEDATA) == 0) {
234                  error = EWOULDBLOCK;
235                  break;
236              }
237              m->m_len = 1;
238              *mtd(m, caddr_t) = tp->t_iebc;
239              if (((int) nam & MSG_PREK) == 0)
240                  tp->t_oobflags ^= (TCPOOB_HAVEDATA | TCPOOB_HADDATA);
241              break;

```

Figure 30.8 tcp_usrreq function: PRU_RCVOOB request.

Verify that reading out-of-band data is appropriate

225-232 It is an error for the process to try to read out-of-band data if any one of the following three conditions is true:

1. if the socket's out-of-band mark is 0 (`so_oobmark`) and the socket is not at the mark (the `SS_RCVATMARK` flag is not set), or

2. if the `SO_OOBINLINE` socket option is set, or
3. if the `TCPOOB_HADDATA` flag is set for the connection (i.e., the connection did have an out-of-band byte, but it has already been read).

The error `EINVAL` is returned if any one of these is true.

Check that out-of-band byte has arrived

229-236 If none of the three conditions above is true, but the `TCPOOB_HAVEDATA` flag is false, this indicates that TCP has received an urgent mode notification from the other end, but the byte whose sequence number is 1 less than the urgent pointer has not been received yet (Figure 29.17). The error `EWOULDBLOCK` is returned. It is possible for TCP to send an urgent notification with an urgent offset referencing a byte that the sender has not been able to send yet. Figure 26.7 of Volume 1 shows an example of this scenario, which often happens if the sender's data transmission has been stopped by a zero-window advertisement.

Return out-of-band byte

237-238 The single byte of out-of-band data that was stored in `t_ioobc` by `tcp_pulloutofband` is returned to the process.

Flip flags

239-241 If the process is actually reading the out-of-band byte (as compared to peeking at it with the `MSG_PEEK` flag), this exclusive OR turns the `HAVE` flag off and the `HAD` flag on. We are guaranteed at this point in the `case` statement that the `HAVE` flag is set and the `HAD` flag is cleared. The purpose of the `HAD` flag is to prevent the process from trying to read the out-of-band byte more than once. Once the `HAD` flag is set, it is not cleared until a new urgent pointer is received from the other end (Figure 29.17).

The reason for this hard-to-understand exclusive OR, instead of the simpler

```
tp->t_oobflags = TCPOOB_HADDATA;
```

is to allow additional bits in `t_oobflags` to be used. Net/3, however, only uses the 2 bits that we've described.

The `PRU_SENDOOB` request, shown in Figure 30.9, is issued by `so_send` when the process writes data and specifies the `MSG_OOB` flag.

Check for room and append to send buffer

242-249 The process is allowed to exceed the size of the send buffer by up to 512 bytes when sending out-of-band data. The socket layer is more permissive, allowing out-of-band data to exceed the size of the send buffer by 1024 bytes (Figure 16.24). `sbappend` adds the data to the end of the send buffer.

Calculate urgent pointer

249-257 The urgent pointer (`snd_up`) points to the byte following the final byte from the write request. We showed this in Figure 26.30, assuming the process writes 3 bytes of data with the `MSG_OOB` flag set and that the send buffer was empty. Realize that if the

```

242     case PRU_SENDOOB:                                     tcp_usrreq.c
243         if (subspace(&so->so_snd) < -512) {
244             m_freem(m);
245             error = ENOBUFS;
246             break;
247         }
248         /*
249          * According to RFC961 (Assigned Protocols),
250          * the urgent pointer points to the last octet
251          * of urgent data. We continue, however,
252          * to consider it to indicate the first octet
253          * of data past the urgent section.
254          * Otherwise, snd_up should be one lower.
255          */
256         sbappend(&so->so_snd, m);
257         tp->snd_up = tp->snd_una + so->so_snd.sb_cc;
258         tp->t_force = 1;
259         error = tcp_output(tp);
260         tp->t_force = 0;
261     break;                                               tcp_usrreq.c

```

Figure 30.9 tcp_usrreq function: PRU_SENDOOB request.

process writes more than 1 byte of data with the MSG_OOB flag set, only the final byte is considered the out-of-band byte when the data is received by a Berkeley-derived system.

Force TCP output

258-261 `t_force` is set to 1 and `tcp_output` is called. This causes a segment to be sent with the URG flag set and with a nonzero urgent offset, even if no data can be sent because of a zero-window advertisement. Figure 26.7 of Volume 1 shows the transmission of an urgent segment into a closed window.

The final three requests are shown in Figure 30.10.

262-267 The `getsockname` and `getpeername` system calls issue the `PRU_SOCKADDR` and `PRU_PEERADDR` requests, respectively. The functions `in_setsockaddr` and `in_setpeeraddr` fetch the information from the PCB, storing the result in the `addr` argument.

268-275 The `PRU_SLOWTIMO` request is issued by the `tcp_slowtimo` function. As the comment indicates, the only reason `tcp_slowtimo` doesn't call `tcp_timers` directly is to allow the timer expiration to be traced by the call to `tcp_trace` at the end of the function (Figure 30.1). For the trace record to show which one of the four TCP timer counters expired, `tcp_slowtimo` passes the index into the `t_timer` array (Figure 25.1) as the `nam` argument, and this is left shifted 8 bits and logically ORed into the request value (`req`). The `trpt` program knows about this hack and handles it accordingly.

```

262     case PRU_SOCKADDR:
263         in_setsockaddr(inp, nam);
264         break;
265
266     case PRU_PEERADDR:
267         in_setpeeraddr(inp, nam);
268         break;
269
270     /*
271      * TCP slow timer went off; going through this
272      * routine for tracing's sake.
273      */
274     case PRU_SLOWTMR:
275         tp = tcp_timers[tp, (int) nam];
276         req |= (int) nam << 8; /* for debug's sake */
277         break;

```

Figure 30.10 `tcp_usrreq` function: `PRU_SOCKADDR`, `PRU_PEERADDR`, and `PRU_SLOWTMR` requests.

30.3 `tcp_attach` Function

The `tcp_attach` function is called by `tcp_usrreq` to process the `PRU_ATTACH` request (i.e., when the socket system call is issued or when a new connection request arrives for a listening socket). Figure 30.11 shows the code.

Allocate space for send buffer and receive buffer

483-492 If space has not been allocated for the socket's send and receive buffers, `sbr@serve` sets them both to 8192, the default values of the global variables `tcp_sendspace` and `tcp_recvspace` (Figure 24.3).

Whether these defaults are adequate depends on the MSS for each direction of the connection, which depends on the MTU. For example, [Comer and Lin 1994] show that anomalous behavior occurs if the send buffer is less than three times the MSS, which drastically reduces performance. Some implementations have much higher defaults, such as 61,444 bytes, realizing the effect these defaults have on performance, especially with higher MTUs (e.g., FDDI and ATM).

Allocate Internet PCB and TCP control block

373-377 `in_pcballoc` allocates an Internet PCB and `tcp_newtcpcb` allocates a TCP control block and links it to the PCB.

378-384 The code with the comment XXX is executed if the call to `malloc` in `tcp_newtcpcb` fails. Remember that the `PRU_ATTACH` request is issued as a result of the socket system call, and when a connection request arrives for a listening socket (`sonewconn`). In the latter case the socket flag `SS_NOFDREF` is set. If this flag is left on, the call to `sofree` by `in_pcbdetach` releases the socket structure. As we saw in `tcp_input`, this structure should not be released until that function is done with the received segment (the `dropsocket` flag in Figure 29.27). Therefore the current value of the `SS_NOFDREF` flag is saved in the variable `nofd` when `in_pcbdetach` is called, and reset before `tcp_attach` returns.

385-386 The TCP connection state is initialized to `CLOSED`.


```

361 int
362 tcp_attach(so)
363 struct socket *so;
364 {
365     struct tcpcb *tp;
366     struct inpcb *inp;
367     int     error;
368
369     if (so->so_snd.sb_hiwat == 0 || so->so_rcv.sb_hiwat == 0) {
370         error = sorereserve(so, tcp_sendspace, tcp_recvspace);
371         if (error)
372             return (error);
373     }
374     error = in_pcballoc(so, &tcp);
375     if (error)
376         return (error);
377     inp = sotoinpcb(so);
378     tp = tcp_newtcpcb(inp);
379     if (tp == 0) {
380         int     nofd = so->so_state & SS_NOFDREF;    /* XXX */
381         so->so_state &= TSS_NOFDREF;    /* don't free the socket yet */
382         in_pcbdetach(inp);
383         so->so_state |= nofd;
384         return (ENOBUFFS);
385     }
386     tp->t_state = TCPS_CLOSED;
387     return (0);

```

Figure 30.11 tcp_attach function: create a new TCP socket.

30.4 tcp_disconnect Function

tcp_disconnect, shown in Figure 30.12, initiates a TCP disconnect.

Connection not yet synchronized

396-402 If the socket is not yet in the ESTABLISHED state (i.e., LISTEN, SYN_SENT, or SYN_RCVD), tcp_close just releases the PCB and the TCP control block. Nothing needs to be sent to the other end since the connection has not been synchronized.

Hard disconnect

403-404 If the connection is synchronized, the SO_LINGER socket option is set, and the linger time (so_linger) is set to 0, the connection is dropped by tcp_drop. This sets the connection state to CLOSED, sends an RST to the other end, and releases the PCB and TCP control block. The connection does not pass through the TIME_WAIT state. The call to close that caused the PRU_DISCONNECT request will discard any data still in the send or receive buffers.

If the SO_LINGER socket option has been set with a nonzero linger time, it is handled by soclose.

```

396 struct tcpb *
397 tcp_disconnect(tp)
398 struct tcpb *tp;
399 {
400     struct socket *so = tp->t_inpcb->inp_socket;
401     if (tp->t_state < TCP_ESTABLISHED)
402         tp = tcp_close(tp);
403     else if ((so->so_options & SO_LINGER) && so->so_linger == 0)
404         tp = tcp_drop(tp, 0);
405     else {
406         soisdisconnecting(so);
407         sbflush(&so->so_rcv);
408         tp = tcp_usrclosed(tp);
409         if (tp)
410             (void) tcp_output(tp);
411     }
412     return (tp);
413 }

```

Figure 30.12 tcp_disconnect function: initiate TCP disconnect.

Graceful disconnect

405-406 This code is executed when the connection has been synchronized but the `SO_LINGER` option either was not set or was set with a nonzero linger time. TCP's normal connection termination steps must be followed. `soisdisconnecting` sets the socket's state.

Discard pending receive data

407 Any pending data in the receive buffer is discarded by `sbflush`, since the process has closed the socket. The send buffer is left alone, however, and `tcp_output` will try to send what remains. We say "try" because there's no guarantee that the data still to be sent will be transmitted successfully. The other end might crash before it receives and acknowledges the data, or even if the TCP module at the other end receives and acknowledges the data, the system might crash before the application at the other end reads the data. Since the local process has closed the socket, if TCP gives up trying to send what remains in the send buffer (because its retransmission timer finally expires), there is no way to notify the process of the error.

Change connection state

408-410 `tcp_usrclosed` moves the connection into the next state, based on the current state. This normally moves the connection to the `FIN_WAIT_1` state, since the connection is typically closed from the `ESTABLISHED` state. We'll see that `tcp_usrclosed` always returns the current control block pointer (`tp`), since the state must be synchronized to get to this point in the code, so `tcp_output` is always called to send a segment. If the connection moves from the `ESTABLISHED` to the `FIN_WAIT_1` state, this causes a FIN to be sent.

30.5 tcp_usrclosed Function

This function, shown in Figure 30.13, is called from `tcp_disconnect` and when the `PRU_SHUTDOWN` request is processed.

```

424 struct tcpcb *
425 tcp_usrclosed(tp)
426 struct tcpcb *tp;
427 {
428     switch (tp->t_state) {
429     case TCPS_CLOSED:
430     case TCPS_LISTEN:
431     case TCPS_SYN_SENT:
432         tp->t_state = TCPS_CLOSED;
433         tp = tcp_close(tp);
434         break;
435     case TCPS_SYN_RECEIVED:
436     case TCPS_ESTABLISHED:
437         tp->t_state = TCPS_FIN_WAIT_1;
438         break;
439     case TCPS_CLOSE_WAIT:
440         tp->t_state = TCPS_LAST_ACK;
441         break;
442     }
443     if (tp && tp->t_state >= TCPS_FIN_WAIT_2)
444         soisdisconnected(tp->t_inpcb->inp_socket);
445     return (tp);
446 }

```

tcp_usrreq.c

Figure 30.13 tcp_usrclosed function: move connection to next state, based on process close.

Simple close when SYN not received

429-434 If a SYN has not been received on the connection, a FIN need not be sent. The new state is CLOSED and `tcp_close` releases the Internet PCB and the TCP control block.

Move to FIN_WAIT_1 state

435-438 In the SYN_RCVD and ESTABLISHED states, the new state is FIN_WAIT_1, which causes the next call to `tcp_output` to send a FIN (the `tcp_outflags` value in Figure 24.16).

Move to LAST_ACK state

439-442 In the CLOSE_WAIT state, the close moves the connection into the LAST_ACK state. The next call to `tcp_output` will cause a FIN to be sent.

443-444 If the connection state is either FIN_WAIT_2 or TIME_WAIT, `soisdisconnected` marks the socket state appropriately.

30.6 tcp_ctloutput Function

The `tcp_ctloutput` function is called by the `getsockopt` and `setsockopt` system calls when the descriptor argument refers to a TCP socket and when the level is not `SOL_SOCKET`. Figure 30.14 shows the two socket options supported by TCP.

optname	Variable	Access	Description
<code>TCP_NODELAY</code>	<code>t_flags</code>	read, write	Nagle algorithm (Figure 26.8)
<code>TCP_MAXSEG</code>	<code>t_maxseg</code>	read, write	maximum segment size TCP will send

Figure 30.14 Socket options supported by TCP.

Figure 30.15 shows the first part of the function.

```

284 int
285 tcp_ctloutput(op, so, level, optname, mp)
286 int op;
287 struct socket *so;
288 int level, optname;
289 struct mbuf **mp;
290 {
291     int error = 0, s;
292     struct inpcb *inp;
293     struct tcpcb *tp;
294     struct mbuf *m;
295     int i;
296
297     s = splnet();
298     inp = sotoinpcb(so);
299     if (inp == NULL) {
300         splx(s);
301         if (op == PRCO_SETOPT && *mp)
302             (void) m_free(*mp);
303         return (ECONNRESET);
304     }
305     if (level != IPPROTO_TCP) {
306         error = ip_ctloutput(op, so, level, optname, mp);
307         splx(s);
308         return (error);
309     }
310     tp = intotcpcb(inp);

```

Figure 30.15 `tcp_ctloutput` function: first part.

296-303 The processor priority is set to `splnet` while the function executes, and `inp` points to the Internet PCB for the socket. If `inp` is null, the `mbuf` is released if the operation was to set a socket option, and an error is returned.

304-308 If the `level` (the second argument to the `getsockopt` and `setsockopt` system calls) is not `IPPROTO_TCP`, the command is for some other protocol (i.e., IP). For example, it is possible to create a TCP socket and set the IP source routing socket option. In

this example IP processes the socket option, not TCP. `ip_ctloutput` handles the command.

309 The command is for TCP, so `tp` is set to the TCP control block.

The remainder of the function is a switch with two cases: one for `PRCO_SETOPT` (shown in Figure 30.16) and one for `PRCO_GETOPT` (shown in Figure 30.17).

```

310     switch (op) {
311     case PRCO_SETOPT:
312         m = *mp;
313         switch (optname) {
314         case TCP_NODELAY:
315             if (m == NULL || m->m_len < sizeof(int))
316                 error = EINVAL;
317             else if (*ntod(m, int *))
318                 tp->t_flags |= TF_NODELAY;
319             else
320                 tp->t_flags &= ~TF_NODELAY;
321             break;
322         case TCP_MAXSEG:
323             if (m && (i = *ntod(m, int *)) > 0 && i <= tp->t_maxseg)
324                 tp->t_maxseg = i;
325             else
326                 error = EINVAL;
327             break;
328         default:
329             error = ENOPROTOPT;
330             break;
331         }
332         if (m)
333             (void) m_free(m);
334         break;

```

tcp_usrreq.c

tcp_usrreq.c

Figure 30.16 `tcp_ctloutput` function: set a socket option.

325-326 `m` is an mbuf containing the fourth argument to `setsockopt`. For both of the TCP options the mbuf must contain an integer value. If either the mbuf pointer is null, or the amount of data in the mbuf is less than the size of an integer, an error is returned.

TCP_NODELAY option

317-323 If the integer value is nonzero, the `TF_NODELAY` flag is set. This disables the Nagle algorithm in Figure 26.8. If the integer value is 0, the Nagle algorithm is enabled (the default) and the `TF_NODELAY` flag is cleared.

TCP_MAXSEG option

322-327 A process can only decrease the MSS. When a TCP socket is created, `tcp_newtcpcb` initializes `t_maxseg` to its default of 512. When a SYN is received from the other end with an MSS option, `tcp_input` calls `tcp_mss`, and `t_maxseg` can

be set as high as the outgoing interface MTU (minus 40 bytes for the default IP and TCP headers), which is 1460 for an Ethernet. Therefore, after a call to `socket` but before a connection is established, a process can only decrease the MSS from its default of 512. After a connection is established, the process can decrease the MSS from whatever value was selected by `tcp_mss`.

4.4BSD was the first Berkeley release to allow the MSS to be set with a `socket` option. Prior releases only allowed a `getsockopt` for the MSS.

Release mbuf

335-337 The mbuf chain is released.

Figure 30.17 shows the processing for the `PRCO_GETOPT` command.

```

335     case PRCO_GETOPT:
336         *mp = m = m_get(M_WAIT, MT_SOOPTS);
337         m->m_len = sizeof(int);
338
339         switch (optname) {
340             case TCP_NODELAY:
341                 *mtod(m, int *) = tp->t_flags & TF_NODELAY;
342                 break;
343             case TCP_MAXSEG:
344                 *mtod(m, int *) = tp->t_maxseg;
345                 break;
346             default:
347                 error = ENOPROTOPT;
348                 break;
349         }
350         break;
351     }
352     splx(s);
353     return (error);

```

tcp_usrreq.c

tcp_usrreq.c

Figure 30.17 `tcp_ctloutput` function: get a socket option.

335-337 Both TCP socket options return an integer to the process, so `m_get` obtains an mbuf and its length is set to the size of an integer.

339-341 `TCP_NODELAY` returns the current status of the `TF_NODELAY` flag: 0 if the flag is not set (the Nagle algorithm is enabled) or `TF_NODELAY` if the flag is set.

342-344 The `TCP_MAXSEG` option returns the current value of `t_maxseg`. As we said in our discussion of the `PRCO_SETOPT` command, the value returned depends whether the socket has been connected yet.

30.7 Summary

The `tcp_usrreq` function is straightforward because most of the required processing is done by other functions. The `PRU_XXX` requests form the glue between the protocol-independent system calls and the TCP protocol processing.

The `tcp_ctloutput` function is also simple because only two socket options are supported by TCP: enable or disable the Nagle algorithm, and set or fetch the maximum segment size.

Exercises

- 30.1 Now that we've covered all of TCP, list the processing steps and the TCP state transitions when a client goes through the normal steps of `socket`, `connect`, `write` (a request to the server), `read` (a reply from the server), and `close`. Do the same exercise for the server end.
- 30.2 If a process sets the `SO_LINGER` socket option with a linger time of 0 and then calls `close`, we showed how `tcp_disconnect` is called, which causes an RST to be sent. What happens if a process sets this socket option with a linger time of 0 but is then killed by a signal instead of calling `close`? Is the RST segment still sent?
- 30.3 The description for `TCP_LINGERTIME` in Figure 25.4 is the "maximum #seconds for `SO_LINGER` socket option." Given the code in Figure 30.2, is this description correct?
- 30.4 A Net/3 client calls `socket` and `connect` to actively open a connection to a server. The server is reached through the client's default router. A total of 1,129 segments are sent by the client host to the server. Assuming the route to the destination does not change, how many routing table lookups are done on the client host for this connection? Explain.
- 30.5 Obtain the `sock` program described in Appendix C of Volume 1. Run it as a sink server with a pause before reading (`r`) and a large receive buffer. Then run the same program on another system as a source client. Watch the data with `tcpdump`. Verify that TCP's ACK-every-other-segment does not occur and that the only ACKs seen from the server are delayed ACKs.
- 30.6 Modify the `SO_KEEPAIVE` socket option so that the parameters can be configured on a per-connection basis.
- 30.7 Read RFC 1122 to determine why it recommends that an implementation should allow an RST to carry data. Modify the Net/3 code to implement this.

[The page contains extremely faint, illegible text, likely bleed-through from the reverse side of the document. The text is too light to transcribe accurately.]

BPF: BSD Packet Filter

31.1 Introduction

The BSD Packet Filter (BPF) is a software device that “taps” network interfaces. A process accesses a BPF device by opening `/dev/bpf0`, `/dev/bpf1`, and so on. Each BPF device can be opened only by one process at a time.

Since each BPF device allocates 8192 bytes of buffer space, the system administrator typically limits the number of BPF devices. If `open` returns `EBUSY`, the device is in use, and a process tries the next device until the `open` succeeds.

The device is configured with several `ioctl` commands that associate the device with a network interface and install filters to receive incoming packets selectively. Packets are received by reading from the device, and packets are queued on the network interface by writing to the device.

We will use the term *packet* even though *frame* is more accurate, since BPF works at the data-link layer and includes the link-layer headers in the frames it sends and receives.

BPF works only with network interfaces that been modified to support BPF. In Chapter 3 we saw that the Ethernet, SLIP, and loopback drivers call `bpfattach`. This call configures the interface for access through the BPF devices. In this section we show how the BPF device driver is organized and how packets move between the driver and the network interfaces.

BPF is normally used as a diagnostic tool to examine the traffic on a locally attached network. The `tcpdump` program is the best example of such a tool and is described in Appendix A of Volume 1. Normally the user is interested in packets between a given set of machines, or for a particular protocol, or even for a particular TCP connection. A BPF device can be configured with a filter that discards or accepts incoming packets according to a filter specification. Filters are specified as instructions to a pseudo-machine. The details of BPF filters are not discussed in this text. For more information about filters, see `bpf(4)` and [McCanne and Jacobson 1993].

31.2 Code Introduction

The code for the portion of the BPF device driver that we describe resides in the two headers and one C file listed in Figure 31.1.

File	Description
<code>net/bpf.h</code>	BPF constants
<code>net/bpfdesc.h</code>	BPF structures
<code>net/bpf.c</code>	BPF device support

Figure 31.1 Files discussed in this chapter.

Global Variables

The global variables introduced in this chapter are shown in Figure 31.2.

Variable	Datatype	Description
<code>bpf_iflist</code>	<code>struct bpf_if *</code>	linked list of BPF-capable interfaces
<code>bpf_dtab</code>	<code>struct bpf_d []</code>	array of BPF descriptor structures
<code>bpf_bufsize</code>	<code>int</code>	default size of BPF buffers

Figure 31.2 Global variables introduced in this chapter.

Statistics

Figure 31.3 shows the two statistics collected in the `bpf_d` structure for every active BPF device.

<code>bpf_d</code> member	Description
<code>bd_rcount</code>	#packets received from network interface
<code>bd_dcount</code>	#packets dropped because of insufficient buffer space

Figure 31.3 Statistics collected in this chapter.

The remainder of this chapter is divided into four sections:

- BPF interface structures,
- BPF device descriptors,
- BPF input processing, and
- BPF output processing.

31.3 bpf_if Structure

BPF keeps a list of the network interfaces that support BPF. Each interface is described by a `bpf_if` structure, and the global pointer `bpf_iflist` points to the first structure in the list. Figure 31.4 shows a BPF interface structure.

```

67 struct bpf_if {
68     struct bpf_if *bif_next; /* list of all interfaces */
69     struct bpf_d *bif_dlist; /* descriptor list */
70     struct bpf_if **bif_driverp; /* pointer into softc */
71     u_int bif_dlt; /* link layer type */
72     u_int bif_hdrlen; /* length of header (with padding) */
73     struct ifnet *bif_ifp; /* corresponding interface */
74 };

```

bpfdesc.h

Figure 31.4 `bpf_if` structure.

67-69 `bif_next` points to the next BPF interface structure in the list. `bif_dlist` points to a list of BPF devices that have been opened and configured to tap this interface.

70 `bif_driverp` points to a `bpf_if` pointer stored in the `ifnet` structure of the tapped interface. When the interface is *not* tapped, `*bif_driverp` is null. When a BPF device is configured to tap an interface, `*bif_driverp` is changed to point back to the `bif_if` structure and tells the interface to begin passing packets to BPF.

71 The type of interface is saved in `bif_dlt`. The values for our example interfaces are shown in Figure 31.5.

<code>bif_dlt</code>	Description
<code>DLT_EN10MB</code>	10Mb Ethernet interface
<code>DLT_SLIP</code>	SLIP interface
<code>DLT_NULL</code>	loopback interface

Figure 31.5 `bif_dlt` values.

72-74 Each packet accepted by BPF has a BPF header prepended to it. `bif_hdrlen` is the size of the header. Finally, `bif_ifp` points to the `ifnet` structure for the associated interface.

Figure 31.6 shows the `bpf_hdr` structure that is prepended to every incoming packet.

```

122 struct bpf_hdr {
123     struct timeval bh_tstamp; /* time stamp */
124     u_long bh_captlen; /* length of captured portion */
125     u_long bh_datalen; /* original length of packet */
126     u_short bh_hdrlen; /* length of bpf header (this struct plus
127                          alignment padding) */
128 };

```

bpf.h

Figure 31.6 `bpf_hdr` structure.

122-128 `bh_timestamp` records the time the packet was captured. `bh_caplen` is the number of bytes saved by BPF, and `bh_datalen` is the number of bytes in the original packet. `bh_headlen` is the size of the `bpf_hdr` structure plus any padding. This value should match `bif_hdrlen` for the receiving interface and is used by processes to interpret the packets read from the BPF device.

Figure 31.7 shows how `bpf_if` structures are connected to the `ifnet` structures for each of our three sample interfaces (`ie_softc[0]`, `sl_softc[0]`, and `loif`).

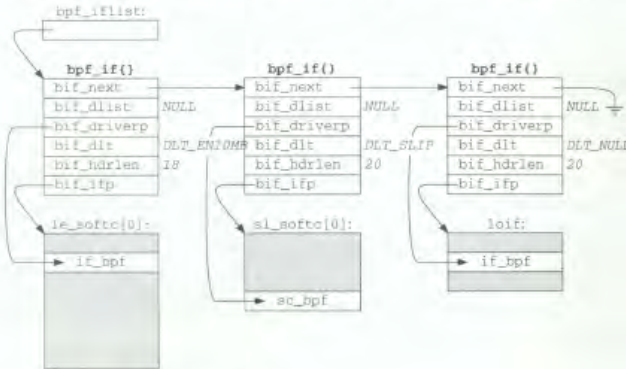


Figure 31.7 `bpf_if` and `ifnet` structures.

Notice that `bif_driverp` points to the `if_bpf` and `sc_bpf` pointers in the network interfaces and *not* to the interface structures.

The SLIP device uses `sc_bpf`, instead of the `if_bpf` member. One reason might be that the SLIP BPF code was written before the `if_bpf` member was added to the `ifnet` structure. The `ifnet` structure in Net/2 does not include a `if_bpf` member.

The link-type and header-length members are initialized for all three interfaces according to the information passed by each driver in the call to `bpfattach`.

In Chapter 3 we saw that `bpfattach` was called by the Ethernet, SLIP, and loopback drivers. The linked list of BPF interface structures is built as each device driver calls `bpfattach` during initialization. The function is shown in Figure 31.8.

1051-1053 `bpfattach` is called by each device driver that supports BPF. The first argument is the pointer saved in `bif_driverp` (described with Figure 31.4). The second argument points to the `ifnet` structure of the interface. The third argument identifies the data-link type, and the fourth argument identifies the size of link-layer header passed with the packet. A new `bpf_if` structure is allocated for the interface.


```

1053 void                                     bpf.c
1054 bpfattach(driverp, ifp, dlt, hdrlen)
1055 caddr_t *driverp;
1056 struct ifnet *ifp;
1057 u_int dlt, hdrlen;
1058 {
1059     struct bpf_if *bp;
1060     int i;
1061     bp = (struct bpf_if *) malloc(sizeof(*bp), M_DEVBUF, M_DONTWAIT);
1062     if (bp == 0)
1063         panic("bpfattach");
1064     bp->bif_dlist = 0;
1065     bp->bif_driverp = (struct bpf_if **) driverp;
1066     bp->bif_ifp = ifp;
1067     bp->bif_dlt = dlt;
1068     bp->bif_next = bpf_iflist;
1069     bpf_iflist = bp;
1070     *bp->bif_driverp = 0;
1071     /*
1072     * Compute the length of the bpf header. This is not necessarily
1073     * equal to sizeof_bpf_hdr because we want to insert spacing such
1074     * that the network layer header begins on a longword boundary (for
1075     * performance reasons and to alleviate alignment restrictions).
1076     */
1077     bp->bif_hdrlen = BPF_WORDALIGN(hdrlen + sizeof_bpf_hdr) - hdrlen;
1078     /*
1079     * Mark all the descriptors free if this hasn't been done.
1080     */
1081     if (!ISPRERE(&bpf_dtab[0]))
1082         for (i = 0; i < NBPFILTER; ++i)
1083             D_MARKFREE(&bpf_dtab[i]);
1084     printf("bpf: %s%d attached\n", ifp->if_name, ifp->if_unit);
1085 }

```

Figure 31.8 bpfattach function.

Initialize bpf_if structure

1064-1070 The `bpf_if` structure is initialized from the arguments and inserted into the front of the BPF interface list, `bpf_iflist`.

Compute BPF header size

1071-1077 `bif_hdrlen` is set to force the *network-layer* header (e.g., the IP header) to start on a longword boundary. This improves performance and avoids unnecessary alignment restrictions for the BPF filter. Figure 31.9 shows the overall organization of the captured BPF packet for each of our three sample interfaces.

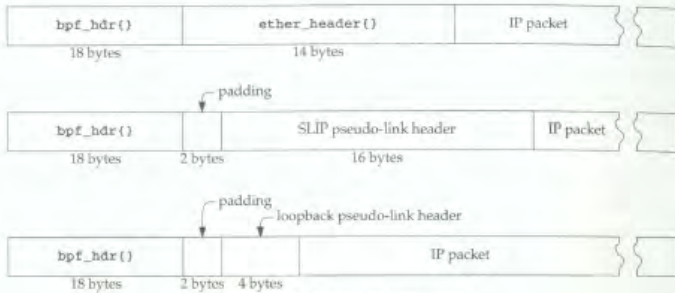


Figure 31.9 BPF packet organization.

The `ether_header` structure was described with Figure 4.10, the SLIP pseudo-link header was described with Figure 5.14, and the loopback pseudo-link header was described with Figure 5.28.

Notice that the SLIP and loopback packets require 2 bytes of padding to force the IP header to appear on a 4-byte boundary.

Initialize `bpf_dtab` table

1076-1083 This code initializes the BPF descriptor table, which is described with Figure 31.10. The initialization occurs the first time `bpfattach` is called and is skipped thereafter.

Print console message

1084-1085 A short message is printed to the console to announce that the interface has been configured for use by BPF.

31.4 `bpf_d` Structure

To begin tapping an interface, a process opens a BPF device and issues `ioctl` commands to select the interface, the read buffer size, and timeouts, and to specify a BPF filter. Each BPF device has an associated `bpf_d` structure, shown in Figure 31.10.

45-46 `bpf_d` structures are placed on a linked list when more than one BPF device is attached to the same network interface. `bd_next` points to the next structure in the list.

Packet buffers

47-52 Each `bpf_d` structure has two packet buffers associated with it. Incoming packets are always stored in the buffer attached to `bd_sbuf` (the store buffer). The other buffer is either attached to `bd_fbuf` (the free buffer), which means it is empty, or to `bd_hbuf` (the hold buffer), which means it contains packets that are being read by a process. `bd_slen` and `bd_hlen` record the number of bytes saved in the store and hold buffer respectively.

```

45 struct bpf_d {
46     struct bpf_d *bd_next; /* linked list of descriptors */
47     caddr_t bd_sbuf; /* store slot */
48     caddr_t bd_hbuf; /* hold slot */
49     caddr_t bd_fbuf; /* free slot */
50     int bd_slen; /* current length of store buffer */
51     int bd_hlen; /* current length of hold buffer */
52
53     int bd_bufsize; /* absolute length of buffers */
54
55     struct bpf_if *bd_bif; /* interface descriptor */
56     u_long bd_rtout; /* read timeout in 'ticks' */
57     struct bpf_insn *bd_filter; /* filter code */
58     u_long bd_rcount; /* number of packets received */
59     u_long bd_dcount; /* number of packets dropped */
60
61     u_char bd_promisc; /* true if listening promiscuously */
62     u_char bd_state; /* idle, waiting, or timed out */
63     u_char bd_immediate; /* true to return on packet arrival */
64     u_char bd_pad; /* explicit alignment */
65     struct selinfo bd_sel; /* bsd select info */
66 };

```

Figure 31.10 bpf_d structure.

When the store buffer becomes full, it is attached to `bd_hbuf` and the free buffer is attached to `bd_sbuf`. When the hold buffer is emptied, it is attached to `bd_fbuf`. The macro `ROTATE_BUFFERS` attaches the store buffer to `bd_hbuf`, attaches the free buffer to `bd_sbuf`, and clears `bd_fbuf`. It is called when the store buffer becomes full, or when the process doesn't want to wait for more packets.

`bd_bufsize` records the size of the two buffers associated with the device. It defaults to 4096 (`BPF_BUF_SIZE`) bytes. The default value can be changed by patching the kernel, or `bd_bufsize` can be changed for a particular BPF device with the `BIOCGBLEN` ioctl command. The `BIOCGBLEN` command returns the current value of `bd_bufsize`, which can never exceed 32768 (`BPF_MAXBUFSIZE`) bytes. There is also a minimum size of 32 (`BPF_MINBUFSIZE`) bytes.

53-57 `bd_bif` points to the `bpf_if` structure associated with the BPF device. The `BIOCSETIF` command specifies the device. `bd_rtout` is the number of clock ticks to delay while waiting for packets to appear. `bd_filter` points to the BPF filter code for this device. Two statistics, which are available to a process through the `BIOCSTATS` command, are kept in `bd_rcount` and `bd_dcount`.

58-63 `bd_promisc` is set with the `BIOCPROMISC` command and causes the interface to operate in promiscuous mode. `bd_state` is unused. `bd_immediate` is set with the `BIOCIMMEDIATE` command and causes the driver to return each packet as it is received instead of waiting for the hold buffer to fill. `bd_pad` pads the `bpf_d` structure to a longword boundary, and `bd_sel` holds the `selinfo` structure for the `select` system call. We don't describe the use of `select` with a BPF device, but `select` itself is described in Section 16.13.

bpffopen Function

When `open` is called for a BPF device, the call is routed to `bpffopen` (Figure 31.11) for processing.

```

256 int
257 bpffopen(dev, flag)
258 dev_t dev;
259 int flag;
260 {
261     struct bpf_d *d;
262     if (minor(dev) >= NBPFILTER)
263         return (ENXIO);
264     /*
265      * Each minor can be opened by only one process. If the requested
266      * minor is in use, return EBUSY.
267      */
268     d = &bpf_dtab[minor(dev)];
269     if (!ID_TOFREE(d))
270         return (EBUSY);
271     /* Mark "free" and do most initialization. */
272     bzero((char *) d, sizeof(*d));
273     d->bd_bufsize = bpf_bufsize;
274     return (0);
275 }

```

Figure 31.11 `bpffopen` function.

255-263 The number of BPF devices is limited at compile time to `NBPFILTER`. The minor device number specifies the device and `ENXIO` is returned if it is too large. This happens when the system administrator creates more `/dev/bpfx` entries than the value `NBPFILTER`.

Allocate `bpf_d` structure

264-275 Only one process is allowed access to a BPF device at a time. If the `bpf_d` structure is already active, `EBUSY` is returned. Programs such as `tcpdump` try the next device when this error is returned. If the device is available, the entry in the `bpf_dtab` table specified by the minor device number is cleared and the size of the packet buffers is set to the default value.

bpfioctl Function

Once the device is opened, it is configured with `ioctl` commands. Figure 31.12 summarizes the `ioctl` commands used with BPF devices. Figure 31.13 shows the `bpfioctl` function. Only the code for `BIOCSETF` and `BIOCSETIF` is shown. We have omitted the `ioctl` commands that are not discussed in this text.

Command	Third argument	Function	Description
<i>FIORREAD</i>	u_int	bpfioctl	return #bytes in hold buffer and store buffers.
<i>BIOCGBLEN</i>	u_int	bpfioctl	return size of packet buffers
<i>BIOCSBLEN</i>	u_int	bpfioctl	set size of packet buffers
<i>BIOCSETF</i>	struct bpf_program	bpf_setf	install BPF program
<i>BIOCFDUSH</i>		reset_d	discard pending packets
<i>BIOCPROMISC</i>		ifpromisc	enable promiscuous mode
<i>BIOCGDLT</i>	u_int	bpfioctl	return bif_dlt
<i>BIOCSETIF</i>	struct ifreq	bpf_ifname	return name of attached interface
<i>BIOCSETIF</i>	struct ifreq	bpf_setif	attach network interface to device
<i>BIOCSRTIMEOUT</i>	struct timeval	bpfioctl	set read timeout value
<i>BIOCGRTIMEOUT</i>	struct timeval	bpfioctl	return read timeout value
<i>BIOCGSTATS</i>	struct bpf_stat	bpfioctl	return BPF statistics
<i>BIOCIMMEDIATE</i>	u_int	bpfioctl	enable immediate mode
<i>BIOCVERSION</i>	struct bpf_version	bpfioctl	return BPF version information

Figure 31.12 BPF ioctl commands.

```

501 int
502 bpfioctl(dev, cmd, addr, flag)
503 dev_t dev;
504 int cmd;
505 caddr_t addr;
506 int flag;
507 {
508     struct bpf_d *d = &bpf_dtab[minor(dev)];
509     int s, error = 0;
510     switch (cmd) {
511         /*
512          * Set link layer read filter.
513          */
514         case BIOCSETF:
515             error = bpf_setf(d, (struct bpf_program *) addr);
516             break;
517         /*
518          * Set interface.
519          */
520         case BIOCSETIF:
521             error = bpf_setif(d, (struct ifreq *) addr);
522             break;
523
524         /* other ioctl commands from Figure 31.12 */
525
526         default:
527             error = EINVAL;
528             break;
529     }
530     return (error);
531 }

```

Figure 31.13 bpfioctl function.

501-509 As with `bpfopen`, the minor device number selects the `bpf_d` structure from the `bpf_dtab` table. The command is processed by the cases within the switch. We show two commands, `BIOCSETF` and `BIOCSETIF`, as well as the default case.

510-522 The `bpf_setf` function installs the filter passed in `addr`, and `bpf_setif` attaches the named interface to the `bpf_d` structure. We don't show the implementation of `bpf_setif` in this text.

558-573 If the command is not recognized, `EINVAL` is returned.

Figure 31.14 shows the `bpf_d` structure after `bpf_setif` has attached it to the LANCE interface in our example system.

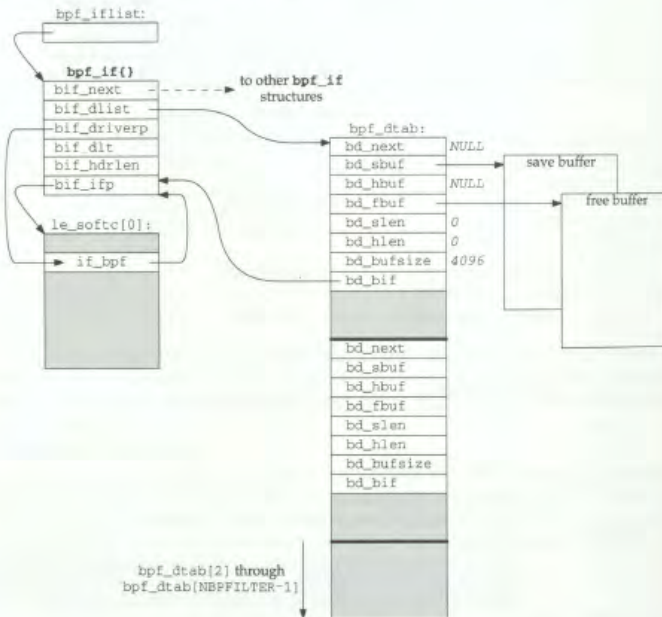


Figure 31.14 BPF device attached to the Ethernet interface.

In the figure, `bif_dlist` points to `bpf_dtab[0]`, the first and only descriptor in the descriptor list for the Ethernet interface. In `bpf_dtab[0]`, the `bd_sbuf` and `bd_hbuf` members point to the store and hold buffers. Each buffer is 4096

(`bd_bufsize`) bytes long. `bd_bif` points back to the `bpf_if` structure for the interface.

`if_bpf` in the `ifnet` structure (`le_softc[0]`) also points back to the `bpf_if` structure. As shown in Figures 4.19 and 4.11, when `if_bpf` is nonnull, the driver begins passing packets to the BPF device by calling `bpf_tap`.

Figure 31.15 shows the same structures after a second BPF device is opened and attached to the same Ethernet network interface as in Figure 31.10.

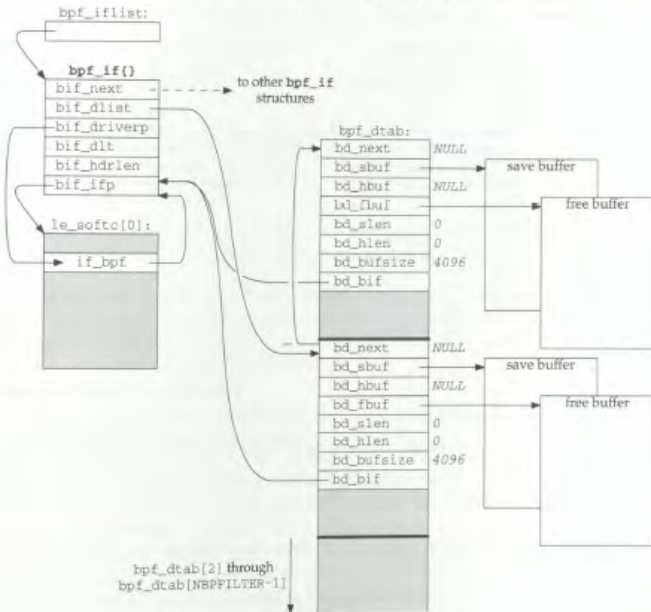


Figure 31.15 Two BPF devices attached to the Ethernet interface.

When the second BPF device is opened, a new `bpf_d` structure is allocated from the `bpf_dtab` table, in this case, `bpf_dtab[1]`. The second BPF device is also attached to the Ethernet interface, so `bif_dlist` points to `bpf_dtab[1]`, and `bpf_dtab[1].bd_next` points to `bpf_dtab[0]`, which is the first BPF descriptor attached to the Ethernet interface. Separate store and hold buffers are allocated and attached to the new descriptor structure.

bpf_setif Function

The `bpf_setif` function, which associates the BPF descriptor with a network interface, is shown in Figure 31.16.

```

721 static int
722 bpf_setif(u, ifr)
723 struct bpf_d *d;
724 struct ifreq *ifr;
725 {
726     struct bpf_if *bp;
727     char *cp;
728     int unit, s, error;

729     /*
730      * Separate string into name part and unit number. Put a null
731      * byte at the end of the name part, and compute the number.
732      * If the a unit number is unspecified, the default is 0,
733      * as initialized above. XXX This should be common code.
734      */
735     unit = 0;
736     cp = ifr->ifr_name;
737     cp[sizeof(ifr->ifr_name) - 1] = '\0';
738     while (*cp++) {
739         if (*cp >= '0' && *cp <= '9') {
740             unit = *cp - '0';
741             *cp++ = '\0';
742             while (*cp)
743                 unit = 10 * unit + *cp++ - '0';
744             break;
745         }
746     }
747     /*
748      * Look through attached interfaces for the named one.
749      */
750     for (bp = bpf_iflist; bp != 0; bp = bp->bif_next) {
751         struct ifnet *ifp = bp->bif_ifp;

752         if ((ifp == 0 || unit != ifp->if_unit
753             || strcmp(ifp->if_name, ifr->ifr_name) != 0)
754             continue;
755         /*
756          * We found the requested interface.
757          * If it's not up, return an error.
758          * Allocate the packet buffers if we need to.
759          * If we're already attached to requested interface,
760          * just flush the buffer.
761          */
762         if ((ifp->if_flags & IFF_UP) == 0)
763             return (ENETDOWN);

```



```

764     if (d->bd_sbuf == 0) {
765         error = bpf_allocbufs(d);
766         if (error != 0)
767             return (error);
768     }
769     s = splimp();
770     if (bp != d->bd_bif) {
771         if (d->bd_bif)
772             /*
773              * Detach if attached to something else.
774              */
775             bpf_detachd(d);
776         bpf_attachd(d, bp);
777     }
778     reset_d(d);
779     splx(s);
780     return (0);
781 }
782 /* Not found. */
783 return (ENXIO);
784 }

```

*bpf.c*Figure 31.16 `bpf_setif` function.

721-746 The first part of `bpf_setif` separates the text portion of the name in the `ifreq` structure (Figure 4.23) from the numeric portion. The numeric portion is saved in `unit`. For example, if the first 4 bytes of `ifr_name` start is "sl1\0", after this code executes they are "sl\0\0" and `unit` is 1.

Locate matching `ifnet` structure

747-754 The `for` loop searches the interfaces that support BPF (the ones in `bpf_iflist`) for the one specified in the `ifreq` structure.

755-768 If the matching interface is not up `ENETDOWN` is returned. If the interface is up, `bpf_allocate` attaches the free and store buffers to the `bpf_d` structure, if they have not already been allocated.

Attach `bpf_d` structure

769-777 If no interface is attached to the BPF device, or if a different interface from the one specified in the `ifreq` structure is attached, `bpf_detachd` discards the previous interface (if any), and `bpf_attachd` attaches the new interface to the device.

778-784 `reset_d` resets the packet buffers, discarding any pending packets in the process. The function returns 0 to indicate success or returns `ENXIO` if the interface was not located.

`bpf_attachd` Function

The `bpf_attachd` function shown in Figure 31.17 associates a BPF descriptor structure with a BPF device and with a network interface.

bpf_setif Function

The `bpf_setif` function, which associates the BPF descriptor with a network interface, is shown in Figure 31.16.

```

721 static int
722 bpf_setif(s, ifr)
723 struct bpf_d *d;
724 struct ifreq *ifr;
725 {
726     struct bpf_if *bp;
727     char *cp;
728     int unit, s, error;
729     /*
730      * Separate string into name part and unit number. Put a null
731      * byte at the end of the name part, and compute the number.
732      * If the a unit number is unspecified, the default is 0,
733      * as initialized above. XXX This should be common code.
734      */
735     unit = 0;
736     cp = ifr->ifr_name;
737     cp[sizeof(ifr->ifr_name) - 1] = '\0';
738     while (*cp++) {
739         if (*cp >= '0' && *cp <= '9') {
740             unit = *cp - '0';
741             *cp++ = '\0';
742             while (*cp)
743                 unit = 10 * unit + *cp++ - '0';
744             break;
745         }
746     }
747     /*
748      * Look through attached interfaces for the named one.
749      */
750     for (bp = bpf_iflist; bp != 0; bp = bp->bif_next) {
751         struct ifnet *ifp = bp->bif_ifp;
752         if (ifp == 0 || unit != ifp->if_unit
753             || strcmp(ifp->if_name, ifr->ifr_name) != 0)
754             continue;
755         /*
756          * We found the requested interface.
757          * If it's not up, return an error.
758          * Allocate the packet buffers if we need to.
759          * If we're already attached to requested interface,
760          * just flush the buffer.
761          */
762         if ((ifp->if_flags & IFP_UP) == 0)
763             return (ENETDOWN);

```

```

764     if (d->bd_sbuf == 0) {
765         error = bpf_allocbufs(d);
766         if (error != 0)
767             return (error);
768     }
769     s = splimp();
770     if (bp != d->bd_bif) {
771         if (d->bd_bif)
772             /*
773              * Detach if attached to something else.
774              */
775             bpf_detachd(d);
776         bpf_attachd(d, bp);
777     }
778     reset_d(d);
779     splx(s);
780     return (0);
781 }
782 /* Not found. */
783 return (ENXIO);
784 }

```

bpf.c

Figure 31.16 bpf_setif function.

722-746 The first part of `bpf_setif` separates the text portion of the name in the `ifreq` structure (Figure 4.23) from the numeric portion. The numeric portion is saved in `unit`. For example, if the first 4 bytes of `ifr_name` start is "sl1\0", after this code executes they are "sl\0\0" and `unit` is 1.

Locate matching `ifnet` structure

747-754 The `for` loop searches the interfaces that support BPF (the ones in `bpf_iflist`) for the one specified in the `ifreq` structure.

755-768 If the matching interface is not up `ENETDOWN` is returned. If the interface is up, `bpf_allocate` attaches the free and store buffers to the `bpf_d` structure, if they have not already been allocated.

Attach `bpf_d` structure

769-777 If no interface is attached to the BPF device, or if a different interface from the one specified in the `ifreq` structure is attached, `bpf_detachd` discards the previous interface (if any), and `bpf_attachd` attaches the new interface to the device.

778-784 `reset_d` resets the packet buffers, discarding any pending packets in the process. The function returns 0 to indicate success or returns `ENXIO` if the interface was not located.

`bpf_attachd` Function

The `bpf_attachd` function shown in Figure 31.17 associates a BPF descriptor structure with a BPF device and with a network interface.

```

189 static void
190 bpf_attachd(d, bp)
191 struct bpf_d *d;
192 struct bpf_if *bp;
193 {
194     /*
195      * Point d at bp, and add d to the interface's list of listeners.
196      * Finally, point the driver's bpf cookie at the interface so
197      * it will divert packets to bpf.
198      */
199     d->bd_bif = bp;
200     d->bd_next = bp->bif_dlist;
201     bp->bif_dlist = d;
202     *bp->bif_driverp = bp;
203 }

```

Figure 31.17 bpf_attachd function.

189-203 First, `bd_bif` is set to point to the BPF interface structure for the network device. Next, the `bpf_d` structure is inserted into the front of the list of `bpf_d` structures associated with the device. Finally, the BPF pointer within the network interface is changed to point to the BPF structure, which causes the interface to begin passing packets to the BPF device.

31.5 BPF Input

Once the BPF device is opened and configured, a process uses the `read` system call to receive packets from the interface. The BPF tap collects *copies* of the incoming packets so BPF does not interfere with normal network processing. Incoming packets are collected in the store and hold buffers associated with each BPF device.

bpf_tap Function

We described the call to `bpf_tap` by the LANCE device driver with Figure 4.11 and use this call to describe the `bpf_tap`. The call (from Figure 4.11) is:

```
bpf_tap(fe->sc_if.if_bpf, buf, len * sizeof(struct ether_header));
```

The `bpf_tap` function is shown in Figure 31.18.

349-382 The first argument is a pointer to the `bpf_if` structure, which is set by `bpfattach`. The second argument is a pointer to the incoming packet, including the Ethernet header. The third argument is the number of bytes contained in the buffer, in this case, the size of the Ethernet header (14 bytes) plus the size of the data portion of the Ethernet frame.


```

869 void
870 bpf_tap(arg, pkt, pktlen)
871 caddr_t arg;
872 u_char *pkt;
873 u_int  pktlen;
874 {
875     struct bpf_if *bp;
876     struct bpf_d *d;
877     u_int  slen;
878     /*
879      * Note that the ipi does not have to be raised at this point.
880      * The only problem that could arise here is that if two different
881      * interfaces shared any data. This is not the case.
882      */
883     bp = (struct bpf_if *) arg;
884     for (d = bp->bif_dlist; d != 0; d = d->bd_next) {
885         ++d->bd_xcount;
886         slen = bpf_filter(d->bd_filter, pkt, pktlen, pktlen);
887         if (slen != 0)
888             catchpacket(d, pkt, pktlen, slen, bcopy);
889     }
890 }

```

Figure 31.18 bpf_tap function.

Pass packet to one or more BPF devices

883-890 The for loop traverses the list of BPF devices attached to the interface. For each device, the packet is passed to `bpf_filter`. If the filter accepts the packet, it returns the number of bytes to capture and `catchpacket` saves a copy of the packet. If the filter rejects the packet, `slen` is 0 and the loop continues. When the loop completes, `bpf_tap` returns. This mechanism enables each BPF device to have a separate filter when multiple BPF devices are associated with the same network interface.

The loopback driver calls `bpf_mtmap` to pass packets to BPF. This function is similar to `bpf_tap` but copies the packet from an mbuf chain instead of from a contiguous area of memory. This function is not described in this text.

catchpacket Function

In Figure 31.18 we saw that `catchpacket` is called when the filter accepts the packet. The function is shown in Figure 31.19.

896-955 The arguments to `catchpacket` are: `d`, a pointer to the BPF device structure; `pkt` a generic pointer to the incoming packet; `pktlen` the length of the packet as it was received; `snplen` the number of bytes to save from the packet; and `cpfn` a pointer to a function that will copy the packet from `pkt` to a contiguous area of memory. When the packet is already in a contiguous area of memory, `cpfn` is `bcopy`. When the packet is stored in an mbuf (i.e., `pkt` points to the first mbuf in a chain such as with the loopback driver), `cpfn` is `bpf_mcopy`.

```

946 static void
947 catchpacket(d, pkt, pktlen, snaplen, cpfn)
948 struct bpf_d *d;
949 u_char *pkt;
950 u_int  pktlen, snaplen;
951 void (*cpfn) (const void *, void *, u_int);
952 {
953     struct bpf_hdr *hp;
954     int  totlen, curlen;
955     int  hdrlen = d->bd_bif->bif_hdrlen;
956     /*
957      * Figure out how many bytes to move.  If the packet is
958      * greater or equal to the snapshot length, transfer that
959      * much.  Otherwise, transfer the whole packet (unless
960      * we hit the buffer size limit).
961      */
962     totlen = hdrlen + min(snaplen, pktlen);
963     if (totlen > d->bd_bufsize)
964         totlen = d->bd_bufsize;
965     /*
966      * Round up the end of the previous packet to the next longword.
967      */
968     curlen = BPF_WORDALIGN(d->bd_slend);
969     if (curlen + totlen > d->bd_bufsize) {
970         /*
971          * This packet will overflow the storage buffer.
972          * Rotate the buffers if we can, then wakeup any
973          * pending reads.
974          */
975         if (d->bd_fbuf == 0) {
976             /*
977              * We haven't completed the previous read yet,
978              * so drop the packet.
979              */
980             ++d->bd_dcount;
981             return;
982         }
983         ROTATE_BUFFERS(d);
984         bpf_wakeup(d);
985         curlen = 0;
986     } else if (d->bd_immediate)
987         /*
988          * Immediate mode is set.  A packet arrived so any
989          * reads should be woken up.
990          */
991         bpf_wakeup(d);
992     /*
993      * Append the bpf header.
994      */
995     hp = (struct bpf_hdr *) (d->bd_sbuf + curlen);
996     microtime(&hp->bh_tstamp);
997     hp->bh_datalen = pktlen;
998     hp->bh_hdrlen = hdrlen;

```

```

999      /*
1000     * Copy the packet data into the store buffer and update its length.
1001     */
1002     (*cpfn) (pkt, (u_char *) bp + hdrlen, (bp->bh_caplen = totlen - hdrlen) );
1003     d->bd_slen = curlen + totlen;
1004 }

```

— *bpf.c*Figure 31.19 `catchpacket` function.

356-864 In addition to the link-layer header and the packet, `catchpacket` appends a `bpf_hdr` to every packet. The number of bytes to save from the packet is the smaller of `snaplen` and `pktlen`. The resulting packet and `bpf_hdr` must fit within the packet buffers (`bd_bufsize` bytes).

Will the packet fit?

365-995 `curlen` is the number of bytes already in the store buffer plus enough bytes to align the next packet on a longword boundary. If the incoming packet doesn't fit in the remaining buffer space, the store buffer is full. If a free buffer is not available (i.e., a process is still reading data from the hold buffer), the incoming packet is discarded. If a free buffer is available, it is rotated into place by `ROTATE_BUFFERS` and any process waiting for incoming data is awakened by `bpf_wakeup`.

Immediate mode processing

996-997 If the device is operating in immediate mode, any waiting processes are awakened to process the incoming packet—there is no buffering of packets in the kernel.

Append BPF header

992-1004 The current time (`microtime`), the packet length, and the header length are saved in a `bpf_hdr`. The function pointed to by `cpfn` is called to copy the packet into the store buffer and the length of the store buffer is updated. Since `bpf_tap` is called directly from `lread` even before the packet is transferred from a device buffer to an `mbuf` chain, the receive timestamp is close to the actual reception time.

bpfread Function

The kernel routes a `read` on a BPF device to `bpfread`. BPF supports a timed read through the `BIOCSRTIMEOUT` command. This “feature” is easily emulated by the more general `select` system call, but `tcpdump`, for example, uses `BIOCSRTIMEOUT` and not `select`. The process must provide a read buffer that matches the size of the hold buffer for the device. The `BIOCGBLEN` command returns the size of the buffer. Normally, a read returns when the store buffer becomes full. The kernel rotates the store buffer to the hold buffer, which is copied to the buffer provided with the `read` system call while the BPF device continues collecting incoming packets in the store buffer. `bpfread` is shown in Figure 31.20.

```

344 int
345 bpfread(dev, uio)
346 dev_t dev;
347 struct uio *uio;
348 {
349     struct bpf_d *d = &bpf_dtab[minor(dev)];
350     int error;
351     int s;
352     /*
353      * Restrict application to use a buffer the same size as
354      * as kernel buffers.
355      */
356     if (uio->uio_resid != d->bd_bufsize)
357         return (EINVAL);
358     s = splimp();
359     /*
360      * If the hold buffer is empty, then do a timed sleep, which
361      * ends when the timeout expires or when enough packets
362      * have arrived to fill the store buffer.
363      */
364     while (d->bd_hbuf == 0) {
365         if (d->bd_immediate && d->bd_slenn != 0) {
366             /*
367              * A packet(s) either arrived since the previous
368              * read or arrived while we were asleep.
369              * Rotate the buffers and return what's here.
370              */
371             ROTATE_BUFFERS(d);
372             break;
373         }
374         error = tsleep((caddr_t) d, PRINET | PCATCH, "bpf", d->bd_rtout);
375         if (error == EINTR || error == ERESTART) {
376             splx(s);
377             return (error);
378         }
379         if (error == EWOULDBLOCK) {
380             /*
381              * On a timeout, return what's in the buffer,
382              * which may be nothing. If there is something
383              * in the store buffer, we can rotate the buffers.
384              */
385             if (d->bd_hbuf)
386                 /*
387                  * We filled up the buffer in between
388                  * getting the timeout and arriving
389                  * here, so we don't need to rotate.
390                  */
391                 break;

```

bpf.c

```

392         if (d->bd_slen == 0) {
393             splx(s);
394             return (0);
395         }
396         ROTATE_BUFFERS(d);
397         break;
398     }
399 }
400 /*
401  * At this point, we know we have something in the hold slot.
402  */
403 splx(s);
404 /*
405  * Move data from hold buffer into user space.
406  * We know the entire buffer is transferred since
407  * we checked above that the read buffer is bpf_bufsize bytes.
408  */
409 error = uiomove(d->bd_hbuf, d->bd_hlen, UIO_READ, uiol);
410 s = uplmp();
411 d->bd_fbuf = d->bd_hbuf;
412 d->bd_hbuf = 0;
413 d->bd_hlen = 0;
414 splx(s);
415 return (error);
416 }

```

—bpf.c

Figure 31.20 bpfread function.

344–357 The minor device number selects the BPF device from the `bpf_dtab` table. If the read buffer doesn't match the size of the BPF device buffers, `EINVAL` is returned.

Wait for data

358–364 Since multiple processes may be reading from the same BPF device, the `while` loop forces the read to continue when some other process gets to the data first. If there is data in the hold buffer, the loop is skipped. This is different from two processes tapping the same network interface through two different BPF devices (Exercise 31.2).

Immediate mode

365–373 If the device is in immediate mode and there is some data in the store buffer, the buffers are rotated and the `while` loop terminates.

No packets available

374–384 If the device is not in the immediate mode, or there is no data in the store buffer, the process sleeps until a signal arrives, the read timer expires, or data arrives in the hold buffer. If a signal arrives, `EINTR` or `ERESTART` is returned.

Remember that a process never sees the `ERESTART` error because the error is handled by the `syscall` function and never returned to a process.

Check hold buffer

385-391 If the timer expired and data is in the hold buffer, the loop terminates.

Check store buffer

392-399 If the timer expired and there is no data in the store buffer, the read returns 0. The process must handle this case when using a timed read. If the timer expired and there is data in the store buffer, it is rotated to the hold buffer and the loop terminates.

If `tsleep` returns without an error and data is present, the `while` loop test is false and the loop terminates.

Packets are available

400-418 At this point, there is data in the hold buffer. `uiomove` moves `bd_hlen` bytes of data from the hold buffer to the process. After the move, the hold buffer is moved to the free buffer, and the buffer counts are cleared before the function returns. The comment before `uiomove` indicates that `uiomove` will always be able to copy `bd_hlen` bytes into the process because the read buffer was checked to ensure it can hold the maximum number of bytes, `bd_bufsize`.

31.6 BPF Output

Finally, we describe how to add packets to the network interface output queues with BPF. An entire data-link frame must be constructed by the process. For Ethernet this includes the source and destination hardware addresses and the frame type (Figure 4.8). The kernel will not modify the frame before putting it on the interface's output queue.

`bpfwrite` Function

The frame is passed to the BPF device with the `write` system call, which the kernel routes to `bpfwrite`, shown in Figure 31.21.

Check device number

437-449 The minor device number selects the BPF device, which must be attached to a network interface. If it isn't, `ENXIO` is returned.

Copy data into mbuf chain

450-457 If the write specified 0 bytes, 0 is returned immediately. `bpf_movein` copies the data from the process into an mbuf chain. Based on the interface type passed from `bif_dlt`, it computes the length of the packet excluding the link-layer header and returns the value in `datlen`. It also returns an initialized `sockaddr` structure in `dst`. For Ethernet, the type of this address structure will be `AF_UNSPEC`, indicating that the mbuf chain contains the data-link header for the outgoing frame. If the packet is larger than the MTU of the interface, `EMSGSIZE` is returned.

Queue packet

458-465 The resulting mbuf chain is passed to the network interface using the `if_output` function specified in the `ifnet` structure. For Ethernet, `if_output` is `ether_output`.

```
437 int bpf.c
438 bpfwrite(dev, uio)
439 dev_t dev;
440 struct uio *uio;
441 {
442     struct bpf_d *d = &bpf_dtab[minor(dev)];
443     struct ifnet *ifp;
444     struct mbuf *m;
445     int error, s;
446     static struct sockaddr dst;
447     int datlen;
448
449     if (d->bd_bif == 0)
450         return (ENXIO);
451
452     ifp = d->bd_bif->bif_ifp;
453
454     if (uio->uio_resid == 0)
455         return (0);
456
457     error = bpf_movein(uio, (int) d->bd_bif->bif_dlt, &m, &dot, &datlen);
458     if (error)
459         return (error);
460
461     if (datlen > ifp->if_mtu)
462         return (EMSGSIZE);
463
464     s = splnet();
465     error = (*ifp->if_output) (ifp, m, &dst, (struct rentry *) 0);
466     splx(s);
467     /*
468      * The driver frees the mbuf.
469      */
470     return (error);
471 }
```

Figure 31.21 bpfwrite function.

31.7 Summary

In this chapter we showed how BPF devices are configured, how incoming frames are passed to BPF devices, and how outgoing frames can be transmitted on a BPF device.

We showed that a single network interface can have multiple BPF taps, each with a separate filter. The store and hold buffers minimize the number of read system calls required to process incoming frames.

We focused only on the major features of BPF in this chapter. For a more detailed description of the filtering code and the other features of the BPF device, the interested reader should examine the source code and the Net/3 manual pages.

Exercises

- 31.1 Why is it OK to call `bpf_wakeup` in `catchpacket` before the packet is stored in the BPF buffers?
- 31.2 With Figure 31.20, we noted that two processes may be waiting for data from the same BPF device. With Figure 31.11, we noted that only one process at a time can open a particular BPF device. How can both of these statements be true?
- 31.3 What happens if the device named in the `BIOCSETIF` command does not support BPF?

Raw IP

32.1 Introduction

A process accesses the raw IP layer by creating a socket of type `SOCK_RAW` in the Internet domain. There are three uses for raw sockets:

1. Raw sockets allow a process to send and receive ICMP and IGMP messages.

The Ping program uses this type of socket to send ICMP echo requests and to receive ICMP echo replies.

Some routing daemons use this feature to track ICMP redirects that are processed by the kernel. We saw in Section 19.7 that Net/3 generates an `RTM_REDIRECT` message on a routing socket when a redirect is processed, obviating the need for this use of raw sockets.

This feature is also used to implement protocols based on ICMP, such as router advertisement and router solicitation (Section 9.6 of Volume 1), which use ICMP but are better implemented as user processes than within the kernel.

The multicast routing daemon uses a raw IGMP socket to send and receive IGMP messages.

2. Raw sockets let a process build its own IP headers. The Traceroute program uses this feature to build its own UDP datagrams, including the IP and UDP headers.

- Raw sockets let a process read and write IP datagrams with an IP protocol type that the kernel doesn't support.

The `gated` program uses this to support three routing protocols that are built directly on IP: EGP, HELLO, and OSPF.

This type of raw socket can also be used to experiment with new transport layers on top of IP, instead of adding support to the kernel. It is usually much easier to debug code within a user process than it is within the kernel.

This chapter examines the implementation of raw IP sockets.

32.2 Code Introduction

There are five raw IP functions in a single C file, shown in Figure 32.1.

File	Description
<code>netinet/raw_ip.c</code>	raw IP functions

Figure 32.1 File discussed in this chapter.

Figure 32.2 shows the relationship of the five raw IP functions to other kernel functions.

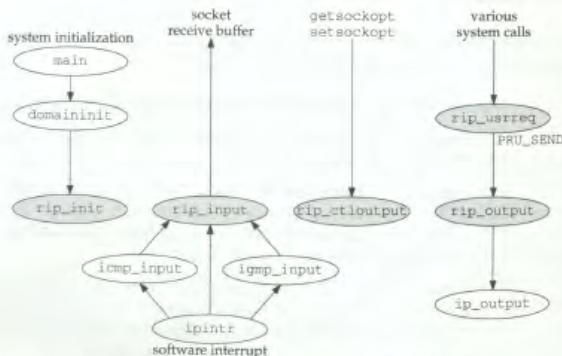


Figure 32.2 Relationship of raw IP functions to rest of kernel.

The shaded ellipses are the five functions that we cover in this chapter. Be aware that the “rip” prefix used within the raw IP functions stands for “raw IP” and not the “Routing Information Protocol,” whose common acronym is RIP.

Global Variables

Four global variables are introduced in this chapter, which are shown in Figure 32.3.

Variable	Datatype	Description
rawinpcb	struct inpcb	head of the raw IP Internet PCB list
ripsrc	struct sockaddr_in	contains sender's IP address on input
rip_recvspace	u_long	default size of socket receive buffer, 8192 bytes
rip_sendspace	u_long	default size of socket send buffer, 8192 bytes

Figure 32.3 Global variables introduced in this chapter.

Statistics

Raw IP maintains two of the counters in the `ipstat` structure (Figure 8.4). We describe these in Figure 32.4.

ipstat member	Description	Used by SNMP
<code>ips_noproto</code>	#packets with an unknown or unsupported protocol	•
<code>ips_rawout</code>	total #raw ip packets generated	

Figure 32.4 Raw IP statistics maintained in the `ipstat` structure.

The use of the `ips_noproto` counter with SNMP is shown in Figure 8.6. Figure 8.5 shows some sample output of these two counters.

32.3 Raw IP protosw Structure

Unlike all other protocols, raw IP is accessed through multiple entries in the `inetsw` array. There are four entries in this structure with a socket type of `SOCK_RAW`, each with a different protocol value:

- `IPPROTO_ICMP` (protocol value of 1),
- `IPPROTO_IGMP` (protocol value of 2),
- `IPPROTO_RAW` (protocol value of 255), and
- raw wildcard entry (protocol value of 0).

The first two entries for ICMP and IGMP were described earlier (Figures 11.12 and 13.9). The difference in these four entries can be summarized as follows:

- If the process creates a raw socket (`SOCK_RAW`) with a nonzero protocol value (the third argument to `socket`), and if that value matches `IPPROTO_ICMP`, `IPPROTO_IGMP`, or `IPPROTO_RAW`, then the corresponding `protosw` entry is used.

- If the process creates a raw socket with a nonzero protocol value that is not known to the kernel, the wildcard entry with a protocol of 0 is matched by `pffindproto`. This allows a process to handle any IP protocol that is not known to the kernel, without making kernel modifications.

We saw in Section 7.8 that all entries in the `ip_protosw` array that are unknown are set to point to the entry for `IPPROTO_RAW`, whose protocol switch entry we show in Figure 32.5.

Member	inetsw(31)	Description
<code>pr_type</code>	<code>SOCK_RAW</code>	raw socket
<code>pr_domain</code>	<code>INETDOMAIN</code>	raw IP is part of the Internet domain
<code>pr_protocol</code>	<code>IPPROTO_RAW (255)</code>	appears in the <code>ip_p</code> field of the IP header
<code>pr_flags</code>	<code>PR_ATOMIC PR_ADDR</code>	socket layer flags, not used by protocol processing
<code>pr_input</code>	<code>rip_input</code>	receives messages from IP layer
<code>pr_output</code>	0	not used by raw IP
<code>pr_ctlinput</code>	0	not used by raw IP
<code>pr_ctloutput</code>	<code>rip_ctloutput</code>	respond to administrative requests from a process
<code>pr_usrreq</code>	<code>rip_usrreq</code>	respond to communication requests from a process
<code>pr_init</code>	0	not used by raw IP
<code>pr_fasttimo</code>	0	not used by raw IP
<code>pr_slowtimo</code>	0	not used by raw IP
<code>pr_drain</code>	0	not used by raw IP
<code>pr_sysctl</code>	0	not used by raw IP

Figure 32.5 The raw IP `protosw` structure.

We describe the three functions that begin with `rip_` in this chapter. We also cover the function `rip_output`, which is not in the protocol switch entry but is called by `rip_usrreq` when a raw IP datagram is output.

The fifth raw IP function, `rip_init`, is contained only in the wildcard entry. The initialization function must be called only once, so it could appear in either the `IPPROTO_RAW` entry or in the wildcard entry.

What Figure 32.5 doesn't show, however, is that other protocols (ICMP and IGMP) also reference some of the raw IP functions in their `protosw` entries. Figure 32.6 compares the relevant fields in the `protosw` entries for the four `SOCK_RAW` protocols. To highlight the differences, values in these rows are in a bolder font when they differ.

protosw entry	SOCK_RAW protocol type			
	IPPROTO_ICMP (1)	IPPROTO_IGMP (2)	IPPROTO_RAW (255)	wildcard (0)
<code>pr_input</code>	<code>icmp_input</code>	<code>igmp_input</code>	<code>rip_input</code>	<code>rip_input</code>
<code>pr_output</code>	<code>rip_output</code>	<code>rip_output</code>	<code>rip_output</code>	<code>rip_output</code>
<code>pr_ctloutput</code>	<code>rip_ctloutput</code>	<code>rip_ctloutput</code>	<code>rip_ctloutput</code>	<code>rip_ctloutput</code>
<code>pr_usrreq</code>	<code>rip_usrreq</code>	<code>rip_usrreq</code>	<code>rip_usrreq</code>	<code>rip_usrreq</code>
<code>pr_init</code>	0	<code>igmp_init</code>	0	<code>rip_init</code>
<code>pr_sysctl</code>	<code>icmp_sysctl</code>	0	0	0
<code>pr_fasttimo</code>	0	<code>igmp_fasttimo</code>	0	0

Figure 32.6 Comparison of protocol switch values for raw sockets.

The implementation of raw sockets has changed with the different BSD releases. The entry with a protocol of `IPPROTO_RAW` has always been used as the wildcard entry in the `ip_protocx` table for unknown IP protocols. The entry with a protocol of 0 has always been the default entry, to allow processes to read and write IP datagrams with a protocol that the kernel doesn't support.

Usage of the `IPPROTO_RAW` entry by a process started when Traceroute was developed by Van Jacobson, because Traceroute was the first process that needed to write its own IP headers (to change the TTL field). The kernel patches to 4.3BSD and Net/1 to support Traceroute included a change to `rip_output` so that if the protocol was `IPPROTO_RAW`, it was assumed the process had passed a complete IP datagram, including the IP header. This was changed with Net/2 when the `IP_HDRINCL` socket option was introduced, removing this overloading of the `IPPROTO_RAW` protocol and allowing a process to send its own IP header with the wildcard entry.

32.4 rip_init Function

The `domaininit` function calls the raw IP initialization function `rip_init` (Figure 32.7) at system initialization time.

```

47 void
48 rip_init()
49 {
50     rawinpcb.inp_next = rawinpcb.inp_prev = &rawinpcb;
51 }

```

raw_ip.c

raw_ip.c

Figure 32.7 `rip_init` function.

The only action performed by this function is to set the next and previous pointers in the head PCB (`rawinpcb`) to point to itself. This is an empty doubly linked list.

Whenever a socket of type `SOCK_RAW` is created by the `socket` system call, we'll see that the raw IP `PRU_ATTACH` function creates an Internet PCB and puts it onto the `rawinpcb` list.

32.5 rip_input Function

Since all entries in the `ip_protocx` array for unknown protocols are set to point to the entry for `IPPROTO_RAW` (Section 7.8), and since the `pr_input` function for this protocol is `rip_input` (Figure 32.6), this function is called for all IP datagrams that have a protocol value that the kernel doesn't recognize. But from Figure 32.2 we see that both ICMP and IGMP also call `rip_input`. This happens under the following conditions:

- `icmp_input` calls `rip_input` for all unknown ICMP message types and for all ICMP messages that are not reflected.
- `igmp_input` calls `rip_input` for all IGMP packets.

One reason for calling `rip_input` in these two cases is to allow a process with a raw socket to handle new ICMP and IGMP messages that might not be supported by the kernel.

Figure 32.8 shows the `rip_input` function.

```

59 void                                     raw_ip.c
60 rip_input(m)
61 struct mbuf *m;
62 {
63     struct ip *ip = mtod(m, struct ip *);
64     struct inpcb *inp;
65     struct socket *last = 0;
66
67     ripsrc.sin_addr = ip->ip_src;
68     for (inp = rawinpcb.inp_next; inp != &rawinpcb; inp = inp->inp_next) {
69         if (inp->inp_ip.ip_p && inp->inp_ip.ip_p != ip->ip_p)
70             continue;
71         if (inp->inp_laddr.s_addr &&
72             inp->inp_laddr.s_addr == ip->ip_dst.s_addr)
73             continue;
74         if (inp->inp_faddr.s_addr &&
75             inp->inp_faddr.s_addr == ip->ip_src.s_addr)
76             continue;
77         if (last) {
78             struct mbuf *n;
79             if (n = m_copy(m, 0, (int) M_COPYALL)) {
80                 if (sbappendaddr(&last->so_rcv, &riprsrc,
81                                 n, (struct mbuf *) 0) == 0)
82                     /* should notify about lost packet */
83                     m_freem(n);
84                 else
85                     sorwakeup(last);
86             }
87             last = inp->inp_socket;
88         }
89         if (last) {
90             if (sbappendaddr(&last->so_rcv, &riprsrc,
91                             m, (struct mbuf *) 0) == 0)
92                 m_freem(m);
93             else
94                 sorwakeup(last);
95         } else {
96             m_freem(m);
97             ipstat.ips_noproto++;
98             ipstat.ips_delivered--;
99         }
100 }

```

Figure 32.8 `rip_input` function.

Save source IP address

89-95 The source address from the IP datagram is put into the global variable `rip_src`, which becomes an argument to `sbappendaddr` whenever a matching PCB is found. Unlike UDP, there is no concept of a port number with raw IP, so the `sin_port` field in the `sockaddr_in` structure is always 0.

Search all raw IP PCBs for one or more matching entries

97-98 Raw IP handles its list of PCBs differently from UDP and TCP. We saw that these two protocols maintain a pointer to the PCB for the most recently received datagram (a one-behind cache) and call the generic function `in_pcblookup` to search for a single “best” match when the received datagram does not equal the cache entry. Raw IP has completely different criteria for a matching PCB, so it searches the PCB list itself. `in_pcblookup` cannot be used because a raw IP datagram can be delivered to multiple sockets, so every PCB on the raw PCB list must be scanned. This is similar to UDP’s handling of a received datagram destined for a broadcast or multicast address (Figure 23.26).

Compare protocols

68-69 If the protocol field in the PCB is nonzero, and if it doesn’t match the protocol field in the IP header, the PCB is ignored. This implies that a raw socket with a protocol value of 0 (the third argument to `socket`) can match any received raw IP datagram.

Compare local and foreign IP addresses

79-75 If the local address in the PCB is nonzero, and if it doesn’t match the destination IP address in the IP header, the PCB is ignored. If the foreign address in the PCB is nonzero, and if it doesn’t match the source IP address in the IP header, the PCB is ignored.

These three tests imply that a process can create a raw socket with a protocol of 0, not bind a local address, and not connect to a foreign address, and the process receives *all* datagrams processed by `rip_input`.

Pass copy of received datagram to processes

76-84 `sbappendaddr` passes a copy of the received datagram to the process. The use of the variable `last` is similar to what we saw in Figure 23.26: since `sbappendaddr` releases the mbuf after placing it onto the appropriate queue, if more than one process receives a copy of the datagram, `rip_input` must make a copy by calling `m_copy`. But if only one process receives the datagram, there’s no need to make a copy.

Undeliverable datagram

85-89 If no matching sockets are found for the datagram, the mbuf is released, `ips_noproto` is incremented, and `ips_delivered` is decremented. This latter counter was incremented by IP just before calling the `rip_input` (Figure 8.15). It must be decremented so that the two SNMP counters, `ipInDiscards` and `ipInDelivers` (Figure 8.6) are correct, since the datagram was not really delivered to a transport layer.

At the beginning of this section we mentioned that `icmp_input` calls `rip_input` for unknown message types and for messages that are not reflected. This means that the receipt of an ICMP host unreachable causes `ips_noproto` to be incremented if there are no raw listeners whose PCB is matched by `rip_input`. That's one reason this counter has such a large value in Figure 8.5. The description of this counter as being "unknown or unsupported protocols" is not entirely accurate.

Net/3 does not generate an ICMP destination unreachable message with code 2 (protocol unreachable) when an IP datagram is received with a protocol field that is not handled by either the kernel or some process through a raw socket. RFC 1122 says an implementation should generate this ICMP error. (See Exercise 32.4.)

32.6 `rip_output` Function

We saw in Figure 32.6 that `rip_output` is called for output for raw sockets by ICMP, IGMP, and raw IP. Output occurs when the application calls one of the five write functions: `send`, `sendto`, `sendmsg`, `write`, or `writew`. 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 destination address must be specified.

The function `rip_output` is shown in Figure 32.9.

Kernel fills in IP header

119-128 If the `IP_HDRINCL` socket option is not defined, `M_PREPEND` allocates room for an IP header, and fields in the IP header are filled in. The fields that are not filled in here are left for `ip_output` to initialize (Figure 8.22). The protocol field is set to the value stored in the PCB, which we'll see in Figure 32.10 is the third argument to the `socket` system call.

The TOS is set to 0 and the TTL to 255. These values are always used for a raw socket when the kernel fills in the header. This differs from UDP and TCP where the process had the capability of setting the `IP_TTL` and `IP_TOS` socket options.

129 Any IP options set by the process with the `IP_OPTIONS` socket options are passed to `ip_output` through the `opts` variable.

Caller fills in IP header: `IP_HDRINCL` socket option

130-133 If the `IP_HDRINCL` socket option is set, the caller supplies a completed IP header at the front of the datagram. The only modification made to this IP header is to set the ID field if the value supplied by the process is 0. The ID field of an IP datagram can be 0. The assignment of the ID field here by `rip_output` is just a convention that allows the process to set it to 0, asking the kernel to assign an ID value based on the kernel's current `ip_id` variable.

134-136 The `opts` variable is set to a null pointer, which ignores any IP options the process may have set with the `IP_OPTIONS` socket option. The convention here is that if the caller builds its own IP header, that header includes any IP options the caller might want. The `flags` variable must also include the `IP_RAWOUTPUT` flag, telling `ip_output` to leave the header alone.

```

105 int                                     raw_ip.c
106 rip_output(m, so, dst)
107 struct mbuf *m;
108 struct socket *so;
109 u_long dst;
110 {
111     struct ip *ip;
112     struct inpcb *inp = sotoinpcb(so);
113     struct mbuf *opts;
114     int flags = (so->so_options & SO_DONTROUTE) | IP_ALLOWBROADCAST;
115     /*
116      * If the user handed us a complete IP packet, use it.
117      * Otherwise, allocate an mbuf for a header and fill it in.
118      */
119     if ((inp->inp_flags & INP_HDRINCL) == 0) {
120         M_PREPEND(m, sizeof(struct ip), M_WAIT);
121         ip = mtod(m, struct ip *);
122         ip->ip_tos = 0;
123         ip->ip_off = 0;
124         ip->ip_p = inp->inp_ip.ip_p;
125         ip->ip_len = m->m_pkthdr.len;
126         ip->ip_src = inp->inp_laddr;
127         ip->ip_dst.s_addr = dst;
128         ip->ip_ttl = MAXTTL;
129         opts = inp->inp_options;
130     } else {
131         ip = mtod(m, struct ip *);
132         if (ip->ip_id == 0)
133             ip->ip_id = htons(ip_id++);
134         opts = NULL;
135         /* XXX prevent ip_output from overwriting header fields */
136         flags |= IP_RAWOUTPUT;
137         ipstat.ips_rawout++;
138     }
139     return (ip_output(m, opts, kinp->inp_route, flags, inp->inp_options));
140 }

```

Figure 32.9 rip_output function.

137 The counter `ips_rawout` is incremented. Running Traceroute causes this variable to be incremented by 1 for each datagram sent by Traceroute.

The operation of `rip_output` has changed over time. When the `IP_HDRINCL` socket option is used in Net/3, the only change made to the IP header by `rip_output` is to set the ID field, if the process sets it to 0. The Net/3 `ip_output` function does nothing to the IP header fields because the `IP_RAWOUTPUT` flag is set. Net/2, however, always set certain fields in the IP header, even if the `IP_HDRINCL` socket option was set; the IP version was set to 4, the fragment offset was set to 0, and the more-fragments flag was cleared.

32.7 rip_usrreq Function

The protocol's user-request function is called for a variety of operations. As with the UDP and TCP user-request functions, `rip_usrreq` is a large `switch` statement, with one `case` for each `PRU_XXX` request.

The `PRU_ATTACH` request, shown in Figure 32.10, is from the `socket` system call

```

194 int                                     raw_ip.c
195 rip_usrreq(so, req, m, nam, control)
196 struct socket *so;
197 int req;
198 struct mbuf *m, *nam, *control;
199 {
200     int error = 0;
201     struct inpcb *inp = sotoinpcb(so);
202     extern struct socket *ip_mrouter;
203     switch (req) {
204     case PRU_ATTACH:
205         if (inp)
206             panic("rip_attach");
207         if ((so->so_state & SS_PRIV) == 0) {
208             error = EACCES;
209             break;
210         }
211         if ((error = soreserve(so, rip_sendspace, rip_recvspace)) ||
212             (error = in_pcballoc(so, &rawinpcb)))
213             break;
214         inp = (struct inpcb *) so->so_pcb;
215         inp->inp_ip.ip_p = (int) nam;
216         break;

```

Figure 32.10 `rip_usrreq` function: `PRU_ATTACH` request.

194-206 Since the `socket` function creates a new `socket` structure each time it is called, that structure cannot point to an Internet PCB.

Verify superuser

207-210 Only the superuser can create a raw socket. This is to prevent random users from writing their own IP datagrams to the network.

Create Internet PCB and reserve buffer space

211-215 Space is reserved for input and output queues, and `in_pcballoc` allocates a new Internet PCB. The PCB is added to the raw IP PCB list (`rawinpcb`). The PCB is linked to the `socket` structure. The `nam` argument to `rip_usrreq` is the third argument to the `socket` system call: the protocol. It is stored in the PCB since it is used by `rip_input` to demultiplex received datagrams, and its value is placed into the protocol field of outgoing datagrams by `rip_output` (if `IP_HDRINCL` is not set).

A raw IP socket can be connected to a foreign IP address similar to a UDP socket being connected to a foreign IP address. This fixes the foreign IP address from which the raw socket receives datagrams, as we saw in `rip_input`. Since raw IP is a

connectionless protocol like UDP, a PRU_DISCONNECT request can occur in two cases:

1. When a connected raw socket is closed, PRU_DISCONNECT is called before PRU_DETACH.
2. When a connect is issued on an already-connected raw socket, soconnect issues the PRU_DISCONNECT request before the PRU_CONNECT request.

Figure 32.11 shows the PRU_DISCONNECT, PRU_ABORT, and PRU_DETACH requests.

```

217     case PRU_DISCONNECT:                                raw_ip.c
218         if ((so->so_state & SS_ISCONNECTED) == 0) {
219             error = ENOTCONN;
220             break;
221         }
222         /* FALLTHROUGH */
223     case PRU_ABORT:
224         oosdisonconnected(so);
225         /* FALLTHROUGH */
226     case PRU_DETACH:
227         if (inp == 0)
228             panic("rip_detach");
229         if (so == ip_mrouter)
230             ip_mrouter_done();
231         in_pcbdetach(inp);
232         break;

```

Figure 32.11 rip_usrreq function: PRU_DISCONNECT, PRU_ABORT, and PRU_DETACH requests.

717-722 The socket must already be connected to disconnect or else an error is returned.
 723-725 A PRU_ABORT abort should never be issued for a raw IP socket, but this case also handles the fall through from PRU_DISCONNECT. The socket is marked as disconnected.

226-230 The close system call issues the PRU_DETACH request, and this case also handles the fall through from the PRU_DISCONNECT request. If the socket structure is the one used for multicast routing (ip_mrouter), multicast routing is disabled by calling ip_mrouter_done. Normally the mouted(8) daemon issues the DVMRP_DONE socket option to disable multicast routing, so this check handles the case of the router daemon terminating (i.e., crashing) without issuing the socket option.

##1 The Internet PCB is released by in_pcbdetach, which also removes the PCB from the list of raw IP PCBs (rawinpcb).

A raw IP socket can be bound to a local IP address with the PRU_BIND request, shown in Figure 32.12. We saw in rip_input that the socket will receive only datagrams sent to this IP address.

213-250 The process fills in a sockaddr_in structure with the local IP address. The following three conditions must all be true, or else the error EADDRNOTAVAIL is returned:

```

233     case PRU_BIND:
234     {
235         struct sockaddr_in *addr = mtod(nam, struct sockaddr_in *);
236
237         if (nam->m_len != sizeof(*addr)) {
238             error = EINVAL;
239             break;
240         }
241         if ((ifnet == 0) ||
242             ((addr->sin_family != AF_INET) &&
243              (addr->sin_family != AF_IMPLINK))) {
244             (addr->sin_addr.s_addr &&
245              ifa_ifwithaddr((struct sockaddr *) addr) == 0)) {
246                 error = EADDRNOTAVAIL;
247                 break;
248             }
249             inp->inp_laddr = addr->sin_addr;
250             break;
251     }

```

Figure 32.12 `rip_usrreq` function: `PRU_BIND` request.

1. at least one IP interface must be configured,
2. the address family must be `AF_INET` (or `AF_IMPLINK`, a historical artifact), and
3. if the IP address being bound is not 0.0.0.0, it must correspond to a local interface. For the call to `ifa_ifwithaddr` to succeed, the port number in the caller's `sockaddr_in` must be 0.

The local IP address is stored in the PCB.

A process can also connect a raw IP socket to a particular foreign IP address. We saw in `rip_input` that this restricts the process so that it receives only IP datagrams with a source IP address equal to the connected IP address. A process has the option of calling `bind`, `connect`, both, or neither, depending on the type of filtering it wants `rip_input` to place on received datagrams. Figure 32.13 shows the `PRU_CONNECT` request.

231-230 If the caller's `sockaddr_in` is initialized correctly and at least one IP interface is configured, the specified foreign IP address is stored in the PCB. Notice that this process differs from the connection of a UDP socket to a foreign address. In the UDP case, `in_pcbconnect` acquires a route to the foreign address and also stores the outgoing interface as the local address (Figure 22.9). With raw IP, only the foreign IP address is stored in the PCB, and unless the process also calls `bind`, only the foreign address is compared by `rip_input`.


```

251     case PRU_CONNECT:
252     {
253         struct sockadr_in *addr = mtod(nam, struct sockadr_in *);
254         if (nam->m_len != sizeof(*addr)) {
255             error = EINVAL;
256             break;
257         }
258         if (ifnet == 0) {
259             error = EADDRNOTAVAIL;
260             break;
261         }
262         if ((addr->sin_family != AF_INET) &&
263             (addr->sin_family != AF_IMPLINK)) {
264             error = EAFNOSUPPORT;
265             break;
266         }
267         inp->inp_faddr = addr->sin_addr;
268         socntconnected(sol);
269         break;
270     }

```

raw_ip.c

Figure 32.13 rip_usrreq function: PRU_CONNECT request.

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 raw IP socket. Figure 32.14 shows the PRU_CONNECT2 and PRU_SHUTDOWN requests.

```

271     case PRU_CONNECT2:
272         error = EOPNOTSUPP;
273         break;
274
275         /*
276          * Mark the connection as being incapable of further input.
277          */
278     case PRU_SHUTDOWN:
279         socntsendmore(sol);
280         break;

```

raw_ip.c

Figure 32.14 rip_usrreq function: PRU_CONNECT2 and PRU_SHUTDOWN requests.

271-273 The PRU_CONNECT2 request is not supported for a raw IP socket.
 274-279 socntsendmore sets the socket's flags to prevent any future output.

In Figure 23.14 we showed how the five write functions call the protocol's `pr_usrreq` function with a `PRU_SEND` request. We show this request in Figure 32.15.

```

280      /*
281      * Ship a packet out. The appropriate raw output
282      * routine handles any massaging necessary.
283      */
284      case PRU_SEND:
285      {
286          u_long dat;
287
288          if (so->so_state & SS_ISCONNECTED) {
289              if (nam) {
290                  error = EISCONN;
291                  break;
292              }
293              dst = inp->inp_faddr.s_addr;
294          } else {
295              if (nam == NULL) {
296                  error = ENOTCONN;
297                  break;
298              }
299              dst = mtod(nam, struct sockaddr_in *)->sin_addr.s_addr;
300          }
301          error = rip_output(m, so, dst);
302          m = NULL;
303          break;
304      }
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000

```

Figure 32.15 `rip_usrreq` function: `PRU_SEND` request.

390-393 If the socket state is connected, the caller cannot specify a destination address (the `nam` argument). Likewise, if the state is unconnected, a destination address is required. If all is OK, in either state, `dst` is set to the destination IP address. `rip_output` sends the datagram. The mbuf pointer `m` is set to a null pointer, to prevent it from being released at the end of the function. This is because the interface output routine will release the mbuf after it has been output. (Remember that `rip_output` passes the mbuf chain to `ip_output`, who appends it to the interface's output queue.)

The final part of `rip_usrreq` is shown in Figure 32.16. The `PRU_SENSE` request, generated by the `fstat` system call, returns nothing. The `PRU_SOCKADDR` and `PRU_PEERADDR` requests are from the `getsockname` and `getpeername` system calls, respectively. The remaining requests are not supported.

398-404 The functions `in_setsockaddr` and `in_setpeeraddr` fetch the information from the PCB, storing the result in the `nam` argument.

```
304     case PRU_SENSE:
305         /*
306          * fstat: don't bother with a blocksize.
307          */
308         return (0);
309
310         /*
311          * Not supported.
312          */
313     case PRU_RCVOOB:
314     case PRU_RCVD:
315     case PRU_LISTEN:
316     case PRU_ACCEPT:
317     case PRU_SENDOOB:
318         error = EOPNOTSUPP;
319         break;
320
321     case PRU_SOCKADDR:
322         in_setsockaddr(inp, nam);
323         break;
324
325     case PRU_PEERADDR:
326         in_setpeeraddr(inp, nam);
327         break;
328
329     default:
330         panic("rip_usrreq");
331 }
332
333 if (m != NULL)
334     m_freem(m);
335
336 return (error);
337 }
```

Figure 32.16 rip_usrreq function: remaining requests.

32.8 rip_ctloutput Function

The `setsockopt` and `getsockopt` system calls invoke the `rip_ctloutput` function. Only one IP socket option is handled here, along with eight socket options related to multicast routing.

Figure 32.17 shows the first part of the `rip_ctloutput` function.

144-172

The size of the mbuf that contains either the new value of the option or will hold the current value of the option must be at least as large as an integer. For the `setsockopt` system call, the flag is set if the integer value in the mbuf is nonzero, or cleared otherwise. For the `getsockopt` system call, the value returned in the mbuf is either 0 or the nonzero value of the flag. The function returns, to avoid the processing at the end of the switch statement for other IP options.

```

144 int
145 rip_ctloutput(op, so, level, optname, m)
146 int op;
147 struct socket *so;
148 int level, optname;
149 struct mbuf **m;
150 {
151     struct inpcb *inp = sotoinpcb(so);
152     int error;
153
154     if (level != IPPROTO_IP)
155         return (EINVAL);
156
157     switch (optname) {
158     case IP_HDRINCL:
159         if (op == PRCO_SETOPT || op == PRCO_GETOPT) {
160             if (m == 0 || *m == 0 || (*m)->m_len < sizeof(int))
161                 return (EINVAL);
162             if (op == PRCO_SETOPT) {
163                 if (*mtod(*m, int *))
164                     inp->inp_flags |= INP_HDRINCL;
165                 else
166                     inp->inp_flags &= ~INP_HDRINCL;
167                 (void) m_free(*m);
168             } else {
169                 (*m)->m_len = sizeof(int);
170                 *mtod(*m, int *) = inp->inp_flags & INP_HDRINCL;
171             }
172             return (0);
173         }
174         break;

```

Figure 32.17 rip_usrreq function: process IP_HDRINCL socket option.

```

173     case DVMRP_INIT:
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:

```

```

/* shown in Figure 14.9 */

```

```

188     }
189     return (ip_ctloutput(op, so, level, optname, m));
190 }

```

Figure 32.18 rip_usrreq function: process multicast routing socket option.

Figure 32.18 shows the last portion of the `rip_etloutput` function. It handles eight multicast routing socket options.

173-188 These eight socket options are valid only for the `setsockopt` system call. They are processed by the `ip_mrouter_cmd` function as discussed with Figure 14.9.

189 Any other IP socket options, such as `IP_OPTIONS` to set the IP options, are processed by `ip_etloutput`.

32.9 Summary

Raw sockets provide three capabilities for an IP host.

1. They are used to send and receive ICMP and IGMP messages.
2. They allow a process to build its own IP headers.
3. They allow additional IP-based protocols to be supported in a user process.

We saw that raw IP output is simple—it just fills in a few fields in the IP header—but it allows a process to supply its own IP header. This allows diagnostic programs to create any type of IP datagram.

Raw IP input provides three types of filtering for incoming IP datagrams. The process chooses to receive datagrams based on (1) the protocol field, (2) the source IP address (set by `connect`), and (3) the destination IP address (set by `bind`). The process chooses which combination of these three filters (if any) to apply.

Exercises

- 32.1 Assume the `IP_HDRINCL` socket option is not set. What value will `rip_output` place into the IP header protocol field (`ip_p`) when the third argument to `socket` is 0? What value will `rip_output` place into this field when the third argument to `socket` is `IPPROTO_RAW` (255)?
- 32.2 A process creates a raw socket with a protocol value of `IPPROTO_RAW` (255). What type of IP datagrams will the process receive on this socket?
- 32.3 A process creates a raw socket with a protocol value of 0. What type of IP datagrams will the process receive on this socket?
- 32.4 Modify `rip_input` to send an ICMP destination unreachable with code 2 (protocol unreachable) when appropriate. Be careful not to generate the error for received ICMP and IGMP packets for which `rip_input` is called.
- 32.5 If a process wants to write its own IP datagrams with its own IP header, what are the differences in using a raw IP socket with the `IP_HDRINCL` option, and using BPF (Chapter 31)?
- 32.6 When would a process read from a raw IP socket, and when would it read from BPF?

Epilogue

"We have come a long way. Nine chapters stuffed with code is a lot to negotiate. If you didn't assimilate all of it the first time through, don't worry—you weren't really expected to. Even the best of code takes time to absorb, and you seldom grasp all the implications until you try to use and modify the program. Much of what you learn about programming comes only from working with the code: reading, revising and rereading."

From the Epilogue of *Software Tools* [Kernighan and Plauger 1976].

"In fact, this RFC will argue that modularity is one of the chief villains in attempting to obtain good performance, so that the designer is faced with a delicate and inevitable tradeoff between good structure and good performance."

From RFC 817 [Clark 1982].

This text has provided a long and detailed examination of a significant piece of a real operating system. Versions of the code presented in the text are shipped as part of the Unix kernel with most flavors of Unix today, along with many non-Unix systems.

The code that we've examined is not perfect and it is not the only way to write a TCP/IP protocol stack. It has been modified, enhanced, tweaked, and maligned over the past 15 years by many people. Large portions of the code that we've presented weren't even written at the U. C. Berkeley Computer Systems Research Group: the multicasting code was written by Steve Deering, the long fat pipe support was added by Thomas Skibo, portions of the TCP code were written by Van Jacobson, and so on. The code contains `gotos` (221 to be exact), many large functions (e.g., `tcp_input` and `tcp_output`), and numerous examples of questionable coding style. (We tried to note these items when discussing the code.) Nevertheless, the code is unquestionably "industrial strength" and continues to be the base upon which new features are added and the standard upon which other implementations are measured.

The Berkeley networking code was designed on VAXes when a VAX-11/780 with 4 megabytes of memory was a big system. For that reason some of the design features (e.g., mbufs) emphasized memory savings over higher performance. This would change if the code were rewritten from scratch today.

There has been a strong push over the last few years toward higher performance of networking software, as the underlying networks become faster (e.g., FDDI and ATM) and as high-bandwidth applications become more prevalent (e.g., voice and video). Whenever designing networking software within the kernel of an operating system, speed normally gives way to clarity [Clark 1982]. This will continue in any real-world implementation.

The research implementation of the Internet protocols described in [Partridge 1993] and [Jacobson 1993] is a move toward much higher performance. [Jacobson 1993] reports the code is 10 to 100 times faster than the implementation described in this book. Mbufs, software interrupts, and much of the protocol layering evident in BSD systems are gone. If widely released, this implementation could become the standard that others are measured against in the future.

In July 1994 the successor to IP version 4, IP version 6 (IPv6), was announced. It uses 128-bit (16-byte) addresses. Many changes will take place with the IP and ICMP protocols, but the transport layers, UDP and TCP, will remain virtually the same. (There is talk of a TCPng, the next generation of TCP, but the authors think just upgrading IP will provide enough of a challenge for the hundreds of vendors and millions of users across the world to put off any changes to TCP.) It will take a year or two for vendor-supported implementations to appear, and many years after that for end users to migrate their hosts and routers to IPv6. Research implementations of IPv6 based on the code in this text should appear in early 1995.

To continue your understanding of the Berkeley networking code, the best course of action at this point is to obtain the source code, and modify it. The source code is easily obtainable (Appendix B) and numerous exercises throughout the text suggest modifications.

Appendix A

Solutions to Selected Exercises

Chapter 1

- 1.2 SLIP drivers execute at `spltty` (Figure 1.13), which must be a priority lower than or equal to `splimp` and must be a priority higher than `splnet`. Therefore the SLIP drivers are blocked from interrupting.

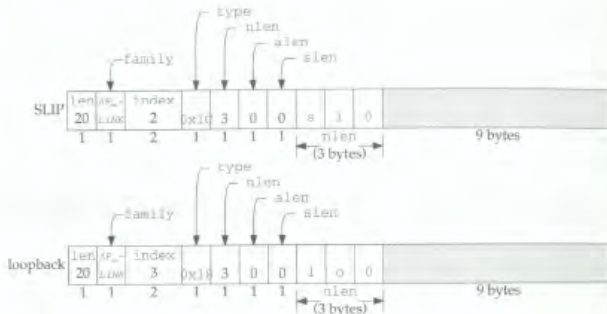
Chapter 2

- 2.1 The `M_EXT` flag is a property of the mbuf itself, not a property of the packet described by the mbuf.
- 2.2 The caller asks for more than 100 (`MHLEN`) contiguous bytes.
- 2.3 This is infeasible since clusters can be pointed to by multiple mbufs (Section 2.9). Also, there is no room in a cluster for a back pointer (Exercise 2.4).
- 2.4 In the macros `MCLALLOC` and `MCLFREE` in `<sys/mbuf.h>` we see that the reference count is an array named `mchrefcnt`. This array is allocated when the kernel is initialized in the file `machdep.c`.

Chapter 3

- 3.3 A large interactive queue would defeat the purpose of the queue by delaying new interactive traffic behind the existing interactive data.
- 3.4 Since the `sl_softc` structures are all declared as global variables, they are initialized to 0 when the kernel starts.

3.5



Chapter 4

- 4.1 `lread` must examine the packet to decide if it needs to be discarded after it is passed to BPF. Since a BPF tap can enable promiscuous mode on the interface, packets may be addressed to some other system on the Ethernet and must be discarded after BPF has processed them.

When the interface is not tapped, the tests must be done in `ether_input`.

- 4.2 If the tests were reversed, the broadcast flag would never be set.

If the second `if` wasn't preceded by an `else`, every broadcast packet would also have the multicast flag set.

Chapter 5

- 5.1 The loopback interface does not need an input function because all its packets are received directly from `looutput`, which performs the "input" functions.
- 5.2 The stack allocation is faster than dynamic memory allocation. Performance is important for BPF processing, since the code is executed for each incoming packet.
- 5.5 The first character that overflows the buffer is discarded, `SC_ERROR` is set, and `sinput` resets the cluster pointers to begin collecting characters at the start of the buffer. Because `SC_ERROR` is set, `sinput` discards the frame when it receives the SLIP END character.
- 5.6 IP discards the packet when the checksum is found to be invalid or when it notices that the IP header length does not match the physical packet size.

- 5.7 Since `ifp` points to the first member of a `ie_softc` structure,

```
sc = (struct ie_softc *)ifp;
```

initializes `sc` correctly.

- 5.8 This is very hard to do. Some routers may send ICMP source quench messages when they begin discarding packets but Net/3 discards these messages for UDP sockets (Figure 23.30). An application would have to begin using the same techniques used by TCP: estimation of the available bandwidth and delay on roundtrip times for acknowledged datagrams.

Chapter 6

- 6.1 Before IP subnetting (RFC 950 [Mogul and Postel 1985]), the network and host portions of IP addresses always appeared on byte boundaries. The definition of an `in_addr` structure was

```
struct in_addr {
    union {
        struct { u_char s_b1, s_b2, s_b3, s_b4; } s_un_b;
        struct { u_short s_w1, s_w2; } s_un_w;
        u_long s_addr;
    } s_un;
#define s_addr s_un.s_addr /* should be used for all code */
#define s_host s_un.s_b.s_b2 /* OBSOLETE: host on imp */
#define s_net s_un.s_un_b.s_b1 /* OBSOLETE: network */
#define s_imp s_un.s_un_w.s_w2 /* OBSOLETE: imp */
#define s_impno s_un.s_un_b.s_b4 /* OBSOLETE: imp # */
#define s_lh s_un.s_un_b.s_b3 /* OBSOLETE: logical host */
};
```

The Internet address could be accessed as 8-bit bytes, 16-bit words, or a single 32-bit address. The macros `s_host`, `s_net`, `s_imp`, and so on have names that correspond to the physical structure of early TCP/IP networks.

The use of subnetting and supernetting makes the byte and word divisions obsolete.

- 6.2 A pointer to the structure labeled `s1_softc[0]` is returned.
- 6.3 The interface output functions, such as `ether_output`, have a pointer only to the `ifnet` structure for the interface, and not to an `ifaddr` structure. Using the IP address in the `arpcom` structure (which is the last IP address assigned to the interface) avoids having to select an address from the `ifaddr` address list.
- 6.4 Only a superuser process can create a raw IP socket. By using a UDP socket, any process can examine the interface configurations but the kernel can still require superuser privileges to modify the interface addresses.
- 6.5 Three functions loop through a `netmask` 1 byte at a time. These are `ifa_ifwithnet`, `ifaof_ifpforaddr`, and `rt_maskedcopy`. A shorter mask improves the performance of these functions.

- 6.6 The Telnet connection is established with the remote system. Net/2 systems shouldn't forward these packets, and other systems should never accept loopback packets that arrive on any interface other than the loopback interface.

Chapter 7

- 7.1 The following call returns a pointer to `inetaw(6)`

```
pf_in@proto(PF_INET, 0, SOCK_RAW)
```

Chapter 8

- 8.1 Probably not. The system could not respond to any broadcasts since it would have no source address to use in the reply.
- 8.4 Since the packet has been damaged, there is no way of knowing if the addresses in the header are correct or not.
- 8.5 If an application selects a source address that differs from the address of the selected outgoing interface, redirects from the selected next-hop router fail. The next-hop router sees a source address different from that of the subnetwork on which it was transmitted and does not send a redirect message. This is a consequence of implementing the weak end system model and is noted in RFC 1122.
- 8.6 The new host thinks the broadcast packet is the address of some other host in the unsubnetted network and tries to send it back out on the network. The network interface begins broadcasting ARP requests for the broadcast address, which are never answered.
- 8.7 The decrement of the TTL is done after the comparison for less than or equal to 1 to avoid the potential error of decrementing a received TTL of 0 to become 255.
- 8.8 If two routers each consider the other the best next-hop for a packet, a routing loop exists. Until the loop is removed, the original packet bounces between the two routers and each one sends an ICMP redirect back to the source host if that host is on the same network as the routers. Loops may exist when the routing tables are temporarily inconsistent during a routing update.
- The TTL of the original packet eventually reaches 0 and the packet is discarded. This is one of the primary reasons why the TTL field exists.
- 8.9 The four Ethernet broadcast addresses would not be checked because they do not belong to the receiving interface. The limited-broadcast addresses would be checked. This implies that a system on a SLIP link can communicate with the system on the other end without knowing the other system's address by utilizing the limited-broadcast address.
- 8.10 ICMP error messages are generated only for the initial fragment of a datagram, which always has an offset of 0. The host and network forms for 0 are the same, so no conversion is necessary.

Chapter 9

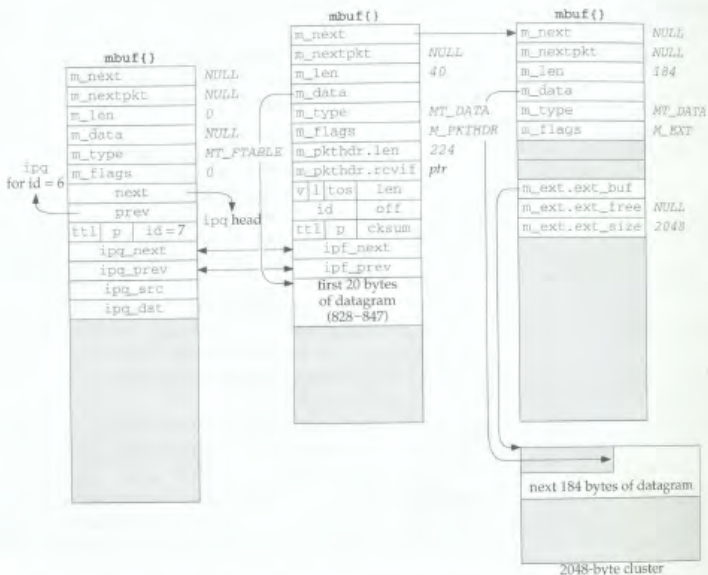
- 9.1 RFC 1122 says that the behavior is implementation dependent when conflicting options appear in a packet. Net/3 processes the first source route option correctly, but since this updates `ip_dst` in the packet header, the second source route processing will be incorrect.
- 9.2 The host within the network can be used as a relay to access other hosts within the network. To communicate with an otherwise-blocked host, the source host need only construct packets with a loose route to the relay host and then to the final destination host. The router does not drop the packets because the destination address is the relay host, which will process the route and forward the packet to the final destination host. The destination host reverses the route and uses the relay host to return packets.
- 9.3 The same principle from the previous exercise applies. We pick a relay router that can communicate with the source and destination hosts and construct source routes to pass through the relay and to the destination. The relay router must be on the same network as the destination host so that a default route is not required for communication.

This technique can be extended to allow two hosts to communicate even if they do not have routes to each other, as long as they can find willing relay hosts.

- 9.4 If the source route is the only IP option, the NOP option causes all the IP addresses to be on a 4-byte boundary in the IP header. This can optimize memory references to these addresses on many architectures. This alignment technique also works when multiple options are present if each option is padded with NOPs to a 4-byte boundary.
- 9.5 A nonstandard time value cannot be confused with a standard value since the largest standard time value is $86,399,999$ ($24 \times 60 \times 60 \times 1000 - 1$) and this value can be represented in 28 bits, which avoids any conflict with the high-order bit since time values are 32 bits long.
- 9.6 The source route option code may change `ip_dst` in the packet during processing. The destination is saved so that the timestamp processing code uses the original destination.

Chapter 10

- 10.2 After reassembly, only the options from the initial fragment are available to the transport protocols.
- 10.3 The fragment is read into a cluster since the data length ($204 + 20$) is greater than 208 (Figure 2.16).
`m_duplup` in Figure 10.11 moves the first 40 bytes into a separate mbuf as in Figure 2.18.



- 10.4 The average number of received fragments per datagram is

$$\frac{72,786 - 349}{16,557} = 4.4$$

The average number fragments created for an outgoing datagram is

$$\frac{796,084}{260,484} = 3.1$$

- 10.5 In Figure 10.11 the packet is initially processed as a fragment. The reserved bit is discarded when `ip_off` is left shifted. The resulting packet is processed as a fragment or as a complete datagram, depending on the values of the MF and off-set bits.

Chapter 11

- 11.1 The outgoing reply uses the source address of the interface on which the request was received. Hosts are not required to recognize 0.0.0.0 as a valid broadcast

address, so the request may be ignored. The recommended broadcast address is 255.255.255.255.

- 11.2 Assume that a host sends link-level broadcast packets with the IP source address of another host and the packet contains errors such as an improperly formed option. Every host receives and detects the error because of the link-level broadcast and because options are processed before a final destination check. Many hosts that detect the error try to send an ICMP message back to the IP source of the packet even though the original packet was sent as a link-level broadcast. The unfortunate host will begin receiving many bogus ICMP error messages. This is one reason why ICMP errors must not be sent in response to link-level broadcasts.
- 11.3 In the first case, such a redirect message can fool the host into sending packets to an arbitrary host on an alternate subnetwork. This host may be masquerading as a router but recording the traffic it receives instead. RFC 1009 requires that routers only generate redirect messages for other routers on the same subnet. Even if the host ignores these messages to redirect packets to a new subnetwork, a host on the same subnetwork can fool the host. The second case guards against this by requiring that the host only accept the redirect advice from the original router that it had (erroneously) selected to receive the traffic. Presumably this incorrect router was a default router specified by an administrator.
- 11.4 By passing the message to `rip_input`, a process-level daemon could respond and old systems that relied on this behavior could continue to be supported.
- 11.5 ICMP errors are sent only for the initial fragment of an IP datagram. Since the offset value of an initial fragment is always 0, the byte ordering of the field is unimportant.
- 11.6 If the ICMP request was received on an interface that was not yet configured with an IP address, `ia` would be null and no reply could be generated.
- 11.7 Net/3 reflects the data along with the timestamp reply.
- 11.10 The high-order bit is reserved and must be 0. If it is sent, `icmp_error` will discard the packet.
- 11.11 The return value is discarded because `icmp_send` does not return an error, but more significantly, errors generated during ICMP processing are discarded to avoid generating an endless series of error messages.

Chapter 12

- 12.1 On an Ethernet, the IP broadcast address 255.255.255.255 translates to the Ethernet broadcast address `ff:ff:ff:ff:ff:ff` and is received by every Ethernet interface on the network. Systems that aren't running IP software must actively receive and discard each of these broadcast packets.

A packet sent to the IP all-hosts multicast group 224.0.0.1 translates to the Ethernet multicast address `01:00:5e:00:00:01` and is received only by systems

that have explicitly instructed their interfaces to receive IP multicast datagrams. Systems that aren't running IP or that aren't level-2 compliant never receive these datagrams, as they are discarded by the Ethernet interface hardware itself.

- 12.2 One alternative would be to specify interfaces by their text name as with the `ifreq` structure and the `ioctl` commands for accessing interface information. `ip_getoptions` and `ip_getmoptions` would have to call `ifunit` instead of `INADDR_TO_IFP` to locate the pointer to the interface's `ifnet` structure.
- 12.3 The high-order 4 bits of a multicast group are always 1110, so only 5 significant bits are discarded by the mapping function.
- 12.4 The entire `ip_moptions` structure must fit within an mbuf, which limits the size of the structure to 108 bytes (remember the 20-byte mbuf header). `IP_MAX_MEMBERSHIPS` can be larger but must be less than or equal to 25. $(4+1+1+2+(4 \times 25) = 108)$
- 12.5 The datagram is duplicated and two copies appear on the IP input queue. A multicast application must be prepared to discard duplicate datagrams.
- 12.6



- 12.8 The process could create a second socket and request another `IP_MAX_MEMBERSHIPS` through the second socket.
- 12.9 Define a new mbuf flag `M_LOCAL` for the `m_flags` member of the mbuf header. The flag can be set on loopback packets by `ip_output` instead of computing the checksum. `ipintr` can skip the checksum verification if the flag is on. SunOS 5.X has an option to do this (`ip_local_cksum`, page 531, Volume 1).
- 12.10 There are $2^{23} - 1$ (8,388,607) unique Ethernet IP multicast addresses. Remember that IP group 224.0.0.0 is reserved.
- 12.11 This assumption is correct since `in_addmulti` rejects all add requests if the interface does not have an `ioctl` function, and this implies that `in_delmulti` is never called if `if_ioctl` is null.
- 12.12 The mbuf is never released. It appears that `ip_getmoptions` contains a memory leak. `ip_getmoptions` is called from `ip_ctloutput`, which allows a call such as:

```
ip_getoptions(IP_ADD_MEMBERSHIP, 0, mp)
```

which exercises the bug in `ip_getmoptions`.

Chapter 13

- 13.1 Responding to an ICMP query from the loopback interface is unnecessary since

the local host is the only system on the loopback network and it already knows its membership status.

- 13.2 $\text{max_linkhdr} + \text{sizeof}(\text{struct ip}) + \text{IGMP_MINLEN} = 16 + 20 + 8 = 44 < 100$
- 13.3 The primary reason for the random delay in reporting memberships is to minimize (ideally to 1) the number of reports that appear on a multicast network. A point-to-point network consists only of two interfaces, so the delay is not necessary to minimize the response to the query. One interface (presumably a multicast router) generates the query, and the other interface responds.

There is another reason not to flood the interface's output queue with all the membership reports. The output queue may have a packet or byte limit that could be exceeded by many IGMP membership reports. For example, in the SLIP driver, if the output queue is full or the device is too busy, the entire queue of pending packets is discarded (Figure 5.16).

Chapter 14

- 14.1 Five. One each for networks A through E.
- 14.2 `grp1st_member` is called only by `ip_mforward`, but `ip_mforward` can be called by `ipintr` during protocol processing, or by `ip_output`, which can be called indirectly from the socket layer. The cache is a shared data structure that must be protected while it is being updated. The membership list itself is protected by `splx` calls in `add_lgrp` and `del_lgrp`, where it is modified.
- 14.3 The `SIOCDELMULTI` command affects only the Ethernet multicast list for the interface. The IP multicast group list remains unchanged, so the interface remains a member of the group. The interface continues accepting multicast datagrams for any groups that are still on the IP group membership list for the interface. Specifically, when `ether_delmulti` returns `ENETRESET` to `leioctl`, the function `iereset` is called to reconfigure the interface (Figure 12.31).
- 14.4 Only one virtual interface is considered to be the parent interface for a multicast spanning tree. If the packet is accepted on the tunnel, then the physical interface cannot be the parent and `ip_mforward` discards the packet.

Chapter 15

- 15.1 The socket could be shared across a fork or passed to a process through a Unix domain socket ([Stevens 1990]).
- 15.2 The `sa_len` member of the structure is larger than the size of the buffer after `accept` returns. This is usually not a problem with the fixed-length Internet address, but it can be when using variable-length addresses supported by the OSI protocols, for example.

- 15.4 The call to `soqremque` is only made when `so_glen` is not equal to 0. If `soqremque` returns a null pointer there must be an error in the socket queueing code so the kernel panics.
- 15.5 The copy is made so that `bzero` can clear the structure while it is locked and so that `dom_dispose` and `sbrlease` can be called after `splx`. This minimizes the amount of time the CPU is kept at `splim` and therefore the amount of time that network interrupts are blocked.
- 15.6 The `sbspace` macro will return 0. As a result, the `sbappendaddr` and `sbappendcontrol` functions (used by UDP) will refuse to queue additional packets. TCP uses `sbappend`, which assumes that the caller has checked for space first. TCP calls `sbappend` even when `sbspace` returns 0. The data placed in the receive queue is not available to a process because the `SS_CANTRCVMORE` flag prevents the read system calls from returning any data.

Chapter 16

- 16.1 When the value is assigned to `uio_resid` in the `uio` structure it becomes a large negative number. `soSEND` rejects the message with `EINVAL`.

Net/2 did not check for a negative value. This problem is described by the comment at the start of `soSEND` (Figure 16.23).

- 16.2 No. The only time the cluster is ever filled with less than `MCLBYTES` is at the end of a message when less than `MCLBYTES` remain. `resid` is 0 at this time and the loop is terminated by the `break` on line 394 before reaching the test for `space > 0`.
- 16.5 The process blocks until the buffer is unlocked. In this case the lock exists only while another process is examining the buffer or passing data to the protocol layer, and not when a process must wait for space in the buffer, which may take an indefinite amount of time.
- 16.6 If the send buffer contained many mbufs, each of which contained only a few bytes of data, `sb_cc` may be well below the limit specified by `sb_hiwat` while a large amount of memory would be allocated for the mbufs. If the kernel didn't limit the number of mbufs attached to each buffer, a process could easily create a memory shortage.
- 16.7 `recvit` is called from `recvfrom` and `recvmsg`. Only `recvmsg` handles control information. The entire `msghdr` structure, including the length of the control message, is copied back to the process by `recvmsg`. For address information, `recvmsg` sets the `namelenp` argument to null because it expects the length in `msg_namelen`. When `recvfrom` calls `recvit`, the `namelenp` is nonnull because it expects the length in `*namelenp`.
- 16.8 `MSG_EOR` is cleared by `soRECEIVE` so that it is not inadvertently returned by `soRECEIVE` before an `M_EOR` mbuf is processed.

- 16.9 There would be a race condition while `select` examined the descriptors. If a selectable event occurred after `select` examined the descriptors but before `select` called `tsleep`, it would not be detected and the process would sleep until another selectable event occurred.

Chapter 17

- 17.1 This simplifies the code that copies data between the kernel and the process. `copyin` and `copyout` can be used for a single mbuf, but `uiomove` is needed to handle multiple mbufs.
- 17.2 The code works correctly because the first member of a `linger` structure is the expected integer flag.

Chapter 18

- 18.1 Write eight rows, one for each possible combination of the bits from the search key, the routing table key, and the routing table mask.

row	1 search key	2 table key	3 table mask	1 & 3	2 == 4?	1 ^ 2	6 & 3
1	0	0	0	0	yes	0	0=yes
2	0	0	1	0	yes	0	0=yes
3	0	1	0	0	no	1	0=yes
4	0	1	1	0	no	1	1=no
5	1	0	0	0	yes	1	0=yes
6	1	0	1	1	no	1	1=no
7	1	1	0	0	no	0	0=yes
8	1	1	1	1	yes	0	0=yes

The column “2 == 4?” should equal the final column “6 & 3.” On first glance they are not the same, but we can ignore rows 3 and 7 because in these two rows the routing table bit is 1 while the same bit in the routing table mask is 1. When the routing table is built the key is logically ANDed with the mask, guaranteeing that for every bit of 0 in the mask, the corresponding bit in the key is also 0.

Another way to look at the exclusive OR and logical AND in Figure 18.40 is that the exclusive OR becomes 1 only if the the search key bit differs from the bit in the routing table key. The logical AND then ignores any differences that correspond to a bit that's 0 in the mask. If the result is still nonzero, the search key does not match the routing table key.

- 18.2 The size of an `rtable` structure is 120 bytes, which includes the two `radix_node` structures. Each entry also requires two `sockaddr_in` structures (Figure 18.28), for 152 bytes per routing table entry. The total is about 3 megabytes.
- 18.3 Since `rn_b` is a short integer, assuming 16 bits for a short imposes a limit of 32767 bits per key (4095 bytes).

Chapter 19

- 19.1 The `RTF_DYNAMIC` flag is set in Figure 19.15 when the route is created by a redirect, and the `RTF_MODIFIED` flag is set when the gateway field of an existing route is modified by a redirect. If a route is created by a redirect and then later modified by another redirect, both flags will be set.
- 19.2 A host route is created for each host accessed through the default route. TCP can then maintain and update routing metrics for each individual host (Figure 27.3).
- 19.3 Each `rt_msghdr` structure requires 76 bytes. Two `sockaddr_in` structures are present for a host route (destination and gateway) giving a message size of 108 bytes. The message size for each ARP entry is 112 bytes: one `sockaddr_in` and one `sockaddr_dl`. The total size is then $(15 \times 112 + 20 \times 108)$ or 3840 bytes. A network route (instead of a host route) requires an additional 8 bytes for the network mask (116 bytes for the message instead of 108), so if the 20 routes are all network routes, the total size is 4000 bytes.

Chapter 20

- 20.1 The return value is returned in the `rtm_errno` member of the message (Figure 20.14) and also as the return value from `write` (Figure 20.22). The latter is more reliable since the former may run into mbuf starvation, causing the reply message to be discarded (Figure 20.17).
- 20.2 For a `SOCK_RAW` socket, the `pf_findproto` function (Figure 7.20) returns the entry with a protocol of 0 (the wildcard) if an exact match isn't found.

Chapter 21

- 21.1 It is assumed that the `ifnet` structure is at the beginning of the `arpcom` structure, which it is (Figure 3.20).
- 21.2 Sending the ICMP echo request does not require ARP, since the destination address is the broadcast address. But the ICMP echo replies are normally unicast, so each sender uses ARP to determine the destination Ethernet address. When the local host receives each ARP request, `in_arpinput` replies and creates an entry for the other host.
- 21.3 When a new ARP entry is created, the `rt_gateway` value, a `sockaddr_dl` structure in this case, is copied from the entry being cloned by `rtrequest` in Figure 19.8. In Figure 21.1 we see that the `sdl_alen` member of this entry is 0.
- 21.4 With Net/3, if the caller of `arpresolve` supplies a pointer to a routing table entry, `arplookup` is not called, and the corresponding Ethernet address is available through the `rt_gateway` pointer (assuming it hasn't expired). This avoids any type of lookup in the common case. In Chapter 22 we'll see that TCP and UDP store a pointer to their routing table entry in their protocol control block.

avoiding a search of the routing table in the case of TCP (where the destination IP address never changes for a connection) and in the case of UDP when the destination doesn't change.

- 21.5 The timeout of an incomplete ARP entry occurs between 0 and 5 minutes after the entry is created. `arpresolve` sets `rt_expire` to the current time when the ARP request is sent. The next time `arptimer` runs, if that entry is not resolved, it is deleted (assuming its reference count is 0).
- 21.6 `ether_output` returns `EHOSTUNREACH` instead of `EHOSTDOWN`, causing an ICMP host unreachable error to be sent to the sending host by `ip_forward`.
- 21.7 The value for 140.252.13.32 is set in Figure 21.28 to the current time when the entry is created. It never changes.

The values for 140.252.13.33 and 140.252.13.34 are copied from the entry for 140.252.13.32 when these two entries are cloned by `rtrequest`. They are then set to the time at which an ARP request is sent by `arpresolve`, and finally set by `in_arpinput` to the time at which an ARP reply is received, plus 20 minutes.

The value for 140.252.13.35 is also copied from the entry for 140.252.13.32 when the entry is cloned, but then set to 0 by the code at the end of Figure 21.29.

- 21.8 Change the call to `arplookup` at the beginning of Figure 21.19 to always specify a second argument of 1 (the create flag).
- 21.9 The first datagram was sent *after* the halfway mark to the next second. Therefore both the first and second datagrams caused ARP requests to be sent, about 500 ms apart, since the kernel's `time.tv_sec` variable had different values when these two datagrams were sent.
- 21.10 Each packet to send is an mbuf chain. The `m_nextpkt` pointer in the first mbuf in each chain could be used to form a list of mbufs awaiting transmission.

Chapter 22

- 22.1 An infinite loop occurs, waiting for a port to become available. This assumes the process is allowed to open enough descriptors to tie up all ephemeral ports.
- 22.2 Few, if any, servers support this option. [Cheswick and Bellonin 1994] mention how this would be nice for implementing firewall systems.
- 22.4 The `udb` structure is initialized to 0 so `udb.inp_lport` starts at 0. The first time through `in_pcbbind` it is incremented to 1, which is less than 1024, so it is set to 1024.
- 22.5 Normally the caller sets the address family (`sa_family`) to `AF_INET`, but we saw in Figure 22.20 that the test for this is commented out. The caller can set the length member (`sa_len`), but we saw in Figure 15.20 that the function `sockargs` always sets this to the third argument to `bind`, which for a `sockaddr_in` structure is specified as 16, normally using C's `sizeof` operator.

The local IP address (*sin_addr*) can be specified as a wildcard address or as a local IP address. The local port number (*sin_port*), can be either 0 (telling the kernel to choose an ephemeral port) or nonzero if the process wants a particular port. Normally a TCP or UDP server specifies a wildcard IP address and a nonzero port, and a UDP client often specifies a wildcard IP address and a port number of 0.

- 22.6 A process is allowed to bind a local broadcast address, because the call to *ifa_ifwithaddr* in Figure 22.22 succeeds. That address is used as the source address for IP datagrams sent on the socket. As noted in Section C.2, this behavior is not allowed by RFC 1122.

An attempt to bind 255.255.255.255, however, fails, since that address is not acceptable to *ifa_ifwithaddr*.

Chapter 23

- 23.1 *send* places the user data into a single mbuf if the size is less than or equal to 100 bytes; into two mbufs if the size is less than or equal to 207 bytes; or into one or more mbufs, each with a cluster, otherwise. Furthermore, *send* calls *MH_ALIGN* if the size is less than 100 bytes, which, it is hoped, will allow room at the beginning of the mbuf for the protocol headers. Since *udp_output* calls *M_PREPEND*, the following five scenarios are possible: (1) If the size of the user data is less than or equal to 72 bytes, a single mbuf contains the IP header, UDP header, and data. (2) If the size is between 73 and 100 bytes, one mbuf is allocated by *send* for the data and another is allocated by *M_PREPEND* for the IP and UDP headers. (3) If the size is between 101 and 207 bytes, two mbufs are allocated by *send* for the data and another by *M_PREPEND* for the IP and UDP headers. (4) If the size is between 208 and *MCLBYTES*, one mbuf with a cluster is allocated by *send* for the data and another by *M_PREPEND* for the IP and UDP headers. (5) Beyond this size, *send* allocates as many mbufs with clusters as necessary to hold the data (up to 64 for a maximum data size of 65507 bytes with 1024-byte clusters), and one mbuf is allocated by *M_PREPEND* for the IP and UDP headers.

- 23.2 IP options are passed to *ip_output*, which calls *ip_insertoptions* to insert the options into the outgoing IP datagram. This function in turn allocates a new mbuf to hold the IP header including options if the first mbuf in the chain points to a cluster (which never happens with UDP output) or if there is not enough room at the beginning of the first mbuf in the chain for the options. In scenario 1 from the previous solution, the size of the options determines whether another mbuf is allocated by *ip_insertoptions*: if the size of the user data is less than $100 - 28 - \text{optlen}$, (where *optlen* is the number of bytes of IP options), there is room in the mbuf for the IP header with options, the UDP header, and the data.

In scenarios 2, 3, 4, and 5, the first mbuf in the chain is always allocated by *M_PREPEND* just for the IP and UDP headers. *M_PREPEND* calls *m_prepend*,

which calls `MH_ALIGN`, moving the 28 bytes of headers to the end of the `mbuf`, hence there is always room for the maximum of 40 bytes of IP options in this first `mbuf` in the chain.

- 23.3 No. The function `in_pcbconnect` is called, either when the application calls `connect` or when the first datagram is sent on an unconnected UDP socket. Since the local address is a wildcard and the local port is 0, `in_pcbconnect` sets the local port to an ephemeral port (by calling `in_pcbbind`) and sets the local address based on the route to the destination.
- 23.4 The processor priority level is left at `spinet`; it is not restored to the saved value. This is a bug.
- 23.5 No. `in_pcbconnect` will not allow a connection to port 0. Even if the process doesn't call `connect` directly, an implicit `connect` is performed, so `in_pcbconnect` is called regardless.
- 23.6 The application must call `ioctl` with the `SIOCGIFCONF` command to return information on all configured IP interfaces. The destination address in the received UDP datagram must then be compared against all the IP addresses and broadcast addresses in the list returned by `ioctl`. (As an alternative to `ioctl`, the `sysctl` system call described in Section 19.14 can also be used to obtain the information on all the configured interfaces.)
- 23.7 `recvit` releases the `mbuf` with the control information.
- 23.8 To disconnect a connected UDP socket, call `connect` with an invalid address, such as 0.0.0.0, and a port of 0. Since the socket is already connected, `sconnect` calls `sodisconnect`, which calls `udp_usrreq` with a `PRU_DISCONNECT` request. This sets the foreign address to 0.0.0.0 and the foreign port to 0, allowing a subsequent call to `sendto` that specifies a destination address to succeed. Specifying the invalid address causes the `PRU_CONNECT` request from `sodisconnect` to fail. We don't want the `connect` to succeed, we just want the `PRU_DISCONNECT` request executed and this back door through `connect` is the only way to execute this request, since the sockets API doesn't provide a `disconnect` function.
- The manual page for `connect(2)` usually contains the following note that hints at this: "Datagram sockets may dissolve the association by connecting to an invalid address, such as a null address." What this note fails to mention is that the call to `connect` for the invalid address is expected to return an error. The term *null address* is also vague: it means the IP address 0.0.0.0, not a null pointer for the second argument to `bind`.
- 23.9 Since an unconnected UDP socket is temporarily connected to the foreign IP address by `in_pcbconnect`, the scenario is the same as if the process calls `connect`: the datagram is sent out the primary interface with a destination IP address corresponding to the broadcast address of that interface.
- 23.10 The server must set the `IP_RECVSTADDR` socket option and use `recvmsg` to obtain the destination IP address from the client's request. For this address to be

the source IP address of the reply requires that this IP address be bound to the socket. Since you cannot bind a socket more than once, the server must create a brand new socket for each reply.

- 23.11 Notice in `ip_output` (Figure 8.22) that IP does not modify the DF bit supplied by the caller. A new socket option could be defined to cause `udp_output` to set the DF bit before passing datagrams to IP.
- 23.12 No. It is used only in the `udp_input` function and should be local to that function.

Chapter 24

- 24.1 The total number of ESTABLISHED connections is 126,820. Dividing this into the total number of bytes transmitted and received yields an average of about 30,000 bytes in each direction.
- 24.2 In `tcp_output`, the mbuf obtained for the IP and TCP headers also contains room for the link-layer headers (`max_linkhdr`). The IP and TCP header prototype is copied into the mbuf using `bcopy`, which won't work if the 40-byte header were split between two mbufs. Although the 40-byte headers must fit into one mbuf, the link-layer header need not. But a performance penalty would occur later (`ether_output`) because a separate mbuf would be required for the link-layer header.
- 24.3 On the author's system `bsd1`, the count was 16, 15 of which were standard system daemons (Telnet, Rlogin, FTP, etc.). On `vangogh.cs.berkeley.edu`, a medium-sized multiuser system with around 20 users, the count was 60. On a large multiuser system (`world.std.com`) with around 150 users, the count was 417 TCP end points and 809 UDP end points.

Chapter 25

- 25.1 In Figure 24.5 there were 531,285 delayed ACKs over 2,592,000 seconds (30 days). This is an average of about one delayed ACK every 5 seconds, or one delayed ACK every 25 times `tcp_fasttimo` is called. This means 96% of the time (24 times out of every 25) every TCP control block is checked for the delayed-ACK flag, when not one is set. On the large multiuser system in the solution to Exercise 24.3, this involves looking at over 400 control blocks, 5 times a second.

One alternative implementation would be to set a global flag when a delayed ACK is needed and only go through the list of control blocks when the flag is set. Alternatively, another list could be maintained that contains only the control blocks that require a delayed ACK. See, for example, the variable `igmp_timers_are_running` in Figure 13.14.

- 25.2 This allows the variable `tcp_keepintvl` to be patched in the running kernel, which then changes the value of `tcp_maxidle` the next time `tcp_slowtimo` is called.

- 25.3 `t_idle` actually counts the time since a segment was last received or transmitted. This is because TCP output must be acknowledged by the other end and the receipt of the ACK clears `t_idle`, as does the receipt of a data segment (Figure 28.8).
- 25.4 Here is one way to rewrite the code:
- ```

case TCPT_2MSL:
 if (tp->t_state == TCPS_TIME_WAIT)
 tp = tcp_close(tp);
 else {
 if (tp->t_idle <= tcp_maxidle)
 tp->t_timer[TCPT_2MSL] = tcp_keepintvl;
 else
 tp = tcp_close(tp);
 }
 break;

```
- 25.5 When the duplicate ACK is received, `t_idle` is 150, but it is reset to 0. When the `FIN_WAIT_2` timer expires, `t_idle` will be 1048 ( $1198 - 150$ ), so the timer is set to 150 ticks. When the timer expires the next time, `t_idle` will be 1198, so the timer is set to 150 ticks. When the timer expires the next time, `t_idle` will be  $1198 + 150$ , so the connection is closed. The duplicate ACK extends the time until the connection is closed.
- 25.6 The first keepalive probe will be sent 1 hour in the future. When the process sets the option, nothing happens other than setting the `SO_KEEPALIVE` option in the `socket` structure. When the timer expires 1 hour in the future, since the option is enabled, the code in Figure 25.15 sends the first probe.
- 25.7 The value of `tcp_rttdfilt` initializes the RTT estimators for every TCP connection. A site can change the default of 3, if desired, by patching the global variable. If the value were a `#define` constant, it could be changed only by recompiling the kernel.

## Chapter 26

- 26.1 The counter `t_idle` is always running for a connection, whereas TCP does not measure the amount of time since the last segment was sent on a connection.
- 26.2 In Figure 25.26 `snd_nxt` is set to `snd_una`, giving a value of 0 for `len`.
- 26.3 If you're running a Net/3 system and encounter a peer that can't handle either of these two newer options (i.e., that peer refuses to establish the connection, even though a host is required to ignore options it doesn't understand), this global can be patched in the kernel to disable one or both of these options.
- 26.4 The timestamp option would have updated the RTT estimators each time an ACK was received for new data: 16 times, twice the number of times without the option. The value calculated when the ACK of 6145 was received at time 217.944, however, would have been bogus—either the data segment with bytes



- 5633 through 6144 that was sent at time 3,740, or the received ACK of 6145, was delayed somewhere for about 200 seconds.
- 26.5 There is no guarantee that the 2-byte MSS value is correctly aligned for such a memory reference.
- 26.6 (This solution is from Dave Borman.) The maximum amount of TCP data in a segment is 65495 bytes, which is 65535 minus the minimum IP and TCP headers (40). Hence there are 39 values of the urgent offset that make no sense: 65496 through and including 65535. Whenever the sender has a 32-bit urgent offset that exceeds 65495, 65535 is sent as the urgent offset instead, and the URG flag is set. This puts the receiver into urgent mode and tells the receiver that the urgent offset points to data that has not been sent yet. The special value of 65535 continues to be sent as the urgent offset (with the URG flag set) until the urgent offset is less than or equal to 65495, at which point the real urgent offset is sent.
- 26.7 We've mentioned that data segments are transmitted reliably (i.e., the retransmission timer is set) but ACKs are not. RST segments are not transmitted reliably either. RST segments are generated when a bogus segment arrives (either a segment that is wrong for a connection, or a segment for a nonexistent connection). If the RST segment is discarded by `ip_output`, when the other end retransmits the segment that caused the RST to be generated, another RST will be generated.
- 26.8 The application does eight writes of 1024 bytes. The first four times `so_send` is called, `tcp_output` is called, and a segment is sent. Since these four segments each contain the final bytes of data in the send buffer, the PSH flag is set for each segment (Figure 26.25). The send buffer is also full, so the next write by the process puts the process to sleep in `so_send`. When the ACK is returned with an advertised window of 0, the 4096 bytes of data in the send buffer have been acknowledged and are discarded, and the process wakes up and continues filling the send buffer with the next four writes. But nothing can be sent until a nonzero window is advertised by the receiver. When this happens, the next four segments are sent, but only the final segment contains the PSH flag, since the first three segments do not empty the send buffer.
- 26.9 The `tp` argument to `tcp_respond` can be a null pointer if the segment being sent does not correspond to a connection. The code should check the value of `tp` and use the default only if the pointer is null.
- 26.10 `tcp_output` always allocates an mbuf just to contain the IP and TCP headers, by calling `MGETHDR` in Figures 26.25 and 26.26. This code allocates room at the front of the new mbuf only for the link-layer header (`max_linkhdr`). If IP options are in use and the size of the options exceeds `max_linkhdr`, another mbuf is allocated by `ip_insertoptions`. If the size of the IP options is less than or equal to `max_linkhdr`, then even though `ip_insertoptions` will use the space at the beginning of the mbuf, this will cause `ether_output` to allocate another mbuf for the link-layer header (assuming Ethernet output).
- To try to avoid the extra mbuf, Figures 26.25 and 26.26 could call `MH_ALIGN` if the segment will contain IP options.

- 26.11 About 80 lines of C code, assuming RFC 1323 timestamps are in use and the segment is timed.

The macro `MGETHDR` invokes the macro `MALLOC`, which might call the function `malloc`. The function `m_copy` is also called, but a full-sized segment will be in a cluster, so the mbuf is not copied, a reference is made to the cluster. The call to `MGET` by `m_copy` might call `malloc`. The function `bcopy` copies the header template and `in_cksum` calculates the TCP checksum.

- 26.12 Nothing changes with `writew` because of the logic in `sosend`. Since the total size of the data (150) is less than `MINCLSIZE` (208), one mbuf is allocated for the first 100 bytes, and since the protocol is not atomic, the `PRU_SEND` request is issued. Another mbuf is allocated for the next 50 bytes, and another `PRU_SEND` is issued. TCP still generates two segments. (`writew` only generates a single "record," that is, a single `PRU_SEND` request, for `PR_ATOMIC` protocols such as UDP.)

With two buffers of length 200 and 300 the total size now exceeds `MINCLSIZE`. An mbuf cluster is allocated and only one `PRU_SEND` is issued. One 500-byte segment is generated by TCP.

## Chapter 27

- 27.1 The first six rows of the table are asynchronous errors that are generated by the receipt of a segment or the expiration of a timer. By storing the nonzero error code in `so_error`, the process receives the error on the next read or write. The call from `tcp_disconnect`, however, occurs when the process calls `close`, or when the descriptor is closed automatically on process termination. In either case of the descriptor being closed, the process won't issue a read or write call to fetch the error. Also, since the process had to set the socket option explicitly to force the RST, returning an error provides no useful information to the process.
- 27.2 Assuming a 32-bit `u_long`, the maximum value is just under 4298 seconds (1.2 hours).
- 27.3 The statistics in the routing table are updated by `tcp_close` and it is called only when the connection enters the CLOSED state. Since the sending of data to the other end is terminated by the FTP client (it does the active close), the local end point enters the TIME\_WAIT state. The routing table statistics won't be updated until twice the MSL has elapsed.

## Chapter 28

- 28.1 0, 1, 2, and 3.
- 28.2 34.9 Mbits/sec. For higher speeds, larger buffers are required on both ends.
- 28.3 In the general case, `tcp_dooptions` doesn't know whether the two timestamp values are aligned on 32-bit boundaries or not. The special code in Figure 28.4,

- however, knows that the values are on 32-bit boundaries, and avoids calling `copy`.
- 28.4 The "options prediction" code in Figure 28.4 handles only the recommended format, so systems that send other than the recommended format cause the slower processing of `tcp_dooptions` to occur for every received segment.
- 28.5 If `tcp_doamp1st` were called every time a socket were created, instead of every time a connection is established, each listening server on a system would have one allocated, which it would never use.
- 28.6 The timestamp clock frequency should be between 1 bit/ms and 1 bit/sec. (Net/3 uses 2 bits/sec.) With the highest frequency of 1 bit/ms, a 32-bit timestamp wraps its sign bit in  $2^{31} / (24 \times 60 \times 60 \times 1000)$  days, which is 24.8 days.
- 28.7 With a frequency of 1 bit per 500 ms, a 32-bit timestamp wraps its sign bit in  $2^{31} / (24 \times 60 \times 60 \times 2)$  days, which is 12,427 days, or about 34 years, longer than the uptime of current computer systems.
- 28.8 The cleanup function of an RST should take precedence over timestamps, and it is recommended that RSTs not carry timestamps (which is enforced by `tcp_input` in Figure 26.24).
- 28.9 Since the client is in the ESTABLISHED state, processing ends up in Figure 28.24. `to_drop` is 1 because `rcv_nxt` was incremented over the SYN when it was first received. The SYN flag is cleared (since it is a duplicate), `tl_seq` is incremented, and `to_drop` is decremented to 0. The `if` statement at the top of Figure 28.25 is executed since `to_drop` and `tl_len` are both 0. The next `if` statement is skipped, and processing continues with the call to `m_adj`. But `tcp_output` is not called in the continuation of `tcp_input` in the next chapter, therefore the client does not respond to the duplicate SYN/ACK. The server will time out and resend the SYN/ACK (recall the timer set in Figure 28.17 when a passive socket receives a SYN), which will also be ignored. This is another bug in the code in Figure 28.25 and this one is also fixed with the code shown in Figure 28.30.
- 28.10 The client's SYN arrives at the server and is delivered to the socket in the TIME\_WAIT state. The code in Figure 28.24 turns off the SYN flag and the code in Figure 28.25 jumps to `dropAfterack`, dropping the segment but generating an ACK with an acknowledgment field of `rcv_nxt` (Figure 26.27). This is called a *resynchronization ACK* because its purpose is to tell the other end what sequence number it expects. When this ACK is received at the client (which is in the SYN\_SENT state), its acknowledgment field is not the expected value (Figure 28.18), causing an RST to be sent. The sequence number of the RST is the acknowledgment field from the resynchronization ACK, and the ACK flag of the RST segment is off (Figure 29.28). When the server receives the RST, its TIME\_WAIT state is prematurely terminated and the socket is closed on the server's host (Figure 28.36). The client times out after 6 seconds and retransmits its SYN. Assuming a listening server process is running on the server host, the new connection is established. Because of this form of TIME\_WAIT



assassination, a new connection is established not only when a SYN arrives with a higher sequence number (as checked for in Figure 28.28), but also when a SYN with a lower sequence number arrives.

## Chapter 29

- 29.1 Assume a 2-second RTT. The server has a passive open pending and the client issues its active open at time 0. The server receives the SYN at time 1 and responds with its own SYN and an ACK of the client's SYN. The client receives this segment at time 2, and the code in Figure 28.20 completes the active open with the call to `soisconnected` (waking up the client process) and an ACK will be sent back to the server. The server receives the ACK at time 3, and the code in Figure 29.2 completes the server's passive open, returning control to the server process. In general, the client process receives control about one-half RTT before the server.
- 29.2 Assume the sequence number of the SYN is 1000 and the 50 bytes of data are numbered 1001–1050. When the SYN is processed by `tcp_input`, first the case starting in Figure 28.15 is executed, which sets `rcv_nxt` to 1001, and then a jump is made to `step6`. Figure 29.22 calls `tcp_reass` and the data is placed onto the socket's reassembly queue. But the data cannot be appended to the socket's receive buffer yet (Figure 27.23) so `rcv_nxt` is left at 1001. When `tcp_output` is called to generate the immediate ACK, `seq_nxt` (1001) is sent as the acknowledgment field. In summary, the SYN is acknowledged, but not the 50 bytes of data. Since the client will retransmit the 30 bytes of data, there is no advantage in sending data with a SYN generated by an active open.
- 29.3 The server's socket is in the SYN\_RCVD state when the client's ACK/FIN arrives, so `tcp_input` ends up processing the ACK in Figure 29.2. The connection moves to the ESTABLISHED state and `tcp_reass` appends the already-queued data to the socket's receive buffer. `rcv_nxt` is incremented to 1051. `tcp_input` continues and the FIN is handled in Figure 29.24 where the `TF_ACKNOW` flag is set and `rcv_nxt` becomes 1052. `socketrcvmore` sets the socket's state so that after the server reads the 50 bytes of data, the server will receive an end-of-file. The server's socket also moves to the CLOSE\_WAIT state, `tcp_output` will be called to ACK the client's FIN (since `rcv_nxt` equals 1052). Assuming the server process closes its socket when it reads the end-of-file, the server will then send a FIN for the client to ACK.

In this example six segments requiring three round trips are required to pass the 50 bytes from the client to server. To reduce the number of segments requires the TCP extensions for transactions [Braden 1994].

- 29.4 The client's socket is in the SYN\_SENT state when the server's response is received. Figure 28.20 processes the segment and moves the connection to the ESTABLISHED state. A jump is made to `step6` and the data is processed in Figure 29.22. `TCP_REASS` appends the data to the socket's receive buffer and

`rcv_nxt` is incremented to acknowledge the data. The FIN is then processed in Figure 29.24, incrementing `rcv_nxt` again and moving the connection to the CLOSE\_WAIT state. When `tcp_output` is called, the acknowledgment field ACKs the SYN, the 50 bytes of data, and the FIN. The client process then reads the 50 bytes of data, followed by the end-of-file, and then probably closes its socket. This moves the connection to the LAST\_ACK state and causes a FIN to be sent by the client, which the server should acknowledge.

- 29.5 The bug is in the entry `tcp_outflags[TCPS_CLOSING]` shown in Figure 24.16. It specifies the `TH_FIN` flag, whereas the state transition diagram (Figure 24.15) doesn't specify that the FIN should be retransmitted. To fix this, remove `TH_FIN` from the `tcp_outflags` entry for this state. The bug is relatively harmless—it just causes two extra segments to be exchanged—and a simultaneous close or a close following a self-connect is rare.
- 29.6 No. An OK return from a write system call only means the data has been copied into the socket buffer. Net/3 does not notify the process when that data is acknowledged by the other end. An application-level acknowledgment is required to obtain this information.
- 29.7 RFC 1323 timestamps defeat header compression because whenever the timestamps change, the TCP options change, and the segment is sent uncompressed. The window scale option has no effect because the value in the TCP header is still a 16-bit value.
- 29.8 IP assigns the ID field from a global variable that is incremented each time *any* IP datagram is sent. This increases the probability that two consecutive TCP segments sent on the same connection will have ID values that differ by more than 1. A difference other than 1 causes the  $\Delta ipid$  field in Figure 29.34 to be transmitted, increasing the size of the compressed header. A better scheme would be for TCP to maintain its own counter for assigning IDs.

## Chapter 30

- 30.2 Yes, the RST is still sent. Part of process termination is the closing of all open descriptors. The same function (`socklose`) is eventually called, regardless of whether the process explicitly closes the socket descriptor or implicitly closes it (by terminating first).
- 30.3 No. The only use of this constant is when a listening socket sets the `SO_LINGER` socket option with a linger time of 0. Normally this causes an RST to be sent when the connection is closed (Figure 30.12), but Figure 30.2 changes this value of 0 to 120 (clock ticks) for a listening socket that receives a connection request.
- 30.4 Two if this is the first use of the default route; otherwise one. When the socket is created the Internet PCB is set to 0 by `in_pcballoc`. This sets the `route` structure in the PCB to 0. When the first segment is sent (the SYN), `tcp_output` calls `ip_output`. Since the `ro_rte` pointer is null, `ro_det` is filled in with the destination address of the IP datagram and `rtaalloc` is called. The pointer to the



default route is saved in the `ro_rt` member of the `route` structure within the PCB for this connection. When `ether_output` is called by `ip_output`, it checks whether the `rt_gwroute` member of the routing table entry is null, and, if so, `rtalloc` is called. Assuming the route doesn't change, each time `tcp_output` is called for this connection, the cached `ro_rt` pointer is used, avoiding any additional routing table lookups.

## Chapter 31

- 31.1 Because `catchpacket` will always run to completion before any sleeping processes are awakened by the `bpf_wakeup` call.
- 31.2 A process that opens a BPF device may call `fork` resulting in multiple processes with access to the same BPF device.
- 31.3 Only supported devices are on the BPF interface list (`bpf_iflist`), so `bpf_setif` returns `ENXIO` when the interface is not found.

## Chapter 32

- 32.1 0 in the first example, and 255 in the second. Both of these values are reserved in RFC 1700 [Reynolds and Postel 1994] and should not appear in datagrams. This means, for example, that a socket created with a protocol of `IPPROTO_RAW` should always have the `IP_HDRINCL` socket option set, and datagrams written to the socket should have a valid protocol value.
- 32.2 Since the IP protocol value of 255 is reserved, datagrams should never appear on the wire with this protocol value. Since this is a nonzero protocol value, the first of the three tests in `rip_input` will ignore every received datagram that does not have this protocol value. Therefore the process should not receive any datagrams on the socket.
- 32.3 Even though this protocol value is reserved and datagrams should never appear on the wire with this value, the first of the three tests in `rip_input` allows datagrams with any protocol value to be received by sockets of this type. The only input filtering that occurs for this type of raw socket is based on the source and destination IP addresses, if the process calls either `connect` or `bind`, or both.
- 32.4 Since the array `ip_protosx` (Figure 7.22) contains information about which protocol the kernel supports, the ICMP error should be generated only when there are no raw listeners for the protocol and the pointer `inetsw[ip_protosx[ip->ip_p]]_pr_input` equals `rip_input`.
- 32.5 In both cases the process must build its own IP header, in addition to whatever follows the IP header (UDP datagram, TCP segment, or whatever). With a raw IP socket, output is normally done using `sendto` specifying the destination address as an Internet socket address structure containing an IP address. `ip_output` is called and normal IP routing is done based on the destination IP address.

BPF requires the process to supply a complete data-link header, such as an Ethernet header. Output is normally done by calling `write`, since a destination address cannot be specified. The packet is passed directly to the interface output function, bypassing `ip_output` (Figure 31.20). The process selects the outgoing interface using the `BIOCSETIF ioctl` (Figure 31.16). Since IP routing is not performed, the destination of the packet is limited to another system on an attached network (unless the process duplicates the IP routing function and sends the packet to a router on an attached network, for the router to forward based on the destination IP address).

- 32.6 A raw IP socket receives only IP datagrams destined for an IP protocol that the kernel does not process itself. A process cannot receive TCP segments or UDP datagrams on a raw socket, for example.

BPF can receive *all* frames received on a specified interface, regardless of whether they are IP datagrams or not. The `BIOCPRMISC ioctl` can put the interface into a promiscuous mode, to receive datagrams that are not even destined for this host.

# Appendix B

## Source Code Availability

### URLs: Uniform Resource Locators

This text uses URLs to specify the location and method of access of resources on the Internet. For example, the common “anonymous FTP” technique is designated as

```
ftp://ftp.cdrom.com/pub/bsd-sources/4.4BSD-Lite.tar.gz
```

This specifies anonymous FTP to the host `ftp.cdrom.com`. The filename is `4.4BSD-Lite.tar.gz` in the directory `pub/bsd-sources`. The suffix `.tar` implies the standard Unix `tar(1)` format, and the additional `.gz` suffix implies that the file has been compressed with the GNU `gzip(1)` program.

### 4.4BSD-Lite

There are numerous ways to obtain the 4.4BSD-Lite release. The entire 4.4BSD-Lite release is available from Walnut Creek CD-ROM as

```
ftp://ftp.cdrom.com/pub/bsd-sources/4.4BSD-Lite.tar.gz
```

You can also obtain this release on CD-ROM. Contact 1 800 786 9907 or +1 510 674 0783.

O'Reilly & Associates publishes the entire set of 4.4BSD manuals along with the 4.4BSD-Lite release on CD-ROM. Contact 1 800 889 8969 or +1 707 829 0515.

### Operating Systems that Run the 4.4BSD-Lite Networking Software

The 4.4BSD-Lite release is *not* a complete operating system. To experiment with the networking software described in this text you need an operating system that is built from

1093

the 4.4BSD-Lite release or an environment that supports the 4.4BSD-Lite networking code.

The operating system used by the authors is commercially available from Berkeley Software Design, Inc. Contact 1 800 ITS BSD8, +1 719 260 8114, or [info@bsd.i.com](mailto:info@bsd.i.com) for additional information.

There are also freely available operating systems built on 4.4BSD-Lite. These are known by the names NetBSD, 386BSD, and FreeBSD. Additional information is available from Walnut Creek CD-ROM ([ftp.cdrom.com](http://ftp.cdrom.com)) or on the various `comp.os.386bsd` Usenet newsgroups.

## RFCs

All RFCs are available at no charge through electronic mail or by using anonymous FTP across the Internet. Sending electronic mail as shown here:

```
To: rfc-info@ISI.EDU
Subject: getting rfc#
help: ways_to_get_rfcs
```

returns a detailed listing of various ways to obtain the RFCs using either email or anonymous FTP.

Remember that the starting place is to obtain the current index and look up the RFC that you want in the index. This entry tells you if that RFC has been made obsolete or updated by a newer RFC.

## GNU Software

The GNU Indent program was used to format all the source code presented in the text, and the GNU Gzip program is often used on the Internet to compress files. These programs are available as

```
http://prep.ai.mit.edu/pub/gnu/indent-1.9.1.tar.gz
http://prep.ai.mit.edu/pub/gnu/gzip-1.2.2.tar
```

The numbers in the filenames will change as newer versions are released. There are also versions of the Gzip program for other operating systems, such as MS-DOS.

There are many sites around the world that also provide the GNU archives, and the FTP greeting on `prep.ai.mit.edu` displays their names.

## PPP Software

There are several freely available implementations of PPP. Part 5 of the `comp.protocols.ppp` FAQ is a good place to start:

```
http://gw.uni-bonn.de/ppp/part5.html
```

**mrouterd Software**

Current releases of the mrouterd software as well as other multicast applications can be found at the Xerox Palo Alto Research Center:

<ftp://parcftp.xerox.com/pub/net-research/>

**ISODE Software**

An SNMP agent implementation compatible with Net/3 is part of the ISODE software package. For more information, start with the ISODE Consortium's World Wide Web page at

<http://www.isode.com/>





# Appendix C

## RFC 1122 Compliance

This appendix summarizes the compliance of the Net/3 implementation with RFC 1122 [Braden 1989a]. This RFC summarizes these requirements in four categories

- link layer
- internet layer
- UDP
- TCP

We have chosen to present these requirements in the same breakdown and order as the chapters of this text.

### C.1 Link-Layer Requirements

This section summarizes the link-layer requirements from Section 2.5 of RFC 1122 and the compliance of the Net/3 code that we've examined to those requirements.

- *May* support trailer encapsulation.  
Partially: Net/3 does not send IP datagrams with trailer encapsulation but some Net/3 device drivers may be able to receive such datagrams. We have omitted all the trailer encapsulation code in this text. Interested readers are referred to RFC 893 and Section 11.8 of [Leffler et al. 1989] for additional details.
- *Must* not send trailers by default without negotiation.  
Not applicable: Net/2 would negotiate the use of trailers but Net/3 ignores requests to send trailers and does not request trailers itself.

1097

- *Must* be able to send and receive RFC 894 Ethernet encapsulation.  
Yes: Net/3 supports RFC 894 Ethernet encapsulation.
- *Should* be able to receive RFC 1042 (IEEE 802) encapsulation.  
No: Net/3 processes packets received with 802.3 encapsulation but only for use with OSI protocols. IP packets that arrive with 802.3 encapsulation are discarded by `ether_input` (Figure 4.13).
- *May* send RFC 1042 encapsulation, in which case there must be a software configuration switch to select the encapsulation method and RFC 894 must be the default.  
No: Net/3 does not send IP packets in RFC 1042 encapsulation.
- *Must* report link-layer broadcasts to the IP layer.  
Yes: The link layer reports link-layer broadcasts by setting the `M_BCAST` flag (or the `M_MCAST` flag for multicasts) in the mbuf packet header.
- *Must* pass the IP TOS value to the link layer.  
Yes: The TOS value is not passed explicitly, but is part of the IP header available to the link layer.

## C.2 IP Requirements

This section summarizes the IP requirements from Section 3.5 of RFC 1122 and the compliance of the Net/3 code that we've examined to those requirements.

- *Must* implement IP and ICMP.  
Yes: `inetsw[0]` implements the IP protocol and `inetsw[4]` implements ICMP.
- *Must* handle remote multihoming in application layer.  
Yes: The kernel is unaware of communication to remote multihomed hosts and neither hinders nor supports such communication by an application.
- *May* support local multihoming.  
Yes: Net/3 supports multiple IP interfaces with the `ifnet` list and multiple addresses per IP interface with the `ifaaddr` list for each `ifnet` structure.
- *Must* meet router specifications if forwarding datagrams.  
Partially: See Chapter 18 for a discussion of the router requirements.
- *Must* provide configuration switch for embedded router functionality. The switch must default to host operation.  
Yes: The `ipforwarding` variable defaults to false and controls the IP packet forwarding mechanism in Net/3.

- *Must not* enable routing based on number of interfaces.  
Yes: The `if_attach` function does not modify `ipforwarding` according to the number of interfaces configured at system initialization time.
- *Should* log discarded datagrams, including the contents of the datagram, and record the event in a statistics counter.  
Partially: Net/3 does not provide a mechanism for logging the contents of discarded datagrams but maintains a variety of statistics counters.
- *Must* silently discard datagrams that arrive with an IP version other than 4.  
Yes: `ipintr` implements this requirement.
- *Must* verify IP checksum and silently discard an invalid datagram.  
Yes: `ipintr` calls `ip_cksum` and implements this requirement.
- *Must* support subnet addressing (RFC 950).  
Yes: Every IP address has an associated subnet mask in the `in_ifaddr` structure.
- *Must* transmit packets with host's own IP address as the source address.  
Partially: When the transport layer sends an IP datagram with all-0 bits as the source address, IP inserts the IP address of the outgoing interface in its place. A process can bind one of the local IP broadcast addresses to the local socket, and IP will transmit it as an invalid source address.
- *Must* silently discard datagrams not destined for the host.  
Yes: If the system is not configured as a router, `ipintr` discards datagrams that arrive with a bad destination address (i.e., an unrecognized unicast, broadcast, or multicast address).
- *Must* silently discard datagrams with bad source address (nonunicast address).  
No: `ipintr` does not examine the source address of incoming datagrams before delivering the datagram to the transport protocols.
- *Must* support reassembly.  
Yes: `ip_reass` implements reassembly.
- *May* retain same ID field in identical datagrams.  
No: `ip_output` assigns a new ID to every outgoing datagram and does not allow the ID to be specified by the transport protocols. See Chapter 32.
- *Must* allow the transport layer to set TOS.  
Yes: `ip_output` accepts any TOS value set in the IP header by the transport protocols. The transport layer must default TOS to all 0s. The TOS value for a particular datagram or connection may be set by the application through the `IP_TOS` socket option.

- *Must* pass received TOS up to transport layer.  
Yes: Net/3 preserves the TOS field during input processing. The entire IP header is made available to the transport layer when IP calls the `pr_input` function for the receiving protocol. Unfortunately, the UDP and TCP transport layers ignore it.
- *Should not* use RFC 795 [Postel 1981d] link-layer mappings for TOS.  
Yes: Net/3 does not use these mappings.
- *Must not* send packet with TTL of 0.  
Partially: The IP layer (`ip_output`) in Net/3 does not check this requirement and relies on the transport layers not to construct an IP header with a TTL of 0. UDP, TCP, ICMP, and IGMP all select a nonzero TTL default value. The default value can be overridden by the `IP_TTL` option.
- *Must not* discard received packets with a TTL less than 2.  
Yes: If the system is the final destination of the packet, `ipintr` accepts it regardless of the TTL value. The TTL is examined only when the packet is being forwarded.
- *Must* allow transport layer to set TTL.  
Yes: The transport layer must set TTL before calling `ip_output`.
- *Must* enable configuration of a fixed TTL.  
Yes: The default TTL is specified by the global integer `ip_defttl`, which defaults to 64 (`IPDEFTTL`). Both UDP and TCP use this value unless the `IP_TTL` socket option has specified a different value for a particular socket. `ip_defttl` can be modified through the `IPCTL_DEFTTL` name for `sysctl`.

## Multihoming

- *Should* select, as the source address for a reply, the specific address received as the destination address of the request.  
Yes: Responses generated by the kernel (ICMP reply messages) include the correct source address (Section C.5). Responses generated by the transport protocols are described in their respective chapters.
- *Must* allow application to choose local IP address.  
Yes: An application can bind a socket to a specific local IP address (Section 15.8).
- *May* silently discard datagrams addressed to an interface other than the one on which it is received.  
No: Net/3 implements the weak end system model and `ipintr` accepts such packets.
- *May* require packets to exit the system through the interface with an IP address that corresponds to the source address of the packet. This requirement pertains only to packets that are not source routed.



No: Net/3 allows packets to exit the system through any interface—another weak end system characteristic.

### Broadcast

- *Must* not select an IP broadcast address as a source address.  
Partially: If an application explicitly selects a source address, the IP layer does not override the selection. Otherwise, IP selects as a source address the specific IP address associated with the outgoing interface.
- *Should* accept an all-0s or all-1s broadcast address.  
Yes: `ipintr` accepts packets sent to either address.
- *May* support a configurable option to send all 0s or all 1s as the broadcast address on an interface. If provided, the configurable broadcast address *must* default to all 1s.  
No: A process must explicitly send to either the all-0s (`INADDR_ANY`) or all-1s broadcast address (`INADDR_BROADCAST`). There is no configurable default.
- *Must* recognize all broadcast address formats.  
Yes: `ipintr` recognizes the limited (all-1s and all-0s) and the network-directed and subnet-directed broadcast addresses.
- *Must* use an IP broadcast or IP multicast destination address in a link-layer broadcast.  
Yes: `ip_output` enables the link-layer multicast or broadcast flags only when the destination is an IP multicast or broadcast address.
- *Should* silently discard link-layer broadcasts when the packet does not specify an IP broadcast address as its destination.  
No: There is no explicit test for the `M_BCAST` or `M_MCAST` flags on incoming packets in Net/3, but `ip_forward` will discard these packets before forwarding them.
- *Should* use limited broadcast address for connected networks.  
Partially: The decision to use the limited broadcast address (versus a subnet-directed or network-directed broadcast) is left to the application level by Net/3.

### IP Interface

- *Must* allow transport layer to use all IP mechanisms (e.g., IP options, TTL, TOS).  
Yes: All the IP mechanisms are available to the transport layer in Net/3.
- *Must* pass interface identification up to transport layer.  
Yes: The `m_pkthdr.revif` member of each mbuf containing an incoming packet points to the `ifnet` structure of the interface that received the packet.

- *Must* pass all IP options to transport layer.  
Yes: The entire IP header, including options, is present in the packet passed to the `pr_input` function of the receiving transport protocol by `ipintr`.
- *Must* allow transport layer to send ICMP port unreachable and any of the ICMP query messages.  
Yes: The transport layer may send any ICMP error messages by calling `icmp_error` or may format and send any type of IP datagram by calling the `ip_output` function.
- *Must* pass the following ICMP messages to the transport layer: destination unreachable, source quench, echo reply, timestamp reply, and time exceeded.  
Yes: These messages are distributed by ICMP to other transport protocols or to any waiting processes using the raw IP socket mechanism.
- *Must* include contents of ICMP message (IP header plus the data bytes present) in ICMP message passed to the transport layer.  
Yes: `icmp_input` passes the portion of the original IP packet contained within the ICMP message to the transport layers.
- *Should* be able to leap tall buildings at a single bound.  
No: The next version of IP may meet this requirement.

### C.3 IP Options Requirements

This section summarizes the IP option processing requirements from Section 3.5 of RFC 1122 and the compliance of the Net/3 code that we've examined to those requirements.

- *Must* allow transport layer to send IP options.  
Yes: The second argument to `ip_output` is a list of IP options to include in the outgoing IP datagram.
- *Must* pass all IP options received to higher layer.  
Yes: The IP header and options are passed to the `pr_input` function of the receiving transport protocol.
- *Must* silently ignore unknown options.  
Yes: The default case in `ip_dooptions` skips over unknown options.
- *May* support the security option.  
No: Net/3 does not support the IP security option.

- *Should not* send the stream identifier option and *must ignore* it in received datagrams.  
Yes: Net/3 does not support the stream identifier option and ignores it on incoming datagrams.
- *May* support the record route option.  
Yes: Net/3 supports the record route option.
- *May* support the timestamp option.  
Partially: Net/3 supports the timestamp option but does not implement it exactly as specified. The originating host does not insert a timestamp when required but the destination host records a timestamp before passing the datagram to the transport layer. The timestamp value follows the rules regarding standard values as specified in Section 3.2.2.8 of RFC 1122 for the ICMP timestamp message.
- *Must* support originating a source route and *must* be able to act as the final destination of a source route.  
Yes: A source route may be included in the options passed to `ip_output`, and `ip_doptions` correctly terminates a source route and saves it for use in constructing return routes.
- *Must* pass a datagram with completed source route up to the transport layer.  
Yes: The source route option is passed up with any other options that may have appeared in the datagram.
- *Must* build correct (nonredundant) return route.  
No: Net/3 blindly reverses the source route and does not check or correct for a route that was built incorrectly with a redundant hop for the original source host.
- *Must* not send multiple source route options in one header.  
No: The IP layer in Net/3 does not prohibit a transport protocol from constructing and sending multiple source route options in a single datagram.

### Source Route Forwarding

- *May* support packet forwarding with the source route option.  
Yes: Net/3 supports the source route options. `ip_doptions` does all the work.
- *Must* obey corresponding router rules while processing source routes.  
Yes: Net/3 follows the router rules whether or not the packet contains a source route.
- *Must* update TTL according to gateway rules.  
Yes: `ip_forward` implements this requirement.

- *Must* generate ICMP error codes 4 and 5 (fragmentation required and source route failed).  
Yes: `ip_output` is able to generate a fragmentation required message, and `ip_dooptions` is able to generate the source route failed message.
- *Must* allow the IP source address of a source routed packet to not be an IP address of the forwarding host.  
Yes: `ip_output` transmits such packets.  

RFC 1122 lists this as a *may* requirement because the addresses *may* be different, which *must* be allowed.
- *Must* update timestamp and record route options.  
Yes: `ip_dooptions` processes these options for source routed packets.
- *Must* support a configurable switch for *nonlocal source routing*. The switch *must* default to off.  
No: Net/3 always allows nonlocal source routing and does not provide a switch to disable this function. Nonlocal source routing is routing packets between two different interfaces instead of receiving and sending the packet on the same interface.
- *Must* satisfy gateway access rules for nonlocal source routing.  
Yes: Net/3 follows the forwarding rules for nonlocal source routing.
- *Should* send an ICMP destination unreachable error (source route failed) if a source routed packet cannot be forwarded (except for ICMP error messages).  
Yes: `ip_dooptions` sends the ICMP destination unreachable error. `icmp_error` discards it if the original datagram was an ICMP error message.

## C.4 IP Fragmentation and Reassembly Requirements

This section summarizes the IP fragmentation and reassembly requirements from Section 3.5 of RFC 1122 and the compliance of the Net/3 code that we've examined to those requirements.

- *Must* be able to reassemble incoming datagrams of at least 576 bytes.  
Yes: `ip_reass` supports reassembly of datagrams of indefinite size.
- *Should* support a configurable or indefinite maximum size for incoming datagrams.  
Yes: Net/3 supports an indefinite maximum size for incoming datagrams.
- *Must* provide a mechanism for the transport layer to learn the maximum datagram size to receive.  
Not applicable: Net/3 has an indefinite limit based on available memory.

- *Must* send ICMP time exceeded error on reassembly timeout.  
No: Net/3 does not send an ICMP time exceeded error. See Figure 10.30 and Exercise 10.1.
- *Should* support a fixed reassembly timeout value. The remaining TTL value in a received IP fragment *should not* be used as a reassembly timeout value.  
Yes: Net/3 uses a compile-time value of 30 seconds (`IPFRAGTTL` is 60 slow-timeout intervals, which equals 30 seconds).
- *Must* provide the `MMS_S` (maximum message size to send) to higher layers.  
Partially: TCP derives the `MMS_S` from the maximum MTU found in the route entry for the destination or from the MTU of the outgoing interface. A UDP application does not have access to this information.
- *May* support local fragmentation of outgoing packets.  
Yes: `ip_output` fragments an outgoing packet if it is too large for the selected interface.
- *Must not* allow transport layer to send a message larger than `MMS_S` if local fragmentation is not supported.  
Not applicable: This is a transport-level requirement that does not apply to Net/3 since local fragmentation is supported.
- *Should not* send messages larger than 576 bytes to a remote destination in the absence of other information regarding the minimum path MTU to the destination.  
Partially: Net/3 TCP defaults to a segment size of 552 (512 data bytes + 40 header bytes). Net/3 UDP applications cannot determine if a destination is local or remote and so they often restrict their messages to 540 bytes (512 + 20 + 8). There is no kernel mechanism that prohibits sending larger messages.
- *May* support an all-subnets-MTU configuration flag.  
Yes: The global integer `subnetsarelocal` defaults to `true`. TCP uses this flag to select a larger segment size (the size of the outgoing interface's MTU) instead of the default segment size for destinations on a subnet of the local network.

## C.5 ICMP Requirements

This section summarizes the ICMP requirements from Section 3.5 of RFC 1122 and the compliance of the Net/3 code that we've examined to those requirements.

- *Must* silently discard ICMP messages with unknown type.  
Partially: `icmp_input` ignores these messages and passes them to `rip_input`, which delivers the message to any waiting processes or silently discards the message if no process is prepared to receive the message.



- *May* include more than 8 bytes of the original datagram.  
No: The `icmp_error` function returns only a maximum of 8 bytes of the original datagram in the ICMP error message, Exercise 11.9.
- *Must* return the header and data unchanged from the received datagram.  
Partially: Net/3 converts the ID, offset, and length fields of an IP packet from network byte order to host byte order in `ipintr`. This facilitates processing the packet, but Net/3 neglects to return the offset and length fields to network byte order before including the header in an ICMP error message. If the system operates with the same byte ordering as the network, this error is harmless. If it operates with a different ordering, the IP header contained within the ICMP error message has incorrect offset and length values.

The authors found that an Intel implementation of SVR4 and AIX 3.2 (Net/2 based) both return the length byte-swapped. Implementations other than Net/2 or Net/3 that were tried (Cisco, NetBlazer, VM, and Solaris 2.3) did not have this bug.

Another error occurs when an ICMP port unreachable error is sent from the UDP code: the header length of the received datagram is changed incorrectly (Section 23.7). The authors found this error in Net/2 and Net/3 implementations. Net/1, however, did not have the bug.

- *Must* demultiplex received ICMP error message to transport protocol.  
Yes: `icmp_error` uses the protocol field from the original header to select the appropriate transport protocol to respond to the error.
- *Should* send ICMP error messages with a TOS field of 0.  
Yes: All ICMP error messages are constructed with a TOS of 0 by `icmp_error`.
- *Must not* send an ICMP error message caused by a previous ICMP error message.  
Partially: `icmp_error` sends an error for an ICMP redirect message, which Section 3.2.2 of RFC 1122 classifies as an ICMP error message.
- *Must not* send an ICMP error message caused by an IP broadcast or IP multicast datagram.  
No: `icmp_error` does not check for this case.

The `icmp_error` function from the original Deering multicast code for BSD checks for this case.

- *Must not* send an ICMP error message caused by an link-layer broadcast.  
Yes: `icmp_error` discards ICMP messages that arrive as link-layer broadcasts or multicasts.
- *Must not* send an ICMP error message caused by a noninitial fragment.  
Yes: `icmp_error` discards errors generated in this case.
- *Must not* send an ICMP error message caused by a datagram with nonunique source address.

Yes: `icmp_reflect` checks for experimental and multicast addresses. `ip_output` discards messages sent from a broadcast address.

- *Must* return ICMP error messages when not prohibited.  
Partially: In general, Net/3 sends appropriate ICMP error messages. It fails to send an ICMP reassembly timeout message at the appropriate time (Exercise 10.1).
- *Should* generate ICMP destination unreachable (protocol and port).  
No: Datagrams for unsupported protocols are delivered to `rip_input` where they are silently discarded if there are no processes registered to accept the datagrams. UDP generates an ICMP port unreachable error.
- *Must* pass ICMP destination unreachable to higher layer.  
Yes: `icmp_input` passes the message to the `prctlinput` function defined for the protocol (`udp_ctlinput` and `tcp_ctlinput` for UDP and TCP, respectively).
- *Should* respond to destination unreachable error.  
See Sections 23.9 and 27.6.
- *Must* interpret destination unreachable as only a hint, as it may indicate a transient condition.  
See Sections 23.9 and 27.6.
- *Must not* send an ICMP redirect when configured as a host.  
Yes: `ip_forward`, the only function that detects and sends redirects, is not called unless the system is configured as a router.
- *Must* update route cache when an ICMP redirect is received.  
Yes: `ipintr` calls `rtredirect` to process the message.
- *Must* handle both host and network redirects. Furthermore, network redirects must be treated as host redirects.  
Yes: `ipintr` calls `rtredirect` for both types of messages.
- *Should* discard illegal redirects.  
Yes: `rtredirect` discards illegal redirects (Section 19.7).
- *May* send source quench if memory is unavailable.  
Yes: `ip_forward` sends a source quench if `ip_output` returns `ENOBUFS`. This occurs when there is a shortage of mbufs or when an interface output queue is full.
- *Must* pass source quench to higher layer.  
Yes: `icmp_input` passes source quench errors to the transport layers.
- *Should* respond to source quench in higher layer.  
See Sections 23.9 and 27.6 for UDP and TCP processing. Neither ICMP nor IGMP

accept ICMP error messages (they don't define a `pr_err` function), in which case they are discarded by IP.

- *Must* pass time exceeded error to transport layer.  
Yes: `icmp_input` passes this message to the transport layers.
- *Should* send parameter problem errors.  
Yes: `ip_dooptions` complains about incorrectly formed options.
- *Must* pass parameter problem errors to transport layer.  
Yes: `icmp_input` passes parameter problem errors to the transport layer.
- *May* report parameter problem errors to process.  
See Sections 23.9 and 27.6 for UDP and TCP processing. Neither ICMP nor IGMP accept ICMP error messages.
- *Must* support an echo server and *should* support an echo client.  
Yes: `icmp_input` implements the echo server and the `ping` program implements the echo client using a raw IP socket.
- *May* discard echo requests to a broadcast address.  
No: The reply is sent by `icmp_reflect`.
- *May* discard echo request to multicast address.  
No: Net/3 responds to multicast echo requests. Both `icmp_reflect` and `ip_output` permit multicast destination addresses.
- *Must* use specific destination address as echo reply source.  
Yes: `icmp_reflect` converts a broadcast or multicast destination to the specific address of the receiving interface and uses the result as the source address for the echo reply.
- *Must* return echo request data in echo reply.  
Yes: The data portion of the echo request is not altered by `icmp_reflect`.
- *Must* pass echo reply to higher layer.  
Yes: ICMP echo replies are passed to `rip_input` for receipt by registered processes.
- *Must* reflect record route and timestamp options in ICMP echo request message.  
Yes: `icmp_reflect` includes the record route and timestamp options in the echo reply message.
- *Must* reverse and reflect source route option.  
Yes: `icmp_reflect` retrieves the reversed source route with `ip_srcroute` and includes it in the outgoing echo reply.

- *Should not* support the ICMP information request or reply.  
Partially: The kernel does not generate or respond to either message, but a process may send or receive the messages through the raw IP mechanism.
- *May* implement the ICMP timestamp request and timestamp reply messages.  
Yes: `icmp_input` implements the timestamp server functionality. The timestamp client may be implemented through the raw IP mechanism.
- *Must* minimize timestamp delay variability (if implementing the timestamp messages).  
Partially: The receive timestamp is applied after the message is taken off the IP input queue and the transmit timestamp is applied before the message is placed in the interface output queue.
- *May* silently discard broadcast timestamp request.  
No: `icmp_input` responds to broadcast timestamp requests.
- *May* silently discard multicast timestamp requests.  
No: `icmp_input` responds to broadcast timestamp requests.
- *Must* use specific destination address as timestamp reply source address.  
Yes: `icmp_reflect` converts a broadcast or multicast destination to the specific address of the receiving interface and uses the result as the source address for the timestamp reply.
- *Should* reflect record route and timestamp options in an ICMP timestamp request.  
Yes: `icmp_reflect` includes the record route and timestamp options in the timestamp reply message.
- *Must* reverse and reflect source route option in ICMP timestamp request.  
Yes: `icmp_reflect` retrieves the reversed source route with `ip_srcroute` and includes it in the outgoing timestamp reply.
- *Must* pass timestamp reply to higher layer.  
Yes: ICMP timestamp replies are passed to `rip_input` for receipt by registered processes.
- *Must* obey rules for standard timestamp value.  
Yes: `icmp_input` calls `iptime`, which returns a standard time value.
- *Must* provide a configurable method for selecting the address mask selection method for an interface.  
No: Net/3 supports only static configuration of address masks through the `ifconfig` program.

- *Must* support static configuration of address mask.  
Yes: This is accomplished indirectly by specifying static information when the `ifconfig` program configures an interface during system initialization, typically in the `/etc/netstart` start-up script.
- *May* get address mask dynamically during system initialization.  
No: Net/3 does not support the use of BOOTP or DHCP to acquire address mask information.
- *May* get address with an ICMP address mask request and reply messages.  
No: Net/3 does not support the use ICMP messages to acquire address mask information.
- *Must* retransmit address mask request if no reply.  
Not Applicable: Not required since this method is not implemented by Net/3.
- *Should* assume default mask if no reply is received.  
Not Applicable: Not required since this method is not implemented by Net/3.
- *Must* update address mask from first reply only.  
Not Applicable: Not required since this method is not implemented by Net/3.
- *Should* perform reasonableness check on any installed address mask.  
No: Net/3 performs no reasonableness check on address masks.
- *Must not* send unauthorized address mask reply messages and *must* be explicitly configured to be agent.  
Yes: `icmp_input` only responds to address mask requests if `icmpmaskrepl` is nonzero (it defaults to 0).
- *Should* support an associated address mask authority flag with each static address mask configuration.  
No: Net/3 consults a global authority flag (`icmpmaskrepl`) to determine if it should send address mask replies for *any* interface.
- *Must* broadcast address mask reply when initialized.  
No: Net/3 does not broadcast an address mask reply when an interface is configured.

## C.6 Multicasting Requirements

This section summarizes the IP multicast requirements from Section 3.5 of RFC 1122 and the compliance of the Net/3 code that we've examined to those requirements.



- *Should* support local IP multicasting (RFC 1112).  
Yes: Net/3 supports IP multicasting.
- *Should* join the all-hosts group at start-up.  
Yes: `in_ifinit` joins the all-hosts group while initializing an interface.
- *Should* provide a mechanism for higher layers to discover an interface's IP multicast capability.  
Yes: The `IFF_MULTICAST` flag in the interface's `ifnet` structure is available directly to kernel code and by the `STOOGIFFLAGS` command for processes.

## C.7 IGMP Requirements

This section summarizes the IGMP requirements from Section 3.5 of RFC 1122 and the compliance of the Net/3 code that we've examined to those requirements.

- *May* support IGMP (RFC 1112).  
Yes: Net/3 supports IGMP.

## C.8 Routing Requirements

This section summarizes the routing requirements from Section 3.5 of RFC 1122 and the compliance of the Net/3 code that we've examined to those requirements. Be aware that the requirements of this RFC apply to a host and not necessarily the kernel implementation. Some items are not explicitly handled by the kernel routing functions in Net/3, but they are expected to be provided by a routing daemon such as `routed` or `gated`.

- *Must* use address mask in determining whether a datagram's destination is on a connected network.  
Yes: When an interface for a connected network such as an Ethernet is configured, its address mask is specified (or a default is chosen based on the class of IP address) and stored in the routing table entry. This mask is used by `rn_match` when it checks a leaf for a network match.
- *Must* operate correctly in a minimal environment when there are no routers (all networks are directly connected).  
Yes: The system administrator must not configure a default route in this case.
- *Must* keep a "route cache" of mappings to next-hop routers.  
Yes: The routing table is the cache.

- *Should* treat a received network redirect the same as a host redirect.  
Yes, as described in Section 19.7.
- *Must* use a default router when no entry exists for the destination in the routing table.  
Yes, if a default route has been entered into the routing table.
- *Must* support multiple default routers.  
Multiple defaults are not supported by the kernel. Instead, this should be provided by a routing daemon.
- *May* implement a table of static routes.  
Yes: These can be created at system initialization time with the `route` command.
- *May* include a flag with each static route specifying whether or not the route can be overridden by a redirect.  
No.
- *May* allow the routing table key to be a complete host address and not just a network address.  
Yes: Host routes take priority over a network route to the same network.
- *Should* include the TOS in the routing table entry.  
No: There is a TOS field in the `sockaddr_inarp` that we describe in the next chapter, but it is not currently used.
- *Must* be able to detect the failure of a next-hop router that appears as the gateway field in the routing table and be able to choose an alternate next-hop router.  
Negative advice, the `RTM_LOSING` message generated by `in_losing`, is passed to any processes reading from a routing socket, which allows the process (e.g., a routing daemon) to handle this event.
- *Should not* assume that a route is good forever.  
Yes: There are no timeouts on routing table entries in the kernel other than those created by ARP. Again, the standard Unix routing daemons time out routes and replace them with alternatives when possible.
- *Must not* ping routers continuously (ICMP echo request).  
Yes: The Net/3 kernel does not do this. The routing daemons don't generate ICMP echo requests either.
- *Must* use ping of a router only when traffic is being sent to that router.  
The Net/3 kernel never generates pings to a next-hop router.
- *Should* allow higher and lower layers to give positive and negative advice.  
Partially: The only information passed by other layers to the Net/3 routing functions

is by `in_losing`, which is called only from TCP. The only action performed by the routing layer is to generate the `RTM_LOSING` message.

- *Must* switch to another default router when the existing default fails.  
Yes, although the Net/3 kernel does not do this, it is supported by the routing daemons.
- *Must* allow the following information to be configured manually in the routing table: IP address, network mask, list of defaults.  
Yes, but only one default is supported in the kernel.

## C.9 ARP Requirements

This section summarizes the ARP requirements from Section 2.5 of RFC 1122 and the compliance of the Net/3 code that we've examined to those requirements.

- *Must* provide a mechanism to flush out-of-date ARP entries. If this mechanism involves a timeout, it *should* be configurable.  
Yes and yes: `arptimer` provides this mechanism. The timeout is configurable (the `arpt_prune` and `arpt_keep` globals) but the only ways to change their values are to recompile the kernel or modify the kernel with a debugger.
- *Must* include a mechanism to prevent ARP flooding.  
Yes, as we described with Figure 21.24.
- *Should* save (rather than discard) at least one (the latest) packet of each set of packets destined to the same unresolved IP address, and transmit the saved packet when the address has been resolved.  
Yes: This is the purpose of the `la_hold` member of the `llinfo_arp` structure.

## C.10 UDP Requirements

This section summarizes the UDP requirements from Section 4.1.5 of RFC 1122 and the compliance of the Net/3 code that we've examined to those requirements.

- *Should* send ICMP port unreachable.  
Yes: `udp_input` does this.
- *Must* pass received IP options to application.  
No: The code to do this is commented out in `udp_input`. This means that a process that receives a UDP datagram with a source route option cannot send a reply using the reversed route.
- *Must* allow application to specify IP options to send.  
Yes: The `IP_OPTIONS` socket option does this. The options are saved in the PCB and placed into the outgoing IP datagram by `ip_output`.

- *Should* treat a received network redirect the same as a host redirect.  
Yes, as described in Section 19.7.
- *Must* use a default router when no entry exists for the destination in the routing table.  
Yes, if a default route has been entered into the routing table.
- *Must* support multiple default routers.  
Multiple defaults are not supported by the kernel. Instead, this should be provided by a routing daemon.
- *May* implement a table of static routes.  
Yes: These can be created at system initialization time with the `route` command.
- *May* include a flag with each static route specifying whether or not the route can be overridden by a redirect.  
No.
- *May* allow the routing table key to be a complete host address and not just a network address.  
Yes: Host routes take priority over a network route to the same network.
- *Should* include the TOS in the routing table entry.  
No: There is a TOS field in the `sockaddr_inarp` that we describe in the next chapter, but it is not currently used.
- *Must* be able to detect the failure of a next-hop router that appears as the gateway field in the routing table and be able to choose an alternate next-hop router.  
Negative advice, the `RTM_LOSING` message generated by `in_losing`, is passed to any processes reading from a routing socket, which allows the process (e.g., a routing daemon) to handle this event.
- *Should not* assume that a route is good forever.  
Yes: There are no timeouts on routing table entries in the kernel other than those created by ARP. Again, the standard Unix routing daemons time out routes and replace them with alternatives when possible.
- *Must not* ping routers continuously (ICMP echo request).  
Yes: The Net/3 kernel does not do this. The routing daemons don't generate ICMP echo requests either.
- *Must* use ping of a router only when traffic is being sent to that router.  
The Net/3 kernel never generates pings to a next-hop router.
- *Should* allow higher and lower layers to give positive and negative advice.  
Partially: The only information passed by other layers to the Net/3 routing functions



is by `in_losing`, which is called only from TCP. The only action performed by the routing layer is to generate the `RTM_LOSING` message.

- *Must* switch to another default router when the existing default fails.  
Yes, although the Net/3 kernel does not do this, it is supported by the routing daemons.
- *Must* allow the following information to be configured manually in the routing table: IP address, network mask, list of defaults.  
Yes, but only one default is supported in the kernel.

## C.9 ARP Requirements

This section summarizes the ARP requirements from Section 2.5 of RFC 1122 and the compliance of the Net/3 code that we've examined to those requirements.

- *Must* provide a mechanism to flush out-of-date ARP entries. If this mechanism involves a timeout, it *should* be configurable.  
Yes and yes: `arptimer` provides this mechanism. The timeout is configurable (the `arpt_prune` and `arpt_keep` globals) but the only ways to change their values are to recompile the kernel or modify the kernel with a debugger.
- *Must* include a mechanism to prevent ARP flooding.  
Yes, as we described with Figure 21.24.
- *Should* save (rather than discard) at least one (the latest) packet of each set of packets destined to the same unresolved IP address, and transmit the saved packet when the address has been resolved.  
Yes: This is the purpose of the `la_hold` member of the `llinfo_arp` structure.

## C.10 UDP Requirements

This section summarizes the UDP requirements from Section 4.1.5 of RFC 1122 and the compliance of the Net/3 code that we've examined to those requirements.

- *Should* send ICMP port unreachable.  
Yes: `udp_input` does this.
- *Must* pass received IP options to application.  
No: The code to do this is commented out in `udp_input`. This means that a process that receives a UDP datagram with a source route option cannot send a reply using the reversed route.
- *Must* allow application to specify IP options to send.  
Yes: The `IP_OPTIONS` socket option does this. The options are saved in the PCB and placed into the outgoing IP datagram by `ip_output`.



- *Must* pass IP options down to IP layer.  
Yes: As mentioned above, IP places the options into the IP datagram.
- *Must* pass received ICMP messages to application.  
Yes: We must look at the exact wording from the RFC: "A UDP-based application that wants to receive ICMP error messages is responsible for maintaining the state necessary to demultiplex these messages when they arrive; for example, the application may keep a pending receive operation for this purpose." The state required by Berkeley-derived systems is that the socket be connected to the foreign address and port. As the comments at the beginning of Figure 23.26 indicate, some applications create both a connected and an unconnected socket for a given foreign port, using the connected socket to receive asynchronous errors.
- *Must* be able to generate and verify UDP checksum.  
Yes: This is done by `udp_input`, based on the global integer `udpcksum`.
- *Must* silently discard datagrams with bad checksum.  
Yes: This is done only if `udpcksum` is nonzero. As we mentioned earlier, this variable controls both the sending of checksums and the verification of received checksums. If this variable is 0, the kernel does not verify a received nonzero checksum.
- *May* allow sending application to specify whether outgoing checksum is calculated, but *must* default to on.  
No: The application has no control over UDP checksums. Regarding the default, UDP checksums are generated unless the kernel is compiled with 4.2BSD compatibility defined, or unless the administrator has disabled UDP checksums using `sysctl(8)`.
- *May* allow receiving application to specify whether received UDP datagrams without a checksum (i.e., the received checksum is 0) are discarded or passed to the application.  
No: Received datagrams with a checksum field of 0 are passed to the receiving process.
- *Must* pass destination IP address to application.  
Yes: The application must call `recvmsg` and specify the `IP_RECVDSTADDR` socket option. Also recall our discussion following Figure 23.25 noting that 4.4BSD broke this option when the destination address is a multicast or broadcast address.
- *Must* allow application to specify local IP address to be used when sending a UDP datagram.  
Yes: The application can call `bind` to set the local IP address. Recall our discussion at the end of Section 22.8 about the difference between the source IP address and the IP address of the outgoing interface. Net/3 does not allow the application to choose the outgoing interface—that is done by `ip_output`, based on the route to the destination IP address.
- *Must* allow application to specify wildcard local IP address.  
Yes: If the IP address `INADDR_ANY` is specified in the call to `bind`, the local IP address is chosen by `in_pcbsconnect`, based on the route to the destination.

- *Should* allow application to learn of the local address that was chosen.  
Yes: The application must call `connect`. When a datagram is sent on an unconnected socket with a wildcard local address, `ip_output` chooses the outgoing interface, which also becomes the source address. The `inp_laddr` member of the PCB, however, is restored to the wildcard address at the end of `udp_output` before `sendto` returns. Therefore, `getsockname` cannot return the value. But the application can `connect` a UDP socket to the destination, causing `in_pcbeconnect` to determine the local interface and store the address in the PCB. The application can then call `getsockname` to fetch the IP address of the local interface.
- *Must* silently discard a received UDP datagram with an invalid source IP address (broadcast or multicast).  
No: A received UDP datagram with an invalid source address is delivered to a socket, if a socket is bound to the destination port.
- *Must* send a valid IP source address.  
Yes: If the local IP address is set by `bind`, it checks the validity of the address. If the local IP address is wildcarded, `ip_output` chooses the local address.
- *Must* provide the full IP interface from Section 3.4 of RFC 1122.  
Refer to Section C.2.
- *Must* allow application to specify TTL, TOS, and IP options for output datagrams.  
Yes: The application can use the `IP_TTL`, `IP_TOS`, and `IP_OPTIONS` socket options.
- *May* pass received TOS to application.  
No: There is no way for the application to receive this value from the IP header. Notice that a `getsockopt` of `IP_TOS` returns the value used in outgoing datagrams, not the value from a received datagram. The received `ip_tos` value is available to `udp_input`, but is discarded along with the entire IP header.

## C.11 TCP Requirements

This section summarizes the TCP requirements from Section 4.2.5 of RFC 1122 and the compliance of the Net/3 code that we've examined to those requirements.

### PSH Flag

- *May* aggregate data sent by the user without the PSH flag.  
Yes and no: Net/3 does not give the process a way to specify the PSH flag with a write operation, but Net/3 does aggregate data sent by the user in separate write operations.
- *May* queue data received without the PSH flag.  
No: The absence or presence of a PSH flag in a received datagram makes no difference. Received data is placed onto the socket's received queue when it is processed.

- Sender *should* collapse successive PSH flags when it packetizes data.  
No.
- *May* implement PSH flag on write calls.  
No: This is not part of the sockets API.
- Since the PSH flag is not part of the write calls, *must not* buffer data indefinitely and *must* set the PSH flag in the last buffered segment.  
Yes: This is the method used by Berkeley-derived implementations.
- *May* pass received PSH flag to application.  
No: This is not part of the sockets API.
- *Should* send maximum-sized segment whenever possible, to improve performance.  
Yes.

### Window

- *Must* treat window size as an unsigned number. *Should* treat window size as 32-bit value.  
Yes: All the window sizes in Figure 24.13 are unsigned longs, which is also required by the window scale option of RFC 1323.
- Receiver *must not* shrink the window (move the right edge to the left).  
Yes, in Figure 26.29.
- Sender *must* be robust against window shrinking.  
Yes, in Figure 29.15.
- *May* keep offered receive window closed indefinitely.  
Yes.
- Sender *must* probe a zero window.  
Yes, this is the purpose of the persist timer.
- *Should* send first zero-window probe when the window has been closed for the RTO.  
No: Net/3 sets a lower bound for the persist timer of 5 seconds, which is normally greater than the RTO.
- *Should* exponentially increase the interval between successive probes.  
Yes, as shown in Figure 25.14.
- *Must* allow peer's window to stay closed indefinitely.  
Yes, TCP never gives up probing a closed window.
- Sender *must not* timeout a connection just because the other end keeps advertising a zero window.  
Yes.

### Urgent Data

- *Must* have urgent pointer point to last byte of urgent data.  
No: Berkeley-derived implementations continue to interpret the urgent pointer as pointing just beyond the last byte of urgent data.

- *Must* support a sequence of urgent data of any length.  
Yes, with the bug fix discussed in Exercise 26.6.
- *Must* inform the receiving process (1) when TCP receives an urgent pointer and there was no previously pending urgent data, or (2) when the urgent pointer advances in the data stream.  
Yes, in Figure 29.17.
- *Must* be a way for the process to determine how much urgent data remains, or at least whether more urgent data remains to be read.  
Yes, this is the purpose of the out-of-band mark, the `SIOCATMARK` ioctl.

### TCP Options

- *Must* be able to receive TCP options in any segment.  
Yes:
- *Must* ignore any options not supported.  
Yes, in Section 28.3.
- *Must* cope with an illegal option length.  
Yes, in Section 28.3.
- *Must* implement both sending and receiving the MSS option.  
Yes, a received MSS option is handled in Figure 28.10, and Figure 26.23 always sends an MSS option with a SYN.
- *Should* send an MSS option in every SYN when its receive MSS differs from 536, and *may* send it always.  
Yes, as mentioned earlier, an MSS option is always sent by Net/3 with a SYN.
- If an MSS option is not received with a SYN, *must* assume a default MSS of 536.  
No: The default MSS is 512, not 536.  

This is probably a historical artifact because VAXes had a physical page size of 512 bytes and trailer protocols working only with data that is a multiple of 512.
- *Must* calculate the "effective send MSS."  
Yes, in Section 27.5.

### TCP Checksums

- *Must* generate a TCP checksum in outgoing segments and *must* verify received checksums.  
Yes, TCP checksums are always calculated and verified.

### Initial Sequence Number Selection

- *Must* use the specified clock-driven selection from RFC 793.  
No: RFC 793 specifies a clock that changes by 125,000 every half-second, whereas



the Net/3 ISN (the global variable `tcp_isn`) is incremented by 64,000 every half-second, about one-half the specified rate.

### Opening Connections

- *Must* support simultaneous open attempts.  
Yes, although Berkeley-derived systems prior to 4.4BSD did not support this, as described in Section 28.9.
- *Must* keep track of whether it reached the SYN\_RCVD state from the LISTEN or SYN\_SENT states.  
Yes, same result, different technique. The purpose of this requirement is to allow a passive open that receives an RST to return to the LISTEN state (as shown in Figure 24.15), but force an active open that ends up in SYN\_RCVD and then receives an RST to be aborted. This is described following Figure 28.36.
- A passive open *must not* affect previously created connections.  
Yes:
- *Must* allow a listening socket with a given local port at the same time that another socket with the same local port is in the SYN\_SENT or SYN\_RCVD state.  
Yes: The stated purpose of this requirement is to allow a given application to accept multiple connection attempts at about the same time. This is done in Berkeley-derived implementations by cloning new connections from the socket in the LISTEN state when the incoming SYN arrives.
- *Must* ask IP to select a local IP address to be used as the source IP address when the source IP address is not specified by the process performing an active open on a multihomed host.  
Yes, done by `ip_pcbsconnect`.
- *Must* continue to use the same source IP address for all segments sent on a connection.  
Yes: Once `ip_pcbsconnect` selects the source address, it doesn't change.
- *Must not* allow an active open for a broadcast or multicast foreign address.  
Yes and no: TCP will not send segments to a broadcast address because the call to `ip_output` in Figure 26.32 does not specify the `SO_BROADCAST` option. Net/3, however, allows connection attempts to multicast addresses.
- *Must* ignore incoming SYNs with an invalid source address.  
Yes: The code in Figure 28.16 checks for these invalid source addresses.

### Closing Connections

- *Should* allow an RST to contain data.  
No: The RST processing in Figure 28.36 ends up jumping to `drop`, which skips the processing of any segment data in Figure 29.22.
- *Must* inform process whether other end closed the connection normally (e.g., sent a FIN) or aborted the connection with an RST.



Yes: The read system calls return 0 (end-of-file) when the FIN is processed, but -1 with an error of `ECONNRESET` when an RST is received.

- *May* implement a half-close.  
Yes: The process calls `shutdown` with a second argument of 1 to send a FIN. The process can still read from the connection.
- If the process completely closes a connection (i.e., not a half-close) and received data is still pending in TCP, or if new data arrives after the close, TCP *should* send an RST to indicate data was lost.  
No and yes: If a process calls `close` and unread data is in the socket's receive buffer, an RST is not sent. But if data arrives after a socket is closed, an RST is returned to the sender.
- *Must* linger in `TIME_WAIT` state for twice the MSL.  
Yes, although the Net/3 MSL of 30 seconds is much smaller than the RFC 793 recommended value of 2 minutes.
- *May* accept a new SYN from a peer to reopen a connection directly from the `TIME_WAIT` state.  
Yes, as shown in Figure 28.28.

### Retransmissions

- *Must* implement Van Jacobson's slow start and congestion avoidance.  
Yes:
- *May* reuse the same IP identifier field when a retransmission is identical to the original packet.  
No: The IP identifier is assigned by `ip_output` from the global variable `ip_id`, which increments each time an IP datagram is sent. It is not assigned by TCP.
- *Must* implement Jacobson's algorithm for calculating the RTO and Karn's algorithm for selecting the RTT measurements.  
Yes, but realize that when RFC 1323 timestamps are present, the retransmission ambiguity problem is gone, obviating half of Karn's algorithm, as we discussed with Figure 29.6.
- *Must* include an exponential backoff for successive RTO values.  
Yes, as described with Figure 25.22.
- Retransmission of SYN segments *should* use the same algorithm as data segments.  
Yes, as shown in Figure 25.16.
- *Should* initialize estimation parameters to calculate an initial RTO of 3 seconds.  
No: The initial value of `t_rxtcur` calculated by `tcp_newtcb` is 6 seconds. This is also seen in Figure 25.16.
- *Should* have a lower bound on the RTO measured in fractions of a second and an upper bound of twice the MSL.  
No: The lower bound is 1 second and the upper bound is 64 seconds (Figure 25.3).

## Generating ACKs

- *Should* queue out-of-order segments.  
Yes, done by `tcp_reass`.
- *Must* process all queued segments before sending any ACKs  
Yes, but only for in-order segments. `ipintr` calls `tcp_input` for each queued datagram that is a TCP segment. For in-order segments, `tcp_input` schedules a delayed ACK and returns to `ipintr`. If there are additional TCP segments on IP's input queue, `tcp_input` is called by `ipintr` for each one. Only when `ipintr` finds no more IP datagrams on its input queue and returns can `tcp_fasttimo` be called to generate a delayed ACK. This ACK will contain the highest acknowledgment number in all the segments processed by `tcp_input`.  
The problem is with out-of-order segments: `tcp_input` calls `tcp_output` itself, before returning to `ipintr`, to generate the ACK for the out-of-order segment. If there are additional segments on IP's input queue that would have made the out-of-order segment be in order, they are processed after the immediate ACK is sent.
- *May* generate an immediate ACK for an out-of-order segment.  
Yes, this is needed for the fast retransmit and fast recovery algorithms (Section 29.4).
- *Should* implement delayed ACKs and the delay *must* be less than 0.5 seconds.  
Yes: The `TF_DELACK` flag is checked by the `tcp_fasttimo` function every 200 ms.
- *Should* send an ACK for at least every second segment.  
Yes, the code in Figure 26.9 generates an ACK for every second segment. We also discussed that this happens only if the process receiving the data reads the data as it arrives, since the calls to `tcp_output` that cause every other segment to be acknowledged are driven by the `PRU_RCVD` request.
- *Must* include silly window syndrome avoidance in the receiver.  
Yes, as seen in Figure 26.29.

## Sending Data

- The TTL value for TCP segments *must* be configurable.  
Yes: The TTL is initialized to 64 (`IPDEFTTL`) by `tcp_newtcpcb`, but can then be changed by a process using the `IP_TTL` socket option.
- *Must* include sender silly window syndrome avoidance.  
Yes, in Figure 26.8.
- *Should* implement the Nagle algorithm.  
Yes, in Figure 26.8.
- *Must* allow a process to disable the Nagle algorithm on a given connection.  
Yes, with the `TCP_NODELAY` socket option.

### Connection Failures

- *Must* pass negative advice to IP when the number of retransmissions for a given segment exceeds some value R1.  
Yes: The value of R1 is 4, and in Figure 25.26, when the number of retransmissions exceeds 4, `in_lodging` is called.
- *Must* close a connection when the number of retransmissions for a given segment exceeds some value R2.  
Yes: The value of R2 is 12 (Figure 25.26).
- *Must* allow process to set the value of R2.  
No: The value 12 is hardcoded in Figure 25.26.
- *Should* inform the process when R1 is reached and before R2 is reached.  
No:
- *Should* default R1 to at least 3 retransmissions and R2 to at least 100 seconds.  
Yes: R1 is 4 retransmissions, and with a minimum RTO of 1 second, the `tcp_backoff` array (Section 25.9) guarantees a minimum value of R2 of over 500 seconds.
- *Must* handle SYN retransmissions in the same general way as data retransmissions.  
Yes, but R1 is normally not reached for the retransmission of a SYN (Figure 25.16).
- *Must* set R2 to at least 3 minutes for a SYN.  
No: R2 for a SYN is limited to 75 seconds by the connection-establishment timer (Figure 25.16).

### Keepalive Packets

- *May* provide keepalives.  
Yes, they are provided.
- *Must* allow process to turn keepalives on or off, and *must* default to off.  
Yes: Default is off and process must turn them on with the `SO_KEEPAVIVE` socket option.
- *Must* send keepalives only when connection is idle for a given period.  
Yes:
- *Must* allow the keepalive interval to be configurable and *must* default to no less than 2 hours.  
No and yes: The idle time before sending keepalive probes is not easily configurable, but it defaults to 2 hours. If the default idle time is changed (by changing the global variable `tcp_keepidle`), it affects all users of the keepalive option on the host—it cannot be configured on a per-connection basis as many users would like.
- *Must not* interpret the failure to respond to any given probe as a dead connection.  
Yes: Nine probes are sent before the connection is considered dead.

## IP Options

- *Must* ignore received IP options it doesn't understand.  
Yes: This is done by the IP layer.
- *May* support the timestamp and record route options in received segments.  
No: Net/3 only reflects these options for ICMP packets that are reflected back to the sender (`icmp_reflect`). `tcp_input` discards any received IP options by calling `ip_stripoptions` in Figure 28.2.
- *Must* allow process to specify a source route when a connection is actively opened, and this route must take precedence over a source route received for this connection.  
Yes: The source route is specified with the `IP_OPTIONS` socket option. `tcp_input` never looks at a received source route when the connection is actively opened.
- *Must* save a received source route in a connection that is passively opened and use the return route for all segments sent on this connection. If a different source route arrives in a later segment, the later route *should* override the earlier one.  
Yes and no: Figure 28.7 calls `ip_srcroute`, but only when the SYN arrives for a listening socket. If a different source route arrives later, it is not used.

## Receiving ICMP Messages from IP

- Receipt of an ICMP source quench *should* trigger slow start.  
Yes: The function `tcp_quench` is called by `tcp_ctlinput`.
- Receipt of a network unreachable, host unreachable, or source route failed *must not* cause TCP to abort the connection and the process *should* be informed.  
Yes and no: As described following Figure 27.12, Net/3 now completely ignores host unreachable and network unreachable errors for an established connection.
- Receipt of a protocol unreachable, port unreachable, or fragmentation required and DF bit set *should* abort an existing connection.  
No: `tcp_notify` records these ICMP error in `t_softerror`, which is reported to the process if the connection is eventually dropped.
- *Should* handle time exceeded and parameter problem errors the same as required previously for network and host unreachable.  
Yes: ICMP parameter problem errors are just recorded in `t_softerror` by `tcp_notify`. ICMP time exceeded errors are ignored by `tcp_ctlinput`. Neither type of ICMP error causes the connection to be aborted.

## Application Programming Interface

- *Must* be a method for reporting soft errors to the process, normally in an asynchronous fashion.  
No: Soft errors are returned to the process if the connection is aborted.

- *Must* allow process to specify TOS for segments sent on a connection. *Should* let application change this during a connection's lifetime.  
Yes to both, with the `IP_TOS` socket option.
- *May* pass most recently received TOS to process.  
No: There is no way to do this with the sockets API. Calling `getsockopt` for `IP_TOS` returns only the current value being sent; it does not return the most recently received value.
- *May* implement a "flush" call.  
No: TCP sends the data from the process as quickly as it can.
- *Must* allow process to specify local IP address before either an active open or a passive open.  
Yes: This is done by calling `bind` before either `connect` or `accept`.





## Bibliography

All the RFCs are available at no charge through electronic mail or by using anonymous FTP across the Internet as described in Appendix B.

Whenever the authors were able to locate an electronic copy of papers and reports referenced in this bibliography, its URL (Uniform Resource Locator, Appendix B) is included.

Almquist, P. 1992. "Type of Service in the Internet Protocol Suite," RFC 1349, 28 pages (July).

Almquist, P., and Kastenholz, F. J. 1994. "Towards Requirements for IP Routers," RFC 1716, 186 pages (Nov.).

This RFC is an intermediate step to replace RFC 1009 [Braden and Postel 1987].

Auerbach, K. 1994. "Max IP Packet Length and MTU," Message-ID <karl.3.000A4DD7@cavebear.com>, Usenet, comp.protocols.tcp-ip Newsgroup (July).

Boggs, D. R. 1982. "Internet Broadcasting," Xerox PARC CSL-83-3, Stanford University, Palo Alto, Calif. (Jan.).

Braden, R. T., ed. 1989a. "Requirements for Internet Hosts—Communication Layers," RFC 1122, 116 pages (Oct.).

The first half of the Host Requirements RFC. This half covers the link layer, IP, TCP, and UDP.

Braden, R. T., ed. 1989b. "Requirements for Internet Hosts—Application and Support," RFC 1123, 98 pages (Oct.).

The second half of the Host Requirements RFC. This half covers Telnet, FTP, TFTP, SMTP, and the DNS.

Braden, R. T. 1989c. "Perspective on the Host Requirements RFCs," RFC 1127, 20 pages (Oct.).

An informal summary of the discussions and conclusions of the IETF working group that developed the Host Requirements RFC.

- Braden, R. T. 1992. "TIME-WAIT Assassination Hazards in TCP," RFC 1337, 11 pages (May).  
Shows how the receipt of an RST while in the TIME\_WAIT state can lead to problems.
- Braden, R. T. 1993. "TCP Extensions for High Performance: An Update," Internet Draft, 10 pages (June).  
This is an update to RFC 1323 [Jacobson, Braden, and Borman 1992].
- Braden, R. T. 1994. "T/TCP—TCP Extensions for Transactions, Functional Specification," RFC 1644, 38 pages (July).
- Braden, R. T., Borman, D. A., and Partridge, C. 1988. "Computing the Internet Checksum," RFC 1071, 24 pages (Sept.).  
Provides techniques and algorithms for calculating the checksum used by IP, ICMP, IGMP, UDP, and TCP.
- Braden, R. T., and Postel, J. B. 1987. "Requirements for Internet Gateways," RFC 1009, 55 pages (June).  
The equivalent of the Host Requirements RFC for routers. This RFC is being replaced by RFC 1710 [Almquist and Kastenholz 1994].
- Brakmo, L. S., O'Malley, S. W., and Peterson, L. L. 1994. "TCP Vegas: New Techniques for Congestion Detection and Avoidance," *Computer Communication Review*, vol. 24, no. 4, pp. 24–35 (Oct.).  
Describes modifications to the 4.3BSD Reno TCP implementation to improve throughput and reduce retransmissions.  
<http://ftp.cs.arizona.edu/skernel/Papers/vegas.ps>
- Carlson, J. 1993. "Re: Bug in Many Versions of TCP," Message-ID <1993jul12.130854.26176@xylogics.com>, Usenet, comp.protocols.tcp-ip Newsgroup (July).
- Casner, S., *Frequently Asked Questions (FAQ) on the Multicast Backbone (MBONE)*, 1993.  
<http://ftp.isi.edu/mbone/faq.txt>
- Cheswick, W. R., and Bellovin, S. M. 1994. *Firewalls and Internet Security: Repelling the Wily Hacker*. Addison-Wesley, Reading, Mass.  
Describes how to set up and administer a firewall gateway and the security issues involved.
- Clark, D. D. 1982. "Modularity and Efficiency in Protocol Implementation," RFC 817, 26 pages (July).
- Comer, D. E., and Lin, J. C. 1994. "TCP Buffering and Performance Over an ATM Network," Purdue Technical Report CSD-TR 94-026, Purdue University, West Lafayette, In. (Mar.).
- Comer, D. E., and Stevens, D. L. 1993. *Internetworking with TCP/IP: Vol. III: Client-Server Programming and Applications, BSD Socket Version*. Prentice-Hall, Englewood Cliffs, N.J.
- Croff, W., and Gilmore, J. 1985. "Bootstrap Protocol (BOOTP)," RFC 951, 12 pages (Sept.).
- Crowcroft, J., Wakeman, I., Wang, Z., and Sirovica, D. 1992. "Is Layering Harmful?," *IEEE Network*, vol. 6, no. 1, pp. 20–24 (Jan.).  
The seven missing figures from this paper appear in the next issue, vol. 6, no. 2 (March).
- Dalton, C., Watson, G., Banks, D., Calamvokis, C., Edwards, A., and Lumley, J. 1993. "Afterburner," *IEEE Network*, vol. 7, no. 4, pp. 36–43 (July).  
Describes how to speed up TCP by reducing the number of data copies performed, and a special-purpose interface card that supports this design.

- Deering, S. E. 1989. "Host Extensions for IP Multicasting," RFC 1112, 17 pages (Aug.).  
The specification of IP multicasting and IGMP.
- Deering, S. E., ed. 1991a. "ICMP Router Discovery Messages," RFC 1256, 19 pages (Sept.).
- Deering, S. E. 1991b. "Multicast Routing in a Datagram Internetwork," STAN-CS-92-1415, Stanford University, Palo Alto, Calif. (Dec.).  
[ftp://groccrio.stanford.edu/vmtp-ip/edthesis.pdf1.ps.z](http://groccrio.stanford.edu/vmtp-ip/edthesis.pdf1.ps.z)
- Deering, S. E., and Cheriton, D. P. 1990. "Multicast Routing in Datagram Internetworks and Extended LANs," *ACM Transactions on Computer Systems*, vol. 8, no. 2, pp. 85-110 (May).  
Proposes extensions to common routing techniques to support multicasting.
- Deering, S., Estrin, D., Farinacci, D., Jacobson, V., Liu, C., and Wei, L. 1994. "An Architecture for Wide-Area Multicast Routing," *Computer Communication Review*, vol. 24, no. 4, pp. 126-135 (Oct.).
- Droms, R. 1993. "Dynamic Host Configuration Protocol," RFC 1541, 39 pages (Oct.).
- Finlayson, R., Mann, T., Mogul, J. C., and Theimer, M. 1984. "A Reverse Address Resolution Protocol," RFC 903, 4 pages (June).
- Floyd, S. 1994. Private Communication.
- Forge, J. 1979. "ST—A Proposed Internet Stream Protocol," IEN 119, MIT Lincoln Laboratory (Sept.).
- Fuller, V., Li, T., Yu, J. Y., and Varadhan, K. 1993. "Classless Inter-Domain Routing (CIDR): An Address Assignment and Aggregation Strategy," RFC 1519, 24 pages (Sept.).
- Hornig, C. 1984. "Standard for the Transmission of IP Datagrams over Ethernet Networks," RFC 894, 3 pages (Apr.).
- Hutchinson, N. C., and Peterson, L. L. 1991. "The x-Kernel: An Architecture for Implementing Network Protocols," *IEEE Transactions on Software Engineering*, vol. 17, no. 1, pp. 64-76 (Jan.).  
[ftp://ftp.cs.arizona.edu/xkernel1/Papers/architecture.ps](http://ftp.cs.arizona.edu/xkernel1/Papers/architecture.ps)
- Itano, W. M., and Ramsey, N. F. 1993. "Accurate Measurement of Time." *Scientific American*, vol. 269, p. 56 (July).  
Overview of historical and current methods for accurate timekeeping. Includes a short discussion of international time scales including International Atomic Time (TAI) and Coordinated Universal Time (UTC).
- Jacobson, V. 1988a. "Some Interim Notes on the BSD Network Speedup," Message-ID <S807200426.AA01221@helios.ee.lbl.gov>, Usenet, comp.protocols.tcp-ip Newsgroup (July).
- Jacobson, V. 1988b. "Congestion Avoidance and Control," *Computer Communication Review*, vol. 18, no. 4, pp. 314-329 (Aug.).  
A classic paper describing the slow start and congestion avoidance algorithms for TCP.  
[ftp://ftp.ee.lbl.gov/congavoid.ps.z](http://ftp.ee.lbl.gov/congavoid.ps.z)
- Jacobson, V. 1990a. "Compressing TCP/IP Headers for Low-Speed Serial Links," RFC 1144, 43 pages (Feb.).  
Describes CSLIP, a version of SLIP with the TCP and IP headers compressed.

- Jacobson, V. 1990b. "4BSD TCP Header Prediction," *Computer Communication Review*, vol. 20, no. 2, pp. 13–15 (Apr.).
- Jacobson, V. 1990c. "Modified TCP Congestion Avoidance Algorithm," April 30, 1990, end2end-interest mailing list.  
Describes the fast retransmit and fast recovery algorithms.  
<http://ftp.isi.edu/end2end/end2end-interest-1990.maj1>
- Jacobson, V. 1990d. "Berkeley TCP Evolution from 4.3-Tahoe to 4.3-Reno," *Proceedings of the Eighteenth Internet Engineering Task Force*, p. 365 (Sept.), University of British Columbia, Vancouver, B.C.
- Jacobson, V. 1993. "Some Design Issues for High-Speed Networks," *Workshop '93* (Nov.), Melbourne, Australia.  
A set of 21 overheads.  
<http://ftp.isi.edu/pub/papers/vj/4bsd93-1.pdf.2>
- Jacobson, V., and Braden, R. T. 1988. "TCP Extensions for Long-Delay Paths," RFC 1072, 16 pages (Oct.).  
Describes the selective acknowledgment option for TCP, which was removed from the later RFC 1323, and the echo options, which were replaced with the timestamp option in RFC 1323.
- Jacobson, V., Braden, R. T., and Borman, D. A. 1992. "TCP Extensions for High Performance," RFC 1323, 37 pages (May).  
Describes the window scale option, the timestamp option, and the PAWS algorithm, along with the reasons these modifications are needed. [Braden 1993] updates this RFC.
- Jain, R., and Routhier, S. A. 1986. "Packet Trains: Measurements and a New Model for Computer Network Traffic," *IEEE Journal on Selected Areas in Communications*, vol. 4, pp. 1162–1167.
- Karels, M. J., and McKusick, M. K. 1986. "Network Performance and Management with 4.3BSD and IP/TCP," *Proceedings of the 1986 Summer USENIX Conference*, pp. 182–188, Atlanta, Ga.  
Describes the changes made from 4.2BSD to 4.3BSD with regard to TCP/IP.
- Karn, P., and Partridge, C. 1987. "Improving Round-Trip Time Estimates in Reliable Transport Protocols," *Computer Communication Review*, vol. 17, no. 5, pp. 2–7 (Aug.).  
Details of Karn's algorithm to handle the retransmission timeout for segments that have been retransmitted.  
<http://sics.se/users/traig/karn-partridge.ps>
- Kay, J., and Pasquale, J. 1993. "The Importance of Non-Data Touching Processing Overheads in TCP/IP," *Computer Communication Review*, vol. 23, no. 4, pp. 259–268 (Sept.).
- Kent, C. A., and Mogul, J. C. 1987. "Fragmentation Considered Harmful," *Computer Communication Review*, vol. 17, no. 5, pp. 390–401 (Aug.).
- Kernighan, B. W., and Plauger, P. J. 1976. *Software Tools*. Addison-Wesley, Reading, Mass.
- Krol, E. 1994. *The Whole Internet, Second Edition*. O'Reilly & Associates, Sebastopol, Calif.  
An introduction to the Internet, common Internet applications, and various resources available on the Internet.
- Krol, E., and Hoffman, E. 1993. "FYI on 'What is the Internet?'," RFC 1462, 11 pages (May).
- Lanciani, D. 1993. "Re: Bug in Many Versions of TCP," Message-ID <1993Jul10.015938.15961@burthuss.harvard.edu>, Usenet, comp.protocols.tcp-ip Newsgroup (July).



- Leffler, S. J., McKusick, M. K., Karels, M. J., and Quarterman, J. S. 1989. *The Design and Implementation of the 4.3BSD UNIX Operating System*. Addison-Wesley, Reading, Mass.  
An entire book on the 4.3BSD Unix system. This book describes the Tahoe release of 4.3BSD.
- Lynch, D. C. 1993. "Historical Perspective," in *Internet System Handbook*, eds. D. C. Lynch and M. T. Rose, pp. 3-14. Addison-Wesley, Reading, Mass.  
A historical overview of the Internet and its precursor, the ARPANET.
- Mallory, T., and Kullberg, A. 1990. "Incremental Updating of the Internet Checksum," RFC 1141, 2 pages (Jan.).  
This RFC is updated by RFC 1624 [Rijsinghani 1994].
- Mano, M. M. 1993. *Computer System Architecture, Third Edition*. Prentice-Hall, Englewood Cliffs, N.J.
- McCanne, S., and Jacobson, V. 1993. "The BSD Packet Filter: A New Architecture for User-Level Packet Capture," *Proceedings of the 1993 Winter USENIX Conference*, pp. 259-269, San Diego, Calif.  
A detailed description of the BSD Packet Filter (BPF) and comparisons with Sun's Network Interface Tap (NIT).  
<ftp://ftp.ee.lbl.gov/papers/bpf-usenix93.ps.gz>
- McCloghrie, K., and Farinacci, D. 1994a. "Internet Group Management Protocol MIB," Internet Draft, 12 pages (Jul.).
- McCloghrie, K., and Farinacci, D. 1994b. "IP Multicast Routing MIB," Internet Draft, 15 pages (Jul.).
- McCloghrie, K., and Rose, M. T. 1991. "Management Information Base for Network Management of TCP/IP-based Internets: MIB-II," RFC 1213 (Mar.).
- McGregor, G. 1992. "PPP Internet Protocol Control Protocol (IPCP)," RFC 1332, 12 pages (May).
- McKenney, P. E., and Dove, K. F. 1992. "Efficient Demultiplexing of Incoming TCP Packets," *Computer Communication Review*, vol. 22, no. 4, pp. 269-279 (Oct.).
- Mogul, J. C. 1991. "Network Locality at the Scale of Processes," *Computer Communication Review*, vol. 21, no. 4, pp. 273-284 (Sept.).
- Mogul, J. C. 1993. "IP Network Performance," in *Internet System Handbook*, eds. D. C. Lynch and M. T. Rose, pp. 575-675. Addison-Wesley, Reading, Mass.  
Covers numerous topics in the Internet protocols that are candidates for tuning to obtain optimal performance.
- Mogul, J. C., and Deering, S. E. 1990. "Path MTU Discovery," RFC 1191, 19 pages (Apr.).
- Mogul, J. C., and Postel, J. B. 1985. "Internet Standard Subnetting Procedure," RFC 950, 18 pages (Aug.).
- Moy, J. 1994. "Multicast Extensions to OSPF," RFC 1584, 102 pages (Mar.).
- Olivier, G. 1994. "What is the Diameter of the Internet?," Message-ID <1994Jan22.094832@mines.u-nancy.fr>, Usenet, comp.unix.wizards Newsgroup (Jan.).
- Partridge, C. 1987. "Implementing the Reliable Data Protocol (RDP)," *Proceedings of the 1987 Summer USENIX Conference*, pp. 367-379, Phoenix, Ariz.

- Partridge, C. 1993. "Jacobson on TCP in 30 Instructions," Message-ID <1993Sep8.213239.28992@slcs.se>, Usenet, comp.protocols.tcp-ip Newsgroup (Sept.).  
Describes a research implementation of TCP/IP being developed by Van Jacobson that reduces TCP receive packet processing down to 30 instructions on a RISC system.
- Partridge, C., and Hinden, R. 1990. "Version 2 of the Reliable Data Protocol (RDP)," RFC 1151, 4 pages (Apr.).
- Partridge, C., Mendez, T., and Milliken, W. 1993. "Host Anycasting Service," RFC 1546, 9 pages (Nov.).
- Partridge, C., and Pink, S. 1993. "A Faster UDP," *IEEE/ACM Transactions on Networking*, vol. 1, no. 4, pp. 429-440 (Aug.).  
Describes implementation improvements to the Berkeley sources to speed up UDP performance about 30%.
- Paxson, V. 1994. Private Communication.
- Perlman, R. 1992. *Interconnections: Bridges and Routers*. Addison-Wesley, Reading, Mass.
- Piscitello, D. M., and Chapin, A. L. 1993. *Open Systems Networking, TCP/IP and OSI*. Addison-Wesley, Reading, Mass.
- Plummer, D. C. 1982. "An Ethernet Address Resolution Protocol," RFC 826, 10 pages (Nov.).
- Postel, J. B., ed. 1981a. "Internet Protocol," RFC 791, 45 pages (Sept.).
- Postel, J. B. 1981b. "Internet Control Message Protocol," RFC 792, 21 pages (Sept.).
- Postel, J. B., ed. 1981c. "Transmission Control Protocol," RFC 793, 85 pages (Sept.).
- Postel, J. B. 1981d. "Service Mappings," RFC 795, 4 pages (Sept.).
- Postel, J. B., and Reynolds, J. K. 1988. "Standard for the Transmission of IP Datagrams over IEEE 802 Networks," RFC 1042, 15 pages (Apr.).
- Rago, S. A. 1993. *LINUX System V Network Programming*. Addison-Wesley, Reading, Mass.
- Reynolds, J. K., and Postel, J. B. 1994. "Assigned Numbers," RFC 1700, 230 pages (Oct.).
- Rijsinghani, A. 1994. "Computation of the Internet Checksum via Incremental Update," RFC 1624, 6 pages (May).  
An update to RFC 1141 [Mallory and Kullberg 1990].
- Romkey, J. L. 1988. "A Nonstandard for Transmission of IP Datagrams Over Serial Lines: SLIP," RFC 1055, 6 pages (June).
- Rose, M. T. 1990. *The Open Book: A Practical Perspective on OSI*. Prentice-Hall, Englewood Cliffs, N.J.
- Salus, P. H. 1994. *A Quarter Century of Unix*. Addison-Wesley, Reading, Mass.
- Sedgewick, R. 1990. *Algorithms in C*. Addison-Wesley, Reading, Mass.
- Simpson, W. A. 1993. "The Point-to-Point Protocol (PPP)," RFC 1548, 53 pages (Dec.).
- Sklower, K. 1991. "A Tree-Based Packet Routing Table for Berkeley Unix," *Proceedings of the 1991 Winter USENIX Conference*, pp. 93-99, Dallas, Tex.

- Stallings, W. 1987. *Handbook of Computer-Communications Standards, Volume 2: Local Network Standards*. Macmillan, New York.
- Stallings, W. 1993. *Networking Standards: A Guide to OSI, ISDN, LAN, and MAN Standards*. Addison-Wesley, Reading, Mass.
- Stevens, W. R. 1990. *UNIX Network Programming*. Prentice-Hall, Englewood Cliffs, N.J.
- Stevens, W. R. 1992. *Advanced Programming in the UNIX Environment*. Addison-Wesley, Reading, Mass.
- Stevens, W. R. 1994. *TCP/IP Illustrated, Volume 1: The Protocols*. Addison-Wesley, Reading, Mass.  
The first volume in this series, which provides a complete introduction to the Internet protocols.
- Tanenbaum, A. S. 1989. *Computer Networks, Second Edition*. Prentice-Hall, Englewood Cliffs, N.J.
- Topolcic, C. 1990. "Experimental Stream Protocol, Version 2 (SY-II)," RFC 1190, 148 pages (Oct).
- Torek, C. 1992. "Re: A Problem in Bind System Call," Message-ID <27240@dog.ee.lbl.gov>, Usenet, comp.unix.internals Newsgroup (Nov).
- Waitzman, D., Partridge, C., and Deering, S. E. 1988. "Distance Vector Multicast Routing Protocol," RFC 1075, 24 pages (Nov).



# Index

Rather than provide a separate glossary (with most of the entries being acronyms), this index also serves as a glossary for all the acronyms used in the book. The primary entry for the acronym appears under the acronym name. For example, all references to the Address Resolution Protocol appear under ARP. The entry under the compound term "Address Resolution Protocol" refers back to the main entry under ARP.

The two end papers at the back of the book contain a list of all the functions and macros presented or described in the text, along with the starting page number of the source code. Similarly one front end paper contains a list of all the structures presented in the text. These end papers should be the starting point to locate the definition of a function or structure.

The various functions, constants, variables, and the like that appear in this index refer to their appearance in the text. We have not attempted to index all these names when they appear in source code files that are included in the text. The definitive answer to a question such as "where are all the references to the constant `IP_RECVOPTS`" can only be obtained by obtaining the Net/3 source code (Appendix B) and using a tool such as `grep`.

The entries in this index for RFCs refer only to the reference for that RFC in the Bibliography. This is to help locate an RFC if you encounter a reference to it by number within the text.

- 224.0.0.1, 343, 345-346
- 224.0.0.2, 343, 346
- 224.0.1.2, 343, 346
- 2MSL, timer, 818-819, 821-822, 825-827, 893, 967
- 4.1cBSD, 4
- 4.2BSD, 4-6, 887, 1128
  - compatibility, 768, 1114
  - keepalives, 956
- 4.3BSD, 5, 844, 1053, 1128-1129
  - Reno, xix, 5, 191, 562, 569, 678, 776, 934, 970, 977, 1126
  - Tahoe, xix, 5, 773, 834, 970, 1129
- 4.4BSD, xix, 678, 977
- 4.4BSD-Lite, xix
- 4.5BSD, 778
- 802.3 encapsulation, 106, 125



- ac\_enaddr member, 81, 99, 111, 683
- ac\_if member, 81
- ac\_ipaddr member, 81, 183, 683
- ac\_multiaddrs member, 81, 343, 364, 366
- ac\_multicast member, 81
- accept function, 440-441, 443, 445-446, 455, 457-464, 474, 552, 555-556, 713, 727, 787, 969, 1013-1014, 1077, 1123
- accept\_args structure, 459
- acceptable ACK, 808
- access rights, 189, 470, 518
- ACK (acknowledgment flag, TCP header), 461-463, 803, 805, 887
  - acceptable, 808
  - number, 803, 807
  - pure, 831, 851, 937
  - resynchronization, 960, 1088
- acknowledgment flag, TCP header, *see* ACK
- adb program, 200
- add\_lgrp function, 401, 411-412, 1077
- add\_mrc function, 401, 422-423
- add\_vif function, 401, 408-409
- RDDDCMAIN macro, 193-194, 584
- address
  - class A IP, 155
  - class B IP, 155
  - class C IP, 155
  - class D IP, 155, 341
  - class E IP, 155
  - directed broadcast, 157, 162
  - Ethernet broadcast, 97, 100, 103
  - Ethernet destination, 99, 110
  - Ethernet hardware, 75, 81, 87, 91-92, 104, 341, 1046
  - Ethernet multicast, 100, 103-104, 341-342
  - Ethernet source, 99
  - Ethernet unicast, 100
  - IP, 155-183
  - IP broadcast, 182, 234, 1101
  - IP destination, 162, 182
  - IP experimental, 156
  - IP multicast, 155-156, 341
  - IP unicast, 155-156, 182
  - limited broadcast, 753
  - link-level, 77, 85-90, 92, 94, 97, 158, 185
  - network broadcast, 162
  - subnet, 1129
- address assignment, IP, 161-177
- address family, 75, 86, 110, 152, 162, 182, 185, 187
  - Internet, 185
  - OSI, 185
  - routing, 185
  - Unix, 185
- address mapping, IP to Ethernet multicast, 341-342
- address mask, link-level, 77
- address mask reply, ICMP, 319
- address mask request, ICMP, 319
- Address Resolution Protocol, *see* ARP
- ADVANCE macro, 661
- Advanced Research Projects Agency network, *see* ARPANET
- AF\_IMPLINK constant, 1060
- AF\_INET constant, 75, 109-110, 160, 185, 187, 192, 361, 363, 364, 577, 581, 627, 637, 647, 662, 701, 1060, 1081
- AF\_IPG constant, 75, 109, 185, 581, 585
- AF\_LINK constant, 75, 86, 90-91, 121, 185, 698, 702, 1070
- AF\_LOCAL constant, 185
- AF\_NS constant, 581
- AF\_OSI constant, 75, 185
- AF\_ROUTE constant, 75, 185
- AF\_UNIX constant, 75, 185, 581
- AF\_UNSPEC constant, 75, 109-110, 363, 577, 684, 686, 1046
- AIX, 4
- all-hosts group, 168, 170-171, 339, 343, 345, 355, 365, 379, 381, 383, 386, 391, 393, 561, 1075
  - joining, 171, 1111
  - membership report, 394
- all-routers group, 339, 343
- Almqvist, P., 140, 226, 1125-1126
- ambiguity problem, retransmission, 976
- American National Standards Institute, *see* ANSI
- ANSI (American National Standards Institute), 103
- ANSI C function prototypes, 41
- any\_count member, 666
- anycasting, 351, 1130
- API (application program interface), 5, 476, 483
- ARP (Address Resolution Protocol), 67, 77, 86, 97, 100-101, 106, 141, 343, 675-712
  - cache, 572, 675, 680, 682, 691, 703-704, 710-712
  - flooding, 109, 696
  - gratuitous, 178, 683, 707
  - header, 681
  - input queue, 97
  - multicasting, 710-711
  - proxy, 688, 703-704
  - RFC 1122 compliance, 1113
  - routing table, 675-678
  - structures, 681-683
  - timer functions, 694-696
- arp program, 571-572, 635, 641, 679-680, 688, 692, 694, 703-704, 706-707, 709-711

- arp\_allocates variable, 680, 707
- arp\_inuse variable, 680, 707, 710
- arp\_maxtries variable, 680, 699
- arp\_pro member, 686
- arp\_request function, 90, 169, 679, 695-696, 703-710
- ARPANET (Advanced Research Projects Agency network), 24, 901, 982, 1129
- arpcom structure, 77, 80, 86, 92, 99-101, 120, 159, 178, 343-344, 346, 364, 366, 685, 693, 709, 1071, 1080
- arphdr structure, 681, 687
- ARPHRD\_ETHERT constant, 686
- arpinit\_done variable, 680
- arpinput function, 106
- arpinr function, 107, 679, 687-688, 694
- arpinrq variable, 97, 101, 106, 680, 687
- arplookup function, 679, 691, 693, 697, 701-704, 707, 1080-1081
- ARPOP\_REPLY constant, 686, 694
- ARPOP\_REQUEST constant, 686
- ARPOP\_REVREPLY constant, 686
- ARPOP\_REVREQUEST constant, 686
- arpquest function, 679, 683-686, 688, 707
- arpresolve function, 109-110, 234, 378, 679, 683-684, 692, 696-701, 707, 710, 712, 1080-1081
- arpt\_down variable, 680
- arpt\_keep variable, 680, 692, 1113
- arpt\_prune variable, 680, 695, 1113
- arptfree function, 679, 695-696, 699, 704, 709
- arptimer function, 679, 694-695, 706, 711, 1081, 1113
- arptnew function, 702
- arpwhoas function, 679, 683-684, 699
- assassination, TIME\_WAIT, 964, 1089, 1126
- Asynchronous Transfer Mode, *see* ATM
- ATM (Asynchronous Transfer Mode), 1018, 1068
- atoi function, 8
- Auerbach, K., 1125
- 
- b\_to\_q function, 146
- backlog, connection, 463
- backoff, exponential, 836, 1119
- Banks, D., 994, 1126
- BBN (Bolt Beranek and Newman), 5
- Bcmp macro, 585, 596-597
- Bcopy macro, 585
- bd\_bif member, 1033, 1036-1037, 1040
- bd\_bufsize member, 1033, 1036-1037, 1043, 1046
- bd\_dcount member, 1033
- bd\_fbdf member, 1032-1033, 1036-1037
- bd\_filter member, 1033
- bd\_hbuf member, 1032-1033, 1036-1037
- bd\_hfiled member, 1032, 1036-1037, 1046
- bd\_immediate member, 1033
- bd\_next member, 1032, 1036-1037
- bd\_pad member, 1033
- bd\_promisc member, 1033
- bd\_rcount member, 1028, 1033
- bd\_rtout member, 1033
- bd\_sbuf member, 1032-1033, 1036-1037
- bd\_sel member, 1033
- bd\_slen member, 1032, 1036-1037
- bd\_state member, 1033
- Bellovin, S. M., 1081, 1126
- Berkeley fast filesystem, 27
- Berkeley Software Distribution, *see* BSD
- Berkeley-derived implementation, 4
- BGP (Border Gateway Protocol), 574
- bh\_caplen member, 1030
- bh\_dataalen member, 1030
- bh\_headerlen member, 1030
- bh\_tstamp member, 1030
- bibliography, 1125-1131
- bif\_dlist member, 1029-1030, 1036-1037
- bif\_dlt member, 1029-1030, 1035-1037, 1046
- bif\_driverp member, 1029-1030, 1036-1037
- bif\_hdrlen member, 1029-1031, 1036-1037
- bif\_if member, 1029
- bif\_ifp member, 1029-1030, 1036-1037
- bif\_next member, 1029-1030, 1036-1037
- bind, 719-721
  - explicit, 729
  - implicit, 729
- bind function, 8-9, 445-446, 452-454, 464, 554, 664, 666, 670, 719-721, 725, 729-730, 732-733, 740, 748, 750-751, 753, 786, 793, 930, 960, 1010, 1060, 1065, 1081-1084, 1091, 1114-1115, 1123
- bind\_args structure, 453
- BIOCFLOWB constant, 1035
- BIOCGBLN constant, 1033, 1035, 1043
- BIOCGDLT constant, 1035
- BIOCGRTIF constant, 1035
- BIOCGRTIMEOUT constant, 1035
- BIOCGSTATE constant, 1033, 1035
- BIOCIMMEDIATE constant, 1033, 1035
- BIOCPROMISC constant, 1033, 1035, 1092
- BIOCMBLEN constant, 1033, 1035
- BIOCSDF constant, 1034-1036
- BIOCSSETIF constant, 1033-1036, 1048, 1092
- BIOCSRTIMEOUT constant, 1035, 1043
- BIOCSVERSION constant, 1035

- Blindheim, R., xxii
- Boggs, D. R., 351, 1125
- Bolt Beranek and Newman, *see* BBN
- BOOTP (Bootstrap Protocol), 291, 321, 1110, 1126
- Bootstrap Protocol, *see* BOOTP
- Border Gateway Protocol, *see* BGP
- Borman, D. A., 235, 824, 1086, 1126, 1128
- Bostic, K., xxiii
- BPF (BSD Packet Filter), 68, 81, 83–85, 101–102, 104, 106–107, 112, 134, 137–138, 141, 152, 1027–1048, 1065, 1092, 1129
  - header, 134
  - loopback packet, 152
  - SLIP, 104
- BPF\_BUFSIZE constant, 1033
- BPF\_MAXBUFSIZE constant, 1033
- BPF\_MINBUFSIZE constant, 1033
- bpf\_allocate function, 1039
- bpf\_attachd function, 1039–1040
- bpf\_bufsize variable, 1028
- bpf\_d structure, 1028, 1032–1040
- bpf\_detachd function, 1039
- bpf\_dtab variable, 1028, 1032, 1034, 1036–1037, 1045
- bpf\_filter function, 1041
- bpf\_hdr structure, 1029–1030, 1032, 1043
- bpf\_if structure, 1028–1033, 1036–1037, 1040
- bpf\_iflist variable, 1028–1031, 1036–1037, 1039, 1091
- bpf\_ifname function, 1035
- bpf\_mcopy function, 1041
- bpf\_movein function, 1046
- bpf\_mtap function, 152, 1041
- bpf\_program structure, 1035
- bpf\_setif function, 1035–1036
- bpf\_setif function, 1035–1036, 1038–1039, 1091
- bpf\_stat structure, 1035
- bpf\_tap function, 104, 112, 138, 141, 152, 1037, 1040–1041, 1043
- bpf\_version structure, 1035
- bpf\_wakeup function, 1043, 1048, 1091
- bpfattach function, 81, 84–85, 1027, 1030–1032, 1040
- bpfioctl function, 1034–1035
- bpfopen function, 1034, 1036
- bpfread function, 1043–1046
- bpfwrite function, 1046–1047
- Braden, R. T., 205, 235, 252, 291, 301, 824, 866, 868, 870, 964, 1089, 1097, 1125–1126, 1128
- Brakmo, L. S., 845, 1126
- broadcast
  - packet, 99
  - storm, 326
- broadcast address
  - Ethernet, 97, 100, 103
  - IP, 182, 234, 1101
  - limited, 753
- BSD (Berkeley Software Distribution), 3, 68–69, 74, 76, 106, 140, 191, 219, 223, 397, 435, 441, 453, 1027
  - history, 3–5
  - Packet Filter, *see* BPF
- buffers
  - reliable protocol, 490
  - socket, 476–477
  - unreliable protocol, 490–491
- BUFOFFSET constant, 133–134
- bug, 33, 181, 223, 327–328, 442, 473, 548, 667, 692, 774, 879, 948, 956, 959, 1005, 1010, 1076, 1083, 1088, 1090, 1106, 1126, 1128
- Bzero macro, 585
- C function prototypes, ANSI, 41
- cache
  - ARP, 572, 675, 680, 682, 691, 703–704, 710–712
  - hiding, UDP, 791
  - multicast group, 399, 412, 415, 434
  - multicast one-behind, 398–399, 422, 424, 434
  - TCP one-behind, 231, 798, 897, 929, 941
  - UDP one-behind, 231, 757, 773–774, 786, 791, 794
  - unicast one-behind, 223, 253
- cached route, 746–747, 750, 768, 843, 887, 894, 898, 1111
- cached segments, 972
- cached\_mrt variable, 398
- cached\_origin variable, 398, 423
- cached\_originmask variable, 398
- caddr\_t data type, 52
- Calamvokis, C., 994, 1126
- callout function, 94
- Carlson, J., 959, 1126
- carrier sense multiple access, *see* CSMA
- Casner, S., 350, 1126
- catchpacket function, 1041–1043, 1048, 1091
- Chapin, A. L., 9, 1130
- checksum
  - algorithm, 1126
  - ICMP, 309
  - IP, 234–239
  - TCP, 800
  - UDP, 758, 764–768, 792
- Cheriton, D. P., 401, 419, 1127
- Cheswick, W. R., 1081, 1126
- child interface, 418–419, 429

- Clark, D. D., 1067–1068, 1126
- Clark, J. J., xxii
- class
  - A IP address, 155
  - B IP address, 155
  - C IP address, 155
  - D IP address, 155, 341
  - E IP address, 155
- cli&#228;f, high-water mark, 147
- cli&#228;f, structure, 131, 141
- cli&#228;intrq variable, 97, 100–101, 150
- CLNP (Connectionless Network Protocol), 97, 100, 666
  - input queue, 97
- GLOBAL constant, 135
- close function, 10, 13, 440, 442, 445–447, 468, 471, 514, 542, 552, 555, 666–667, 669, 786, 818, 820, 825, 827, 829, 980, 1010, 1012, 1019, 1025, 1059, 1087, 1119
- cluster
  - mbuf, 16, 33
  - reference counts, 56–60
  - SLIP, 131
- cm&#228;g\_data member, 781
- cm&#228;g\_len member, 482–483, 781
- cm&#228;g\_level member, 483
- cm&#228;g\_type member, 482–483
- cm&#228;g\_hdr structure, 482–483, 781
- code, ICMP, 302–303
- collision, 97–98, 143
  - with select, 531–532, 534
- Comer, D. E., 457, 1018, 1126
- compressed SLIP, see CSLIP
- compression, header, 995–1004
- Computer Systems Research Group, see CSRG
- CONCAT macro, 193
- concurrency, interrupt levels and, 23–26
- congestion
  - avoidance, 844, 939, 970, 972, 977, 1119, 1127
  - window, 835, 844, 852, 854–855, 882, 903, 906, 939, 972–975, 977
- connect, breaking association with, 468
- connect function, 8, 439–440, 445–446, 464–465, 481, 494, 552, 664, 666, 721, 725, 727, 729–730, 735, 740–741, 748, 750–751, 763, 787–788, 793, 805, 813, 828, 882, 903, 918–920, 930, 946, 948, 960, 962–963, 1011–1013, 1025, 1059–1060, 1065, 1083, 1091, 1115, 1123
- connect structure, 453, 464–465, 467–468, 481
- connect\_args structure, 465
- connected UDP socket, 721, 755, 779–780
- connection, old incarnation, 814
  - connection-establishment timer, 817, 819, 828–831, 892, 946, 948, 1012, 1121
- Connectionless Network Protocol, see CLNP
- connectionless, transport, OSI, 106
- control block
  - routing, 647
  - TCP, 713, 718, 800, 803–805, 808, 818–819, 821–822, 832–835, 837, 846, 866–867, 871–872, 884, 887–888, 893, 897, 906–907, 909, 916, 930, 932, 944–946, 949, 959–960, 966, 969, 981, 986, 989, 1009–1010, 1018–1019, 1021, 1023, 1084
- control message, 509–510
  - reference counts, 470
- conventions
  - source code, 1–3
  - typographical, 3
- Coordinated Universal Time, see UTC
- copyin function, 453, 483–484, 525, 1079
- copyout function, 460, 642–643, 1079
- copyright, source code, xxi–xxii
- cpu\_startup function, 79–83
- CRC (cyclic redundancy check), 99, 103
- Croft, W., 321, 1126
- Crowcroft, J., xxii, 750, 1126
- CSLIP (compressed SLIP), 147, 995–1004, 1127
- CSMA (carrier sense multiple access), 97
- CSRG (Computer Systems Research Group), xix, xxii, 1067
- CTL\_FW constant, 201
- CTL\_RR&#228;N constant, 201
- CTL\_N&#228;T constant, 201, 637
- CTL\_USER constant, 201
- cyclic redundancy check, see CRC
- Dalton, C., 994, 1126
- data-link frame, 96, 127, 210, 278, 1027, 1046–1047
- Dawley, K. B., xxii
- DECNET, 100
- Deering, S. E., 301, 338, 381, 401, 419, 901, 1067, 1127, 1129, 1131
- default
  - raw protocol, 191
  - route, 181
  - TTL, 207
- deferred carries, 236
- del\_intrp function, 401, 412–414, 1077
- del\_m&#228;t function, 401, 421
- del\_vif function, 401, 409–410
- delayed ACK timer, 817–818, 821, 861, 864
- Delp, G., 750



- demultiplexing
  - IP, 219
  - TCP segments, 721–723, 728
  - UDP datagrams, 723–724
- descriptor, 10–15
  - socket, 6, 445–447
- destination
  - address, Ethernet, 99, 110
  - address, IP, 162, 182
  - group, 405
  - unreachable, ICMP, 279
- /dev/bpf device, 1027
- /dev/kmem device, 916
- /dev/mem device, 37
- /dev/udp device, 8
- device driver, 63–94
  - BPF, 1027–1028
  - Ethernet, 63–64, 77, 81, 92, 95–96, 98–112, 124–125, 161, 1029, 1037, 1040
  - loopback, 64, 128, 150–153, 161, 1029
  - SLIP, 63–64, 69, 83, 128–150, 161, 179, 1029
  - TTY, 129–130, 134–135, 141, 148–149
- DF (don't fragment flag, IP header), 230, 275–276, 278–279, 283, 302, 1084, 1122
- DHCP (Dynamic Host Configuration Protocol), 321, 1110, 1127
- diameter, Internet, 223, 1129
- direct route, 561, 621, 706
- directed broadcast address, 157, 162
- Distance Vector Multicast Routing Protocol, *see* DVMRP
- DLT\_EN10MB constant, 1029–1030
- DLT\_NULL constant, 1029–1030
- DLT\_SLIP constant, 1029–1030
- DNS (Domain Name System), 140, 291, 1125
- Dogfight, SGL, 338
- dom\_attach member, 586
- dom\_dispose member, 187, 470, 646, 1078
- dom\_externalize member, 187, 517–518, 646
- dom\_family member, 187, 581, 646, 671
- dom\_init member, 187, 194, 581, 584, 646
- dom\_maxrtkey member, 187, 581, 585, 646
- dom\_name member, 187, 646
- dom\_next member, 187, 193, 646
- dom\_protosw member, 187, 646
- dom\_protoswnPROTOSW member, 187, 646
- dom\_rattach member, 187, 581, 584, 646
- dom\_rtoffset member, 187, 581, 587, 646
- domain, 185, 323, 445, 447, 449
  - initialization, IP, 199
  - initialization, routing, 646
  - Internet, 75, 160, 193, 309, 316, 385, 460, 483
  - OSI, 75
  - routing, 67, 75, 437, 539, 554, 569–570, 572, 581, 584, 624, 632, 645–673
  - Unix, 75, 189, 450, 460, 470, 510, 518, 1077
- Domain Name System, *see* DNS
- domain structure, 186–188, 191–195, 202, 204, 581, 584–585, 587, 646
  - Internet, 186, 191–196
  - routing, 646
- domaininit function, 79, 193–196, 199, 204, 571, 584, 646, 756, 760, 796, 812, 1050, 1053
- domains variable, 186–187, 193, 195–196, 204
- don't fragment flag, IP header, *see* DF
- dotted-decimal, 7, 156
- Dove, K. F., 750, 791, 994, 1129
- Droms, R., 321, 1127
- dtom macro, 44, 46–48, 50–52, 61, 909, 913
- DTYPE\_SOCKET constant, 13–14, 713
- DTYPE\_VNODE constant, 13
- dup function, 10–11
- duplicate keys, 587–591
- duplicate, wandering, 813
- DVMRP (Distance Vector Multicast Routing Protocol), 337–339, 384, 401, 418–419, 1131
  - DVMRP\_ADD\_LGRP socket option, 401, 411–413
  - DVMRP\_ADD\_MRT socket option, 401, 419, 421–422
  - DVMRP\_ADD\_VIP socket option, 401, 407, 409
  - DVMRP\_DEL\_LGRP socket option, 401, 411–412, 414
  - DVMRP\_DEL\_MRT socket option, 401, 421
  - DVMRP\_DEL\_VIP socket option, 401, 407, 409–410
  - DVMRP\_DONB socket option, 401, 433, 1059
  - DVMRP\_INIT socket option, 401, 403
- Dynamic Host Configuration Protocol, *see* DHCP
- EADDRS error, 234, 403, 453
- EADDRINUSE error, 359, 403, 409, 720, 733–734, 740
- EADDRNOTAVAIL error, 173, 175, 234, 354, 358, 366, 409, 411, 414, 468, 1059
- EAFNOSUPPORT error, 110, 179, 361, 363, 468
- EALREADY error, 465, 467
- HEADF error, 528, 534
- EBUSY error, 1027, 1034
- echo option, 1066
- echo reply, ICMP, 317
- echo request, ICMP, 317
- ECONNABORTED error, 892
- ECONNREFUSED error, 303, 748, 843, 892, 963
- ECONNRESET error, 892, 964–965, 1119
- EDSTADREQ error, 494
- EDOM error, 544



- Edwards, A., 994, 1126
- EXIST error, 611, 650
- EGP (Exterior Gateway Protocol), 65, 1050
- EHOSTDOWN error, 228, 699–700, 712, 905, 1081
- EHOSTUNREACH error, 107–108, 152, 228, 232, 303, 621, 699, 828, 843, 892, 905, 1081
- ENINPROGRESS error, 465
- ENINTR error, 457, 459, 478, 528, 1045
- EINVAL error, 124–125, 172, 178–180, 240, 271, 353–355, 358, 366, 401, 409, 411, 414, 453, 459, 539, 543, 551, 653, 785, 1007, 1016, 1036, 1045, 1078
- EISCONN error, 467–468, 763, 787, 1009
- EMSGSIZE error, 228, 234, 279, 303, 484, 495, 1046
- encapsulation
  - 802.3, 106, 125
  - Ethernet, 99, 106
  - SLIP, 128–129
- END character, SLIP, 129, 144
- end of option list, *see* EOL
- ENETDOWN error, 107, 228, 1039
- ENETRESET error, 140, 362, 364, 369, 1077
- ENETUNREACH error, 152, 232, 376, 905
- err\_no member, 343, 366
- err\_addrhi member, 342, 344
- err\_addrlo member, 342, 344
- err\_next member, 343, 366
- err\_rfcount member, 343, 364, 366
- ENOBUFS error, 111, 226, 228, 279, 351, 453, 479, 789, 882, 889, 892, 906, 1107
- ENOPROTOPT error, 202, 241, 303, 334, 348, 539–540, 546
- ENOTCONN error, 494, 515, 556, 763, 788
- ENOTDIR error, 245, 334
- ENOTSUP error, 554
- ENXIO error, 116, 132, 369, 1034, 1039, 1046, 1091
- EOL (end of option list), 249–250, 282–283, 865, 933–934
- EOPNOTSUPP error, 164, 245, 351, 371, 403, 468, 513, 513
- EPERM error, 166
- ephemeral port, 21, 715, 719, 725, 729–730, 732, 740, 748, 751–753, 760, 813, 945, 1011–1012, 1081–1083
- EPIPE error, 489, 494
- Epoch, Unix, 105, 683, 695
- ERESTART error, 456–457, 478, 528, 1045
- err\_sys function, 6
- errno variable, 6, 650, 744, 748, 783, 892, 921
- error
  - EACCESS, 234, 403, 453
  - EADDRINUSE, 359, 403, 409, 720, 733–734, 740
  - EADDRNOTAVAIL, 173, 175, 234, 354, 358, 366, 409, 411, 414, 468, 1059
  - EAFNOSUPPORT, 110, 179, 361, 363, 468
  - EALREADY, 465, 467
  - EBADF, 528, 534
  - EBUSY, 1027, 1034
  - ECANCELED, 892
  - ECONNREFUSED, 303, 748, 843, 892, 963
  - ECONNRESET, 892, 964–965, 1119
  - EBUSY, 494
  - EDOM, 544
  - EXIST, 611, 650
  - EHOSTDOWN, 228, 699–700, 712, 905, 1081
  - EHOSTUNREACH, 107–108, 152, 228, 232, 303, 621, 699, 828, 843, 892, 905, 1081
  - ENINPROGRESS, 465
  - ENINTR, 457, 459, 478, 528, 1045
  - EINVAL, 124–125, 172, 178–180, 240, 271, 353–355, 358, 366, 401, 409, 411, 414, 453, 459, 539, 543, 551, 653, 785, 1007, 1016, 1036, 1045, 1078
  - EISCONN, 467–468, 763, 787, 1009
  - EMSGSIZE, 228, 234, 279, 303, 484, 495, 1046
  - ENETDOWN, 107, 228, 1039
  - ENETRESET, 140, 362, 364, 369, 1077
  - ENETUNREACH, 152, 232, 376, 905
  - ENOBUFS, 111, 226, 228, 279, 351, 453, 479, 789, 882, 889, 892, 906, 1107
  - ENOPROTOPT, 202, 241, 303, 334, 348, 539–540, 546
  - ENOTCONN, 494, 515, 556, 763, 788
  - ENOTDIR, 245, 334
  - ENOTSUP, 554
  - ENXIO, 116, 132, 369, 1034, 1039, 1046, 1091
  - EOPNOTSUPP, 164, 245, 351, 371, 403, 468, 513
  - EPERM, 166
  - EPIPE, 489, 494
  - ERESTART, 456–457, 478, 528, 1045
  - ERCHR, 421, 654
  - ETIMEDOUT, 828, 843, 892
  - ETOOMANYREFS, 359, 655
  - EMULDBLOCK, 439, 457, 459, 478, 491, 496, 515–516, 528, 792, 1016
  - ICMP, 205, 223, 226, 228, 232, 250, 292, 326
- ESC character, SLIP, 129, 144
- ERCHR error, 421, 654
- Estrin, D., 419, 1127
- /etc/motd.conf file, 407
- /etc/netstat file, 84, 162, 560, 706, 709, 1110
- ETHER\_LOOKUP\_MULTI macro, 343–344, 364, 369
- ETHER\_MAP\_IP\_MULTICAST macro, 342, 363, 378, 697, 710–711

- `@ether_addrmulticast` function, 356, 362–364, 369
- `@ether_arp` structure, 309, 681, 685–686, 689
- `@ether_delmulticast` function, 356, 362, 369–371, 1077
- `@hw_driver` member, 103–104
- `ether_header` structure, 101–104, 111, 681, 1032, 1040
- `ether_ifattach` function, 91–92
- `ether_input` function, 100–101, 104–107, 125, 213, 221, 373, 687, 943, 1070, 1098
- `ether_ipmulticast_max` variable, 340, 363
- `ether_ipmulticast_min` variable, 340, 363
- `ether_multi` structure, 81, 342–344, 346, 356, 363–366, 369, 379
  - reference counts, 343, 346, 364, 369
- `ether_output` function, 96, 101, 107–112, 127, 139, 141, 150, 152, 378, 580, 679, 684–686, 692, 696, 699–701, 789, 1047, 1071, 1081, 1084, 1086, 1091
- `ether_type` member, 101
- `etherbroadcastaddr` variable, 97, 103, 363
- ETHERMTU constant, 81, 92
- Ethernet
  - broadcast address, 97, 100, 103
  - destination address, 99, 110
  - device driver, 63–64, 77, 81, 92, 95–96, 98–112, 124–125, 161, 1029, 1037, 1040
  - encapsulation, 99, 106
  - frame, 66, 92, 95, 99–104, 106–112, 125, 212, 216, 363, 1040, 1046
  - hardware address, 75, 81, 87, 91–92, 104, 341, 1046
  - header, 92, 103, 110
  - initialization, 80–81
  - length, 106
  - MTL, 92
  - multicast address, 100, 103–104, 341–342
  - multicasting, 156
  - `sockaddr_dl` structure, 91
  - source address, 99
  - type, 99, 103, 106–107, 110, 141, 686
  - unicast address, 100
- ETHERTYPE\_ARP constant, 686–687
- ETHERTYPE\_IP constant, 686–687, 694
- ETHERTYPE\_IPTRAILERS constant, 686–687, 694
- ETHERTYPE\_REVARP constant, 686
- ETIMEDOUT error, 828, 843, 892
- ETOOMANYREPS error, 359, 655
- EWOULDLOCK error, 439, 457, 459, 478, 491, 496, 515–516, 528, 792, 1016
- `exec` function, 27, 555
- exercises, solutions to, 1069–1092
- `exit` function, 10
- expanding-ring search, 351
- experimental address, IP, 156
- exponential backoff, 836, 1119
- `ext_buf` member, 34
- `ext_free` member, 34
- `ext_size` member, 34
- Exterior Gateway Protocol, *see* EGP
- external buffer, `mbuf`, 33
- F\_GETFL macro, 552
- F\_SETOWN constant, 550, 552
- F\_SETFL macro, 552
- F\_SETOWN constant, 550, 552
- `f_data` member, 13, 446–447, 471
- `f_ops` member, 13, 446–447, 471
- `f_type` member, 13
- `falloc` function, 447, 460
- FAQ (frequently asked question), 1094
- Farinacci, D., 399, 419, 1127, 1129
- fast
  - filesystem, Berkeley, 27
  - recovery, 970–974, 1120, 1128
  - retransmit, 908, 970–974, 1120, 1128
- FASYNC constant, 549–550, 552
- `fcntl` function, 10, 439, 445, 507, 537–538, 548–551, 557
- FD\_SETSIZE macro, 525
- `fd_ofileflag` member, 13
- `fd_ofiles` member, 13–14, 446–447
- `fd_set` data type, 525
- FDDI (Fiber Distributed Data Interface), 69, 337, 1018, 1068
  - `ffa` function, 528
- Fiber Distributed Data Interface, *see* FDDI
- `file` structure, 13–14, 446–447, 453, 455, 459–460, 471, 487, 503, 528, 539, 554, 713
- File Transfer Protocol, *see* FTP
- `filedesc` structure, 13, 446–447, 451
- `fileops` structure, 13, 437, 446–448, 529
- filesystem, Berkeley fast, 27
- FIN (finish flag, TCP header), 468, 470, 803, 805
- FIN\_WAIT\_2 timer, 818–819, 821–822, 825–827, 980, 991, 1085
- finish flag, TCP header, *see* FIN
- Finlayson, R., 100, 321, 1127
- FIOASYNC constant, 550, 552
- FIONBIO constant, 550, 552
- FIONREAD constant, 550, 552, 1035
- firewall gateway, 1126
- flooding, ARP, 109, 696
- Floyd, S., xxii, 970, 977, 1127

- PNONBLOCK constant, 549–550, 552
- fo\_close member, 448, 471
- fo\_ioctl member, 448, 552
- fo\_read member, 448
- fo\_select member, 448, 528–529
- fo\_write member, 448, 761
- Forgie, J., 215, 1127
- fork function, 10–11, 27, 555, 1091
- forwarding
  - IP, 181, 217–228
  - multicast, 424–433
- fragmentation, 1128
  - broadcast packet, 234
  - ICMP, 314
  - IP, 48–50, 210, 232, 275–283
  - offset, 276
- frame
  - data-link, 96, 127, 210, 278, 1027, 1046–1047
  - Ethernet, 66, 92, 95, 99–104, 106–112, 125, 212, 216, 363, 1040, 1046
  - PPP, 129
  - SLIP, 83, 128, 131, 133–137, 143–144, 146–147, 1070
- FRAME\_END constant, 136
- FREAD constant, 469, 528
- free function, 88, 671, 897
- Free macro, 585, 606
- frequently asked question, *see* FAQ
- fstat function, 670, 789, 1015, 1062
- fstat program, 187
- FTP (File Transfer Protocol), 4, 140, 272, 921, 1002, 1084, 1087, 1093–1094, 1125
- fudge factor, 463, 640
- full-duplex, 808
- Fuller, V., 170, 1127
- function prototypes, ANSI C, 41
- FWRITE constant, 469, 528
  
- garbage byte, 887
- gated program, 559–560, 571–572, 601, 637, 644, 1050, 1111
- GATEWAY constant, 2, 158
- gateway, firewall, 1126
- Gateway Requirements RFC, 1126
- gather, 481–482, 486–487
- getbits macro, 525
- getmsg function, 8
- getpeername function, 443, 445, 501, 514, 537–538, 555–557, 670, 741, 789, 1017, 1062
- getsock function, 451–453, 455, 459, 465, 469, 487, 503, 539, 545, 554, 556
- getsockname function, 443, 445, 537–538, 554, 556–557, 670, 741, 789, 1017, 1062, 1115
- getsockopt function, 239–244, 401, 437, 445, 537–539, 545–548, 557, 1022, 1024, 1063, 1115, 1123
- Gilmore, J., 321, 1126
- Glaser, G., xxii
- GNU software, 1094
- Grandi, S., xxii
- gratuitous ARP, 178, 683, 707
- grep program, 2–3, 1133
- groty, 654
- group, destination, 405
  - SNMP, 65
- grplist\_member function, 412, 415, 429, 1077
- Gulbenkian, J., xxii
- Gurwitz, R., 5
  
- hack, 913, 1017
- half-close, 468, 807, 818, 825, 957, 980, 1119
- hardware address, 68, 85–87, 89, 100
  - Ethernet, 75, 81, 87, 91–92, 104, 341, 1046
- hardware interrupt, 205
- hardware multicasting, 156, 337
- Haverty, J., 5
- HDLC (high-level data link control), 27
- header
  - ARP, 681
  - BPF, 134
  - compression, 995–1004
  - Ethernet, 92, 103, 110
  - ICMP, 309
  - IGMP, 385
  - IP, 210–212
  - prediction, 923, 934–941
  - TCP, 801–803
  - UDP, 759–760
- held route, 606, 659
- HELLO (routing protocol), 1050
- Hering, D., xxii
- hiding, UDP cache, 791
- high-level data link control, *see* HDLC
- high-water mark, 477, 479, 490, 495–496, 508, 534, 539, 543
  - clist, 147
- Hinden, R., 190, 716, 1130
- history, BSD, 3–5
- Hoffman, E., 1128
- Hogue, J. E., xxii
- Hornig, C., 100, 1127
- Host Requirements RFC, 1125

- ether\_addmulticast function, 356, 362–364, 369
- ether\_arp structure, 309, 681, 685–686, 689
- ether\_delmulticast function, 356, 362, 369–371, 1077
- ether\_ghost member, 103–104
- ether\_header structure, 101–104, 111, 681, 1032, 1040
- ether\_ifattach function, 91–92
- ether\_input function, 100–101, 104–107, 125, 213, 221, 373, 687, 943, 1070, 1098
- ether\_output function, 96, 101, 107–112, 127, 139, 141, 150, 152, 378, 580, 679, 684–686, 692, 696, 699–701, 789, 1047, 1071, 1081, 1084, 1086, 1091
- ether\_type member, 101
- etherbroadcastaddr variable, 97, 103, 363
- ETHERMTU constant, 81, 92
- Ethernet
  - broadcast address, 97, 100, 103
  - destination address, 99, 110
  - device driver, 63–64, 77, 81, 92, 95–96, 98–112, 124–125, 161, 1029, 1037, 1040
  - encapsulation, 99, 106
  - frame, 66, 92, 95, 99–104, 106–112, 125, 212, 216, 363, 1040, 1046
  - hardware address, 75, 81, 87, 91–92, 104, 341, 1046
  - header, 92, 103, 110
  - initialization, 80–81
  - length, 106
  - MTU, 92
  - multicast address, 100, 103–104, 341–342
  - multicasting, 156
  - socaddr\_dl structure, 91
  - source address, 99
  - type, 99, 103, 106–107, 110, 141, 686
  - unicast address, 100
- ETHERTYPE\_arp constant, 686–687
- ETHERTYPE\_ip constant, 686–687, 694
- ETHERTYPE\_iptrailer constant, 686–687, 694
- ETHERTYPE\_revarf constant, 686
- ETIMEDOUT error, 828, 843, 892
- ETOOMANYREFS error, 359, 655
- EMOULDRELOC error, 439, 457, 459, 478, 491, 496, 515–516, 528, 792, 1016
- exec function, 27, 555
- exercises, solutions to, 1069–1092
- exit function, 10
- expanding-ring search, 351
- experimental address, IP, 156
- exponential backoff, 836, 1119
- ext\_buf member, 34
- ext\_free member, 34
- ext\_size member, 34
- Exterior Gateway Protocol, *see* EGP
- external buffer, mbuf, 33
- F\_GETFD macro, 552
- F\_GETOWN constant, 550, 552
- F\_SETFD macro, 552
- F\_SETOWN constant, 550, 552
- f\_data member, 13, 446–447, 471
- f\_ops member, 13, 446–447, 471
- f\_type member, 13
- falloc function, 447, 460
- FAQ (frequently asked question), 1094
- Farinacci, D., 399, 419, 1127, 1129
- fast
  - filesystem, Berkeley, 27
  - recovery, 970–974, 1120, 1128
  - retransmit, 908, 970–974, 1120, 1128
- FASYNC constant, 549–550, 552
- fcntl function, 10, 439, 445, 507, 537–538, 548–551, 557
- FD\_SETSIZE macro, 525
- fd\_ofileflags member, 13
- fd\_ofiles member, 13–14, 446–447
- fd\_set data type, 525
- FDDI (Fiber Distributed Data Interface), 69, 337, 1018, 1068
- ffs function, 528
- Fiber Distributed Data Interface, *see* FDDI
- file structure, 13–14, 446–447, 453, 455, 459–460, 471, 487, 503, 528, 539, 554, 713
- File Transfer Protocol, *see* FTP
- filedesc structure, 13, 446–447, 451
- fileops structure, 13, 437, 446–448, 529
- filesystem, Berkeley fast, 27
- FIN (finish flag, TCP header), 468, 470, 803, 805
- FIN\_WAIT\_2 timer, 818–819, 821–822, 825–827, 980, 991, 1085
- finish flag, TCP header, *see* FIN
- Finlayson, R., 100, 321, 1127
- FIOASYNC constant, 550, 552
- FIONBIO constant, 550, 552
- FIONREAD constant, 550, 552, 1035
- firewall gateway, 1126
- flooding, ARP, 109, 696
- Floyd, S., xxii, 970, 977, 1127

- FNONBLOCK constant, 549–550, 552
- fo\_close member, 448, 471
- fo\_ioctl member, 448, 552
- fo\_read member, 448
- fo\_select member, 448, 528–529
- fo\_write member, 448, 761
- Forgie, J., 215, 1127
- fork function, 10–11, 27, 555, 1091
- forwarding
  - IP, 181, 217–228
  - multicast, 424–433
- fragmentation, 1128
  - broadcast packet, 234
  - ICMP, 314
  - IP, 48–50, 210, 232, 275–283
  - offset, 276
- frame
  - data-link, 96, 127, 210, 278, 1027, 1046–1047
  - Ethernet, 66, 92, 95, 99–104, 106–112, 125, 212, 216, 363, 1040, 1046
  - PPP, 129
  - SLIP, 83, 128, 131, 133–137, 143–144, 146–147, 1070
- FRAME\_END constant, 136
- FREAD constant, 469, 528
- free function, 88, 671, 897
- Free macro, 585, 606
- frequently asked question, see FAQ
- fatat function, 670, 789, 1015, 1062
- fatat program, 187
- FTP (File Transfer Protocol), 4, 140, 272, 921, 1002, 1084, 1087, 1093–1094, 1125
- fudge factor, 463, 640
- full-duplex, 806
- Fuller, V., 170, 1127
- function prototypes, ANSI C, 41
- FWRITE constant, 469, 528
  
- garbage byte, 887
- gated program, 559–560, 571–572, 601, 637, 644, 1050, 1111
- GATEWAY constant, 2, 158
- gateway, firewall, 1126
- Gateway Requirements RFC, 1126
- gather, 481–482, 486–487
- getbits macro, 525
- getmsg function, 8
- getpeername function, 443, 445, 501, 514, 537–538, 555–557, 670, 741, 789, 1017, 1062
- getsock function, 451–453, 455, 459, 465, 469, 487, 503, 539, 545, 554, 556
- getsockname function, 443, 445, 537–538, 554, 556–557, 670, 741, 789, 1017, 1062, 1115
- getsockopt function, 239–244, 401, 437, 445, 537–539, 545–548, 557, 1022, 1024, 1063, 1115, 1123
- Gilmore, J., 321, 1126
- Glater, G., xxii
- GNU software, 1094
- Grandj, S., xxii
- gratuitous ARP, 178, 683, 707
- grep program, 2–3, 1133
- grotty, 654
- group, destination, 405
  - SNMP, 65
- grplist\_member function, 412, 415, 429, 1077
- Gulbenkian, J., xxii
- Gurwitz, R., 5
  
- hack, 913, 1017
- half-close, 468, 807, 818, 825, 957, 980, 1119
- hardware address, 68, 85–87, 89, 100
  - Ethernet, 75, 81, 87, 91–92, 104, 341, 1046
- hardware interrupt, 205
- hardware multicasting, 156, 337
- Haverty, J., 5
- HDLCL (high-level data link control), 27
- header
  - ARP, 681
  - BPF, 134
  - compression, 995–1004
  - Ethernet, 92, 103, 110
  - ICMP, 309
  - IGMP, 385
  - IP, 210–212
  - prediction, 923, 934–941
  - TCP, 801–803
  - UDP, 759–760
- held route, 606, 659
- HELLO (routing protocol), 1050
- Hering, D., xxii
- hiding, UDP cache, 791
- high-level data link control, see HDLC
- high-water mark, 477, 479, 490, 495–496, 508, 534, 539, 543
  - clist, 147
- Hinden, R., 190, 716, 1130
- history, BSD, 3–5
- Hoffman, E., 1128
- Hogue, J. E., xxii
- Hornig, C., 100, 1127
- Host Requirements RFC, 1125



- host unreachable, ICMP, 208, 223, 253
- host, vs. router, 157
- hostname macro, 525
- hp\_device structure, 81
- htons function, 7
- Hutchinson, N. C., 60, 750, 1127
- is variable, 64, 94, 544, 548, 1010
- isrc function, 525
- is\_in member, 387
  - is\_inn member, 387
  - is\_in\_addr member, 162, 168, 219
  - is\_in\_broadcast member, 162, 166, 219, 320
  - is\_in\_netaddr member, 162, 168, 320
  - is\_in\_flags member, 161, 171
  - is\_in\_ifa member, 161
  - is\_in\_ifp member, 161
  - is\_multiaddrs member, 161, 346, 361, 366
  - is\_net member, 162, 219
  - is\_netbroadcast member, 162, 219
  - is\_netmask member, 162, 170
  - is\_next member, 159, 161, 346
  - is\_socket member, 162, 166, 168, 170, 172, 320
  - is\_subnet member, 162, 219
  - is\_subnetmask member, 162, 170, 172
- IANA (Internet Assigned Numbers Authority), 338, 341, 363
- ICMP (Internet Control Message Protocol), 65, 83, 140, 188–191, 193, 198, 203, 205–206, 228, 240, 259, 301–335, 381, 440, 451, 464, 477, 541, 1130
  - address mask reply, 319
  - address mask request, 319
  - checksum, 309
  - code, 302–303
  - destination unreachable, 279
  - echo reply, 317
  - echo request, 317
  - error, 205, 223, 226, 228, 232, 250, 292, 326
  - errors and UDP sockets, 748–749
  - fragmentation, 314
  - header, 309
  - host unreachable, 208, 223, 253
  - inetaw variable, 203, 309–310, 1098
  - input processing, 310–323
  - output processing, 324–333
  - parameter problem, 252, 257, 262, 314
  - port unreachable, 494
  - protocol structure, 309–310
  - redirect, 221, 223–228, 252, 321
  - redirect and raw sockets, 746–748
  - reply, 234
  - RFC 1122 compliance, 1105–1110
  - router discovery, 1127
  - router solicitation, 339
  - source quench, 226, 314
  - source route failure, 257
  - time exceeded, 223, 292–293, 300, 314
  - timestamp reply, 318
  - timestamp request, 318
  - type, 250, 302–303
  - unreachable, 314
- icmp structure, 308–309
  - ICMP\_ADVLEN macro, 312, 322
  - ICMP\_ADVLENMIN constant, 312, 322
  - ICMP\_BCMO constant, 302–303, 306, 316–317
  - ICMP\_ECHOREPLY constant, 302–303, 306, 316–317
  - ICMP\_INFOTYPE constant, 325
  - ICMP\_IREQ constant, 302–303, 321
  - ICMP\_IREQREPLY constant, 302–303, 321
  - ICMP\_MASKLEN constant, 312
  - ICMP\_MASKREPLY constant, 302–303, 319–321
  - ICMP\_MASKREQ constant, 302–303, 319–321
  - ICMP\_MINLEN constant, 312
  - ICMP\_PARAMPROB constant, 250, 302–303
  - ICMP\_PARAMPROB\_OPTABSENT constant, 302–303
  - ICMP\_REDIRECT constant, 302–303, 306, 321–322
  - ICMP\_REDIRECT\_HOST constant, 302–303
  - ICMP\_REDIRECT\_NET constant, 302–303
  - ICMP\_REDIRECT\_TOSHOST constant, 302–303
  - ICMP\_REDIRECT\_TOSNET constant, 302–303
  - ICMP\_ROUTERADVERT constant, 302–303, 321
  - ICMP\_ROUTERSOLICIT constant, 302–303, 321
  - ICMP\_SOURCEQUENCH constant, 228, 302–303, 306
  - ICMP\_TIMXCEED constant, 302–303, 306
  - ICMP\_TIMXCEED\_INTRANS constant, 302–303
  - ICMP\_TIMXCEED\_REASS constant, 302–303
  - ICMP\_TTLLEN constant, 312
  - ICMP\_TSTAMP constant, 302–303, 318–319
  - ICMP\_TSTAMPREPLY constant, 302–303, 318–319
  - ICMP\_UNREACH constant, 302–303, 306
  - ICMP\_UNREACH\_HOST constant, 228, 302–303
  - ICMP\_UNREACH\_HOST\_PROHIB constant, 302–303
  - ICMP\_UNREACH\_HOST\_UNKNOWN constant, 302–303
  - ICMP\_UNREACH\_ISOLATED constant, 302–303
  - ICMP\_UNREACH\_NEEDFRAG constant, 228, 302–303
  - ICMP\_UNREACH\_NET constant, 302–303

- ICMP\_UNREACH\_NET\_PROHIB constant, 302-303
- ICMP\_UNREACH\_NET\_UNKNOWN constant, 302-303
- ICMP\_UNREACH\_PORT constant, 302-303
- ICMP\_UNREACH\_PROTOCOL constant, 302-303
- ICMP\_UNREACH\_SRCFAIL constant, 302-303
- ICMP\_UNREACH\_TOSHOST constant, 302-303
- ICMP\_UNREACH\_TOSNET constant, 302-303
- icmp\_checksum member, 309, 314, 317-319, 321, 333
- icmp\_code member, 309, 313, 315, 317-319, 321, 327
- icmp\_data member, 317
- icmp\_dun member, 309
- icmp\_error function, 205, 226, 246, 252, 306, 324-329, 332, 335, 774, 1075, 1102, 1104, 1106
- icmp\_gwaddr member, 321, 327
- icmp\_hun member, 309
- icmp\_id member, 317-319
- icmp\_input function, 304, 309-323, 329, 333, 335, 571, 617, 619, 743, 756, 782-783, 796, 1050, 1052-1053, 1056, 1102, 1105, 1107-1110
  - error processing, 313-315
  - redirect processing, 321-323
  - reply processing, 323
  - request processing, 316-321
- icmp\_ip member, 314, 321
- icmp\_leavegroup function, 395
- icmp\_len member, 314
- icmp\_mask member, 319-320
- icmp\_nextmtu member, 314, 327
- icmp\_otime member, 318
- icmp\_output function, 324
- icmp\_pnvoid member, 314
- icmp\_pptr member, 314, 327
- icmp\_reflect function, 45, 303-304, 317, 324, 326, 328-333, 1107-1109, 1122
- icmp\_rtime member, 318, 335
- icmp\_send function, 324, 332-333, 1075
- icmp\_seq member, 317-319
- icmp\_syncctl function, 203, 319, 334, 1052
- icmp\_ttime member, 318-319
- icmp\_type member, 309, 313-315, 317-321, 325, 327
- icmp\_void member, 314
- ICMPCTL\_MASKREPL constant, 202, 334
- icmpdat variable, 319, 322
- icmpgw variable, 322
- icmpInAddrMaskReps variable, 307
- icmpInAddrMasks variable, 307
- icmpInDestUnrechs variable, 307
- icmpInEchoReps variable, 307
- icmpInEchoes variable, 307
- icmpInErrors variable, 307
- icmpInFlags variable, 307
- icmpInFirmProbe variable, 307
- icmpInRedirects variable, 307
- icmpInSrcQuenchs variable, 307
- icmpInTimeExcds variable, 307
- icmpInTimestampReps variable, 307
- icmpInTimestamps variable, 307
- icmpMaskRepl variable, 305, 319, 334, 1110
- icmpOutAddrMaskReps variable, 307
- icmpOutAddrMasks variable, 307
- icmpOutDestUnrechs variable, 307
- icmpOutEchoReps variable, 307
- icmpOutEchoes variable, 307
- icmpOutErrors variable, 307
- icmpOutFlags variable, 307
- icmpOutFirmProbs variable, 307
- icmpOutRedirects variable, 307
- icmpOutSrcQuenchs variable, 307
- icmpOutTimeExcds variable, 307
- icmpOutTimestampReps variable, 307
- icmpOutTimestamps variable, 307
- icmpsrc variable, 312, 315, 322-323
- icmpstat structure, 306-307
- icmpstat variable, 305-307
- icps\_badcode member, 306-307, 315
- icps\_badlen member, 306-307, 315, 319
- icps\_checksum member, 306-307, 313
- icps\_error member, 306, 324
- icps\_inhist member, 306-307
- icps\_oldicmp member, 306-307
- icps\_oldshort member, 306-307
- icps\_outhist member, 306-307, 317
- icps\_reflect member, 306, 317
- icps\_tooshort member, 306-307, 312
- IEEE (Institute of Electrical and Electronics Engineers), 69, 100, 106, 341
- INET (Internet Experiment Notes), 215
- IEIT (Internet Engineering Task Force), 350-351, 1125
- IF\_DEQUEUE macro, 25, 30, 72
- IF\_DROF macro, 69, 72
- IF\_ENQUEUE macro, 72
- IF\_PREFEND macro, 72
- IF\_PULL macro, 72
- if\_addrilen member, 69, 85, 89, 94
- if\_addrlist member, 66, 91, 166, 636, 656
- if\_addr variable, 94
- if\_attach function, 66, 80-81, 84-92, 1099
- if\_baudrate member, 69
- if\_bpf member, 68, 1030, 1036-1037, 1040
- if\_collisions member, 69, 97-98

- if\_data structure, 628
- if\_done member, 70–71, 96, 127
- if\_down function, 122–123, 571, 627
- if\_flags member, 67–68, 83–85, 99, 121–122, 337
- if\_hdrlen member, 69, 85, 94
- if\_ifaces member, 97, 99, 106
- if\_errors member, 97–99, 138
- if\_ismcasts member, 97, 99
- if\_index member, 67, 87, 91, 99, 574, 636, 643, 681
- if\_index variable, 64–65, 86–87
- if\_indexlen variable, 64, 88, 91
- if\_init member, 71, 81, 96, 127
- if\_ioctl member, 71, 81, 96, 122, 124, 127, 163, 165, 168, 172, 177, 344, 356, 361–362
- if\_ipackets member, 97–99
- if\_lgdrops member, 69, 97, 99
- if\_lastchange member, 69, 97–98, 106
- if\_len member, 72
- if\_metric member, 69, 121, 123
- if\_msghdr structure, 570, 621, 627, 630, 643
- if\_atu member, 69, 81, 83, 85, 94, 97, 99
- if\_name member, 67, 81, 87–88, 90–91, 94, 97, 99
- if\_next member, 66, 91
- if\_noproto member, 69, 97, 99
- if\_oerrors member, 97, 99
- if\_opackets member, 97–99, 141
- if\_omcasts member, 97, 99, 141
- if\_opackets member, 97–99
- if\_output member, 71, 81, 96, 101, 107, 127, 129–130, 139–140, 150, 152, 234, 378, 677, 1046
- if\_pcount member, 68
- if\_qflush function, 72, 123
- if\_rasat member, 71, 81, 96, 127
- if\_slowtime function, 93–94
- if\_snd member, 72, 84, 99, 101, 131, 140–141, 148
- if\_start member, 71, 81, 96, 111–112, 127
- if\_timer member, 68, 93–94
- if\_type member, 68–69, 86, 92, 98–99
- if\_unit member, 67, 81, 87–88, 91, 97, 119, 125, 149, 153, 178
- if\_up function, 122–123, 571, 627
- if\_watchdog member, 68, 71, 94, 96, 127
- IPA\_ROUTE constant, 171
- ifa\_addr member, 73, 89–90, 168–169, 629, 656
- ifa\_addrlist member, 90
- ifa\_broadaddr member, 73
- ifa\_braddr member, 73, 168
- ifa\_flags member, 74
- ifa\_ifaddr member, 91
- ifa\_ifp member, 73, 90–91, 158–159, 346
- ifa\_ifwithaddr function, 182, 264, 409, 731–732, 740, 1060, 1082
- ifa\_ifwithaf function, 182
- ifa\_ifwithdst function, 257
- ifa\_ifwithdstaddr function, 182, 232
- ifa\_ifwithnet function, 86, 182, 237, 619, 1071
- ifa\_ifwithroute function, 182, 609, 617
- ifa\_metric member, 74
- ifa\_msghdr structure, 570, 621, 629–630, 643
- ifa\_netmask member, 73, 89–91, 168, 615
- ifa\_next member, 73, 91, 159, 346
- ifa\_refcnt member, 74
- ifa\_rrequest member, 74, 90, 608, 611, 617, 679, 704
- ifa\_withdstaddr function, 738
- ifa\_withnet function, 738
- ifaddr structure, 66, 73–74, 76–78, 86–87, 89–90, 94, 120, 125, 155, 158–161, 166, 168, 178, 182–183, 232, 568–569, 581, 605, 609, 611, 615, 617, 636, 643, 656, 676–677, 704, 1071, 1098
- reference counts, 74, 177
- ifAdminStatus variable, 98–99
- ifafree function, 74, 605
- IPAFREE macro, 74, 177, 605
- ifam\_addrs member, 621–622
- ifaof\_ifpforaddr function, 182, 264, 319–320, 335, 1071
- ifc\_buf member, 117–118, 120–121
- ifc\_len member, 117–118, 120–121
- ifconf function, 115–120
- ifconf structure, 114, 117–118, 120–121
- ifconf program, 86, 105, 123, 162, 174, 183, 561, 679, 1004, 1109–1110
- ifDescr variable, 99
- ifEntry variable, 98
- IFF\_ALLMULTI constant, 67, 363
- IFF\_BROADCAST constant, 67–68, 74, 81, 234
- IFF\_CANTCHANGE constant, 68, 122
- IFF\_DEBUG constant, 67
- IFF\_LINK0 constant, 67, 83
- IFF\_LINK1 constant, 67, 83
- IFF\_LINK2 constant, 67, 83
- IFF\_LOOPBACK constant, 67, 85
- IFF\_MULTICAST constant, 67, 81, 84–85, 337, 1111
- IFF\_NOARP constant, 67
- IFF\_NOTRAILERS constant, 67
- IFF\_OACTIVE constant, 67, 112
- IFF\_POINTOPOINT constant, 67–68, 74, 84
- IFF\_PROMISC constant, 67, 125
- IFF\_RUNNING constant, 67
- IFF\_SIMPLEX constant, 67, 81, 150

- IFF\_UP constant, 67, 97, 99, 123, 125, 179-180
- ifIndex variable, 99, 574, 681
- ifInDiscards variable, 99
- ifInErrors variable, 99
- ifInIt function, 79, 84, 93-94
- ifInNUcastPkts variable, 99
- ifInOctets variable, 99
- ifInUcastPkts variable, 99
- ifInUnknownProtos variable, 99
- ifIoctl function, 115-116, 121-125, 149, 163-164, 166, 344-345, 451, 552, 554
- ifLastChange variable, 98-99
- ifm\_addr member, 621-622, 628
- ifmtu variable, 99
- ifnet structure, 33, 64-74, 76-78, 80-81, 83-87, 90-92, 94-95, 97-101, 105, 107, 112, 116, 120, 128-131, 140, 152, 158-159, 164, 166, 181-182, 232, 337, 340, 343, 345-346, 366, 406-407, 409, 427, 568-569, 581, 611, 636, 640, 643, 676-677, 681, 689, 706, 1029-1030, 1037, 1046, 1071, 1076, 1080, 1098, 1101, 1111
- ifnet utility functions, 182
- ifnet variable, 64, 86-87, 91, 94, 118, 120, 159, 182, 346, 366
- IFNET\_SLOW2 constant, 94
- ifnet\_addr variable, 64, 86-87, 90-91, 120, 158-159
- ifNumber variable, 65
- ifOperStatus variable, 99
- ifOutDiscards variable, 99
- ifOutErrors variable, 99
- ifOutNUcastPkts variable, 99
- ifOutOctets variable, 99
- ifOutQLen variable, 99
- ifOutUcastPkts variable, 99
- IFF\_TO\_IA macro, 346, 371
- ifPhysAddress variable, 99
- ifPromise function, 1035
- IFO\_MAXLEN constant, 72, 200
- ifq\_drops member, 69, 72, 97, 99
- ifq\_head member, 72
- ifq\_len member, 97, 99
- ifq\_maxlen member, 72
- ifq\_tail member, 72
- ifqmaxlen variable, 64, 72, 84, 93
- ifqueue structure, 72, 97, 207, 680
- ifr\_addr member, 117-118, 121, 168
- ifr\_flags member, 121-122
- ifr\_metric member, 121
- ifr\_name member, 116-119, 121, 1039
- ifra\_addr member, 175
- ifra\_name member, 174
- ifreq structure, 114, 116-119, 121, 162, 164, 168, 172, 174-175, 344, 361, 1035, 1039, 1076
- ifSpecific variable, 98-99
- ifSpeed variable, 98-99
- IFT\_BETHER constant, 69, 92, 121, 706-707
- IFT\_FDDI constant, 69
- IFT\_IS086023 constant, 69
- IFT\_IS086025 constant, 69
- IFT\_LOOP constant, 69, 121
- IFT\_OTHER constant, 69
- IFT\_SLIP constant, 69, 121
- ifTable variable, 98-99
- ifType variable, 99
- ifUnit function, 116, 182-183
- IGMP (Internet Group Management Protocol), 188, 191, 193, 228, 240, 337-338, 346, 373, 381-396, 401, 403, 411-415, 417-419, 440, 451, 477, 541, 1127
  - header, 385
  - inflow variable, 385
  - input processing, 391-395
  - protocol structure, 384-385
  - RFC 1122 compliance, 1111
  - igmp structure, 384-385
- IGMP\_HOST\_MEMBERSHIP\_QUERY constant, 384, 391-394, 411
- IGMP\_HOST\_MEMBERSHIP\_REPORT constant, 384, 394-395, 411-412
- IGMP\_MAX\_HOST\_REPORT\_DELAY constant, 386
- IGMP\_MINLEN constant, 391, 1077
- IGMP\_RANDOM\_DELAY macro, 386-387, 393
- igmp\_all\_hosts\_group variable, 383
- igmp\_checksum member, 384-385
- igmp\_code member, 384-385
- igmp\_fasttimo function, 381-382, 385-389, 393, 1052
- igmp\_group member, 384-385
- igmp\_inic function, 194, 385, 1052
- igmp\_input function, 381-382, 385, 391, 393, 412, 419, 1050, 1052-1053
- igmp\_joining function, 356, 361, 381-382, 386-387
- igmp\_leavesgroup function, 356, 368, 395
- igmp\_report function, 395
- igmp\_sendreport function, 111, 381-382, 389-391
- igmp\_timers\_are\_running variable, 383, 386, 388, 1084
- igmp\_type member, 384-385, 391
- igmpstat structure, 383
- igmpstat variable, 382-383
- igps\_rcv\_badqueries member, 383, 393



- igmp\_rcv\_badreports member, 383
- igmp\_rcv\_badsum member, 383
- igmp\_rcv\_outreports member, 383
- igmp\_rcv\_queries member, 383
- igmp\_rcv\_reports member, 383
- igmp\_rcv\_tooshort member, 383
- igmp\_rcv\_total member, 383
- igmp\_snd\_reports member, 383
- imo\_membership member, 348, 359, 366
- imo\_multicast\_ifn member, 347, 366
- imo\_multicast\_loop member, 348, 355, 371, 376
- imo\_multicast\_rtl member, 348, 354–355, 371
- imo\_num\_memberships member, 348
- IMP (Interface Message Processor), 24, 40, 744, 1060
- implementation, Berkeley-derived, 4
- implementation refinements
  - TCP, 994–995
  - UDP, 791–792
- in\_ interface member, 356, 358, 366
- in\_multiaddr member, 356, 358, 366
- IN\_FIRST\_MULTI macro, 387–388
- IN\_LOOKUP\_MULTI macro, 346–347, 359, 373, 376
- IN\_MULTICAST macro, 943–944, 993
- IN\_MULTICAST macro, 387–388
- in\_addmulti function, 171, 356, 359–361, 386, 1076
- in\_addr structure, 160–162, 258, 271, 348, 401, 406–407, 410, 1071
- in\_allareq structure, 162, 164, 174–176
- in\_arppinput function, 109, 679, 687–694, 696, 701, 711, 1080–1081
- in\_broadcast function, 181, 943–944
- in\_forward function, 181, 221, 245
- in\_ifsum function, 216, 234–239, 245, 313, 333, 768, 792, 1087
- VAX, 239
- in\_ifcontrol function, 162–168, 171–177, 451, 615, 785, 1007
- in\_ifmulti function, 356, 359, 366, 368–369, 380, 395, 1076
- in\_ifaddr structure, 77–78, 89, 155, 158–162, 164, 166–168, 171, 174–176, 183, 207, 218, 253, 319, 329, 345–346, 348, 361, 366, 387, 581, 676–677, 689, 738, 1099
- in\_ifaddr variable, 158–159, 177, 207, 215, 218–219, 329, 346, 387
- in\_ifinit function, 168–171, 175–177, 615, 677, 679, 704, 706, 1111
- in\_ifscrub function, 169, 176, 615
- in\_ifwithout function, 232
- in\_interfaces variable, 158, 166
- in\_localaddr function, 181, 901
- in\_losing function, 571, 749–750, 843, 1112–1113, 1121
- in\_options structure, 366
- in\_multi structure, 161, 345–346, 348, 356, 359, 361, 365–366, 368–369, 379, 381, 382, 386–389, 395
  - reference counts, 359–360, 368, 386, 395
- in\_multistep structure, 387
- in\_netof function, 181, 421
- in\_pcballoc function, 715, 717–719, 785, 1018, 1088, 1090
- in\_pcbbind function, 451, 725, 728–735, 740, 751–753, 763, 786, 1010–1012, 1081, 1083
- in\_pcbconflict function, 725
- in\_pcbconnect function, 572, 728–730, 735–741, 751–752, 763, 787, 944, 1012, 1060, 1083, 1114–1115, 1118
- in\_pcbdetach function, 715, 717–719, 741, 786, 1018, 1059
- in\_pcbdisconnect function, 738, 741, 768, 788
- in\_pcblookup function, 722–728, 730, 732–734, 738–740, 750–751, 773–774, 777, 785, 929, 1055
- in\_pcbnotify function, 742–746, 748–749, 783
- in\_pcbopts function, 240
- in\_rtchange function, 303–304, 743–744, 746
- in\_scrub function, 175
- in\_serpeeraddr function, 741–742, 789, 1017, 1062
- in\_setsockaddr function, 741–742, 789, 1017, 1062
- in\_sockmask member, 170
- in\_socketim function, 170, 183
- in\_ulomove function, 792
- in\_uniqueport function, 725
- INADDR\_ALLHOSTS\_GROUP constant, 338–339
- INADDR\_ANY constant, 219, 353, 358, 363, 371, 379, 389, 409, 736, 741, 1101, 1114
- INADDR\_BROADCAST constant, 219, 736, 1101
- INADDR\_MAX\_LOCAL\_GROUP constant, 338–339, 428
- INADDR\_NO\_IFF macro, 354, 358, 1076
- INADDR\_UNSPEC\_GROUP constant, 338–339
- incarnation, connection, old, 814
- indirect route, 561, 569, 580, 608, 615, 706
- inet\_addr function, 7
- inetcterrmap variable, 744, 782–783, 904
- inetd program, 555
- lnetdomain variable, 186–187, 193, 195, 204, 309, 385, 758, 801, 1052



- inetaw variable, 186, 191-192, 195, 199-200, 204, 220, 228, 286, 310, 1051, 1091
- ICMP, 203, 309-310, 1098
- IGMP, 385
- IP, 191, 198, 203, 228, 1098
- raw IP, 191, 193, 197, 199, 204, 1052, 1072
- TCP, 198-199, 801
- UDP, 203, 758
- init\_sysent.c file, 443
- initial send sequence number, *see* ISS
- initial sequence number, *see* ISN
- inm\_addr member, 345
- inm\_ia member, 345-346, 366
- inm\_ifp member, 345-346, 366
- inm\_next member, 346, 366
- inm\_recount member, 346, 366
- inm\_timer member, 346, 366, 386
- INP\_CONTROLOPTS constant, 717, 776
- INP\_HDRINCL constant, 717
- INP\_RECVSTADDER constant, 242, 717
- INP\_RECVOPTS constant, 242, 717
- INP\_RECVRTOPTS constant, 242, 717
- inp\_back member, 732
- inp\_laddr member, 14, 715-716, 748, 752, 774, 800
- inp\_flags member, 717
- inp\_fport member, 14, 716, 748, 752, 774, 800
- inp\_head member, 716, 718
- inp\_ip member, 717
- inp\_laddr member, 14, 716, 740, 752, 758, 774, 800, 1115
- inp\_lport member, 14, 21, 715-716, 752, 758, 760, 800, 813, 1081
- inp\_options member, 347, 717
- inp\_next member, 14, 715-716, 732, 824
- inp\_options member, 347, 717, 768, 932
- inp\_ppcb member, 714-715, 717
- inp\_prev member, 14, 715-716
- inp\_route member, 347, 717
- inp\_socket member, 13-14, 21, 347, 714, 717-718
- inpcb structure, 13-14, 440, 568, 672, 714-717, 732
- INPROCKUP\_WILDCARD constant, 727-728, 730, 732-734, 774
- input processing
  - ICMP, 310-323
  - IGMP, 391-395
  - IP, 212-220
  - IP multicast, 373
  - TCP, 923-1005
  - UDP, 769-780
- inseq function, 55, 291-292, 671, 683, 707, 718, 913
- Institute of Electrical and Electronics Engineers, *see* IEEE
- interface address, IP source address versus outgoing, 740-741
- interface layer, 10, 63-94
- Interface Message Processor, *see* IMP
- internal node, routing table, 564
- International Atomic Time, *see* TAI
- International Organization for Standardization, *see* ISO
- Internet address family, 185
- Internet Assigned Numbers Authority, *see* IANA
- Internet Control Message Protocol, *see* ICMP
- Internet diameter, 223, 1129
- Internet domain, 75, 160, 193, 309, 316, 385, 460, 483
- Internet domain structure, 186, 191-196
- Internet Engineering Task Force, *see* IETF
- Internet Experiment Notes, *see* IEN
- Internet Group Management Protocol, *see* IGMP
- Internet Protocol, *see* IP
- Internet protocol family, 185, 202-203, 361
- interprocess communication, *see* IPC
- interrupt, 95, 100-101, 141, 312, 469-470, 496, 499, 525, 528
  - hardware, 205
  - levels and concurrency, 23-26
  - network, 73, 138, 148, 213, 469, 1078
  - service routine, *see* isr
  - software, 106-107, 138, 153, 205, 212, 436
- IOCBASBCMD macro, 554
- IOCGROUP macro, 554
- IOCPARM\_LEN macro, 554
- ioctl function, 13, 74, 84, 95-96, 114-117, 120, 125, 127, 130, 132, 149, 159, 162-166, 173-174, 177, 183, 344, 348, 380, 439-440, 445, 447, 451, 506-507, 537-538, 548-549, 552, 554, 557, 569, 637, 666, 683, 785, 1007, 1027, 1032-1035, 1076, 1083, 1092, 1117
- iov\_base member, 481, 483, 486
- iov\_len member, 481, 483, 487
- iovent member, 481
- iovec structure, 481, 483-487, 493, 501-503
- IP (Internet Protocol), 65, 191
  - address, 155-183
  - address assignment, 161-177
  - broadcast address, 182, 234, 1101
  - checksum, 234-239
  - demultiplexing, 219
  - destination address, 162, 182

- domain initialization, 399
- experimental address, 156
- forwarding, 181, 217–228
- fragmentation, 48–50, 210, 232, 275–283
- header, 210–212
- `inetaw` variable, 191, 198, 203, 228, 1098
- input processing, 212–220
- input queue, 97, 106
- multicast address, 155–156, 341
- multicast groups, well-known, 338–339
- multicast input processing, 373
- multicast routing, 397–434
- multicasting, 155–156, 337–380
  - option class, 249
  - options, 247–273
  - output processing, 228–234
  - packet, 210
- `netow` structure, 186, 191–196
- `raw`, 183, 191, 197, 230, 240, 276, 301, 304–305, 312–313, 304, 391, 440, 451, 477, 541, 1049–1065, 1071, 1102, 1108
- reassembly, 48–50, 219, 275–277, 283–300
- RFC 1122 compliance, 1098–1105
- route selection, 230–232
- source address, 232
- source address versus outgoing interface address, 740–741
- subnetting, 156, 170, 181, 1071
- subnetting, and redirects, 226
- supernetting, 170, 1071
- to Ethernet multicast address mapping, 341–342
- unicast address, 155–156, 182
- utility functions, 181
- IP address, class A, 155
  - class B, 155
  - class C, 155
  - class D, 155, 341
  - class E, 155
- `IP_ADD_MEMBERSHIP` socket option, 348, 356–357, 434, 451, 1076
- `IP_ALLOWBROADCAST` constant, 229–230, 234, 333
- `IP_DEFAULT_MULTICAST_LOOP` constant, 371
- `IP_DEFAULT_MULTICAST_TTL` constant, 351, 371, 376
- `IP_DF` constant, 276, 325
- `IP_DROP_MEMBERSHIP` socket option, 348, 366, 451
- `IP_FORWARDING` constant, 228–230, 376, 400, 430
- `IP_HDRINCL` socket option, 191, 1053, 1056–1058, 1065, 1091
- `IP_HDR_LEN` constant, 432
- `IP_MAX_MEMBERSHIPS` constant, 380
- `IP_MAXPACKET` constant, 265
- `IP_MF` constant, 276, 325
- `IP_MULTICAST_IF` socket option, 348, 353–354, 371, 710, 738
- `IP_MULTICAST_LOOP` socket option, 348, 354–355, 371
- `IP_MULTICAST_TTL` socket option, 348, 354, 371
- `IP_OPTIONS` socket option, 230, 240, 242–243, 269–270, 717, 1056, 1065, 1113, 1115, 1122
- `IP_RAWOUTPUT` constant, 230, 1056–1057
- `IP_RECVSTADDR` socket option, 240, 242, 776, 781, 793–794, 1083, 1114
- `IP_RECVOPTS` socket option, 240, 242–243, 769, 776, 794
- `IP_RECVRETOPTS` socket option, 240, 242–243, 776
- `IP_RETOPTS` socket option, 794
- `IP_ROUTEIOIF` constant, 229, 232
- `IP_TOS` socket option, 240, 242, 717, 1056, 1099, 1115, 1123
- `IP_TTL` socket option, 240, 242, 717, 1056, 1100, 1115, 1120
- `ip_okaum` member, 211, 1099
- `ip_qloutput` function, 240–243, 348, 538, 541, 758, 1023, 1065
- `ip_defctl` variable, 207, 209, 785, 835, 889, 1100
- `ip_deq` function, 292
- `ip_dboptions` function, 217–218, 220, 249–265, 273, 283, 405, 1102–1104, 1108
- `ip_drain` function, 193, 298–300, 892
- `ip_dst` member, 211, 218–219, 252, 255, 257, 266, 277, 285–286, 293, 298, 329, 405, 1073
- `ip_enq` function, 292, 294
- `ip_forward` function, 205–206, 219–228, 232, 245–246, 250, 252, 265, 279, 572, 1081, 1101, 1103, 1107
- `ip_freef` function, 298–299
- `ip_freemoptions` function, 719
- `ip_getmoptions` function, 348, 371, 380, 1076
- `ip_hl` member, 211–212, 215, 230, 247, 328, 427, 432, 809
- `ip_hlen` member, 283, 285
- `ip_id` member, 200, 211, 216, 230, 275, 277, 285–286, 289, 293, 373, 1056
- `ip_id` variable, 200, 207, 230, 1000–1001, 1119
- `ip_ifmatrix` variable, 186, 200, 223
- `ip_init` function, 186, 193–194, 199–201
- `ip_insertoptions` function, 230, 248, 265–269, 272, 1082, 1086
  - TCP example, 267–268
  - UDP example, 268

- ip\_len member, 211–212, 216, 223, 234, 267, 277, 281, 285, 298, 312, 327–328, 378, 427, 432, 771–772, 774, 809, 926
- ip\_forward function, 221, 228, 373, 376–377, 400, 405, 409, 424–430, 434, 1077
- ip\_loopback function, 376–378, 400, 425, 427
- ip\_options structure, 347–348, 351, 353–354, 357, 365–366, 371, 375–376, 379–380, 389, 399, 430, 1076
- ip\_mreq structure, 348, 356, 358, 366
- ip\_mrouter variable, 340, 373, 403, 424, 434, 1059
- ip\_mrouter\_end function, 401, 403, 412, 1065
- ip\_mrouter\_done function, 401, 433–434, 1059
- ip\_mrouter\_init function, 401, 403–404
- ip\_nhops variable, 218, 248, 257–259, 261
- ip\_off member, 211, 216, 230, 234, 275–277, 279, 281, 283, 285–286, 293, 300, 325, 335, 378, 1074
- ip\_optcopy function, 279, 282–283
- ip\_output function, 107, 193, 206–208, 209, 216, 220–221, 223, 226, 228–234, 245, 265, 269, 278–281, 293, 300, 324, 326, 329, 332–333, 335, 347, 351, 354, 375–378, 381–382, 389, 399–400, 424–425, 427, 430, 433–434, 569, 572, 578–579, 606, 679, 684, 699–700, 710, 717, 738, 741, 746, 756, 761, 764, 767–768, 789, 793, 796, 852, 882, 885, 887–889, 906, 917, 932, 962, 1050, 1056–1057, 1062, 1076–1077, 1082, 1084, 1086, 1090–1092, 1099–1105, 1107–1108, 1113–1115, 1118–1119
- ip\_p member, 199, 211–212, 220, 277, 285–286, 293, 309–310, 315, 328, 385, 801, 1052, 1065, 1091
- ip\_pcbopts function, 242, 269–272, 717
- ip\_protox variable, 199–200, 204, 207, 220, 286, 310, 1052–1053, 1091
- ip\_reasm function, 283, 285–298, 300, 1099, 1104
- ip\_rtaddr function, 253–254, 257, 405, 572
- ip\_setmoptions function, 348, 351–359, 366, 381–382, 572, 1076
- ip\_slowtimo function, 94, 193, 292, 298–300
- ip\_src member, 211, 233, 255, 259, 277, 285–286, 293, 298, 321, 329, 376, 405
- ip\_srcroute function, 258–261, 265, 272, 332, 932, 1108–1109, 1122
- ip\_srcrt structure, 248, 258–259, 261, 265, 273
- ip\_srcrt variable, 248, 257–259, 261
- ip\_striptoptions function, 769, 925, 1122
- ip\_sum member, 212, 216
- ip\_sysctl function, 193, 203, 244–245
- ip\_timestamp structure, 248, 261
- ip\_tos member, 140, 211, 242, 244, 285, 287, 328, 1115
- ip\_ttl member, 211, 221, 223, 242, 244, 246, 329, 350, 376, 428
- ip\_v member, 211–212, 215, 230
- ip structure, 210–211, 250, 261, 286–289, 298, 324, 1077
- ipasfrag structure, 285–289, 292, 297–298
- IPC (interprocess communication), 9, 185
- IPCTL\_DEPTTL constant, 201–202, 244, 1100
- IPCTL\_FORWARDING constant, 201–202, 244
- IPCTL\_SSHUREDIRECTS constant, 201–202, 244
- ipDefaultTTL variable, 209
- IPDEFTTL constant, 1100, 1120
- ipEff member, 285–288, 294, 298
- ipf\_next member, 287–289
- ipf\_prev member, 287–289
- ipf\_tos member, 298
- ipforward\_rt variable, 223, 226
- ipforwarding variable, 157–158, 207, 209, 219–220, 226, 246, 252, 1098–1099
- ipForwarding variable, 209
- ipForwDatagrams variable, 209
- ipFragCreate variable, 209
- ipFragFail variable, 209
- ipFragOK variable, 209
- IPFRAGTTL constant, 209, 292, 1105
- ipInAddrErrors variable, 209
- ipInDelivers variable, 209, 1055
- ipInDiscards variable, 209, 1055
- ipInHdrErrors variable, 209
- ipInReceives variable, 209
- ipintr function, 101, 107, 131, 150, 205–206, 212–221, 223, 226, 245, 247, 249, 252, 258–259, 283, 285–286, 290, 300, 310, 312, 332, 373, 379, 382, 391, 399–400, 405, 409, 412, 424–425, 427, 743, 756, 771, 776, 796, 923, 962, 1050, 1076–1077, 1099–1102, 1106–1107, 1120
- ipintrq variable, 25, 97, 101, 106, 130–131, 138, 150, 200, 205–207, 212–213, 220, 373, 400, 424–425, 962
- ipInUnknownProto variable, 209
- ipNetToMediaIfIndex variable, 681
- ipNetToMediaNetAddress variable, 681
- ipNetToMediaPhysAddress variable, 681
- ipNetToMediaTable variable, 681
- ipNetToMediaType variable, 681
- IPOPT\_COPIED constant, 283
- IPOPT\_EOL constant, 249
- IPOPT\_LSRR constant, 249, 257
- IPOPT\_MINOFF constant, 250, 257

- IPDPT\_NOP constant, 249, 261
- IPDPT\_OFFSET constant, 250
- IPOPT\_OLEN constant, 250
- IPOPT\_OPTVAL constant, 250
- IPOPT\_RF constant, 249
- IPOPT\_SATID constant, 249
- IPOPT\_SECURITY constant, 249
- IPOPT\_SSRR constant, 249, 257, 261
- IPOPT\_TS constant, 249
- IPOPT\_TS\_PRESPFC constant, 261, 264
- IPOPT\_TS\_TSANDADDR constant, 261, 264
- IPOPT\_TS\_TSONLY constant, 261, 264
- ipopt\_dst member, 265–266
- ipopt\_list member, 265, 267
- ipoption structure, 248, 259, 265, 267, 269, 271, 332
- ipOutDiscards variable, 209
- ipOutRoutes variable, 209
- ipOutRequests variable, 209
- ipovly structure, 760, 764, 772, 803, 906
- IPPORT\_RESERVED constant, 732–733
- IPPORT\_USERRESERVED constant, 733
- IPPROTO\_ICMP constant, 191, 203, 309, 1051
- IPPROTO\_IGMP constant, 191, 201, 385, 1051
- IPPROTO\_IP constant, 240
- IPPROTO\_RAW constant, 191, 193, 200, 1051–1053, 1065, 1091
- IPPROTO\_TCP constant, 191, 196, 240, 801, 1022
- IPPROTO\_UDP constant, 191, 203, 758
- ipq structure, 277, 286–291, 293–294, 296, 298, 300
- ipq variable, 200, 277, 285–289
- ipq\_dst member, 288
- ipq\_id member, 288
- ipq\_next member, 286, 288–289, 293
- ipq\_p member, 288
- ipq\_prev member, 286, 288–289, 293
- ipq\_src member, 288
- ipq\_ttl member, 288, 292
- ipqmaxlen variable, 107, 200, 207
- ipReasmFails variable, 209
- ipReasmOKs variable, 209
- ipReasmReqds variable, 209
- ipReasmTimeout variable, 209
- ipRouteAge variable, 574
- ipRouteDest variable, 574
- ipRouteIfIndex variable, 574
- ipRouteInfo variable, 574
- ipRouteMask variable, 574
- ipRouteMetric1 variable, 574
- ipRouteMetric2 variable, 574
- ipRouteMetric3 variable, 574
- ipRouteMetric4 variable, 574
- ipRouteMetric5 variable, 574
- ipRouteNextHop variable, 574
- ipRouteProto variable, 573–574
- ipRouteTable variable, 573–574
- ipRouteType variable, 573–574
- ips\_badhlen member, 207–209
- ips\_badlen member, 207–209
- ips\_badoptions member, 207–209, 248
- ips\_badsum member, 207–209
- ips\_badvers member, 207–209
- ips\_cantforward member, 207–209, 340, 373
- ips\_cantfrag member, 207–209, 278
- ips\_delivered member, 207–209, 1055
- ips\_forward member, 207–209, 340
- ips\_fragdropped member, 207–209, 291
- ips\_fragmented member, 207–209, 278
- ips\_fragments member, 207–209
- ips\_fragtimeout member, 207–209
- ips\_localout member, 207–209
- ips\_noprots member, 207–209, 1051, 1055–1056
- ips\_noroute member, 207–209, 340
- ips\_ogropped member, 207–209, 278
- ips\_ofragments member, 207–209, 278
- ips\_rawout member, 207–208, 1051, 1057
- ips\_reassembled member, 207–209
- ips\_redirectsent member, 207–208
- ips\_tooshort member, 207–209
- ips\_toosmall member, 48, 207–209
- ips\_toomsi member, 207–209
- ipsendredirects variable, 207, 225
- ipstat structure, 207–209, 248, 278, 340, 1051
- ipstat variable, 207, 278, 340
- ipt\_code member, 262
- ipt\_flg member, 261–262, 264
- ipt\_len member, 262
- ipt\_oflw member, 262
- ipt\_pkt member, 262
- ipt\_ts member, 262
- ipt\_time member, 262
- iptime function, 264, 1109
- IPTOS\_LOWDELAY constant, 140
- IPTTLDEC constant, 223
- IPv4, 1068
- IPv5, 215
- IPv6, 215, 1068
- IPVERSION constant, 215, 230
- ISN (initial sequence number), 1118
- ISO (International Organization for Standardization), 100
- ISO Development Environment, see ISODE
- iso\_ifaddr structure, 77–78, 159



- ISODE (ISO Development Environment), 65, 69, 98, 1095
- isodomain variable, 187, 195
- isow variable, 195
- isr (interrupt service routine), 73, 95, 106, 291
- ISS (initial send sequence number), 797, 812-814, 945, 949, 959, 968, 1012
- Itano, W. M., 106, 1127
- itimexfix function, 525
- Jacobson, V., 5, 60, 112, 147, 419, 500, 562, 800, 824, 831-832, 866, 934, 970, 995, 1027, 1067-1068, 1126-1130
- Jain, R., 223, 750, 1128
- Kacker, M., xxii
- Karels, M. J., 5, 24, 33-34, 94, 129, 445, 457, 470, 562, 1097, 1128-1129
- Karr, P., 144, 844, 1128
- Karn's algorithm, 844, 976, 1119, 1128
- Kastenholz, F. J., 226, 1125-1126
- Kay, J., 234, 1128
- keepalive
  - probe, 818, 820, 822, 828-830, 848-849, 887-888, 1085, 1121
  - timer, 818-819, 821, 828-831, 842, 887, 892, 918, 933
- Kent, C. A., 1128
- Kercheval, B., xxii
- kernel statistics, 37
- Kernighan, B. W., xxii, 1067, 1128
- Kieber, U., xxii
- Krol, E., 1128
- ktrace program, 30
- Kullberg, A., 235, 1129-1130
- kvm function, 37
- l\_linger member, 542, 547
- l\_sndeff member, 542, 547
- la\_asked member, 683, 692, 695, 699-700, 707
- la\_hold member, 678, 683, 692, 696, 699-700, 707, 1113
- la\_next member, 677
- la\_prev member, 677
- la\_rt member, 677, 683
- LAN (local area network), 711, 827, 841, 845, 903, 940
- Lanciani, D., 959, 1128
- last\_ack\_sent member, 867, 869, 871, 884
- Laubach, M., xxii
- layering, 1068, 1126
- le\_softc structure, 77, 80, 86, 91, 97, 100-101, 120, 125, 159, 178, 309, 343, 346, 366, 581, 677, 683, 1030, 1036-1037, 1071
- le\_softc variable, 77, 80, 86, 97, 100-101, 120, 153, 159, 343, 346, 1030, 1036-1037
- leaf node, routing table, 564
- leattach function, 80-83, 86, 91-92
- leaves, in multicast spanning trees, 417, 419
- Leffler, S. J., 24, 33-34, 94, 129, 445, 457, 470, 562, 1097, 1129
- leinit function, 96, 127, 178
- leintr function, 100-103
- leioctl function, 96, 115, 122, 124-125, 127, 163, 168, 177-178, 362-366, 1077
- len member, 16, 20-21, 32, 39, 53, 58, 987
- leput function, 112
- leread function, 101-104, 106, 125, 1043, 1070
- Leres, C., 725
- lerset function, 96, 127, 362-363, 1077
- lstart function, 96, 101, 112, 127
- LETFBUF constant, 112
- lge\_gaddr member, 411
- lge\_vifid member, 411
- lgrpletl structure, 401, 411
- Li, T., 170, 1127
- library functions, system calls and, 7-8
- limited broadcast address, 753
- Lin, J. C., 1018, 1126
- line discipline, 129-130, 148-149
  - SLIP, 129-132, 134, 149
- linger structure, 539, 542, 1079
- link layer, RFC 1122 compliance, 1097-1098
- link\_rtrequest function, 90
- link-level address, 77, 85-90, 92, 94, 97, 158, 185
- link-level address mask, 77
- listen function, 437, 440, 445-446, 455, 457, 459, 463-464, 468, 524, 729-730, 786, 805, 930-931, 1010-1011
- Liu, C., 419, 1127
- LLADDR macro, 87
- llinfo\_arp structure, 677, 680, 682, 691, 695-697, 701, 704, 706-707, 710-711, 1113
- llinfo\_arp variable, 677, 680
- local area network, see LAN
- local multicast group, 339
- locking, mbuf, 43
- loioctl function, 361
- loif variable, 77, 86, 120, 128, 159, 1030
- loioctl function, 115, 124, 163, 168, 177, 180, 356
- LOHTU constant, 85
- long fat pipe, 866, 1128
- loop\_rtrequest function, 90



- loopattach function, 78, 80, 85–86
- loopback
  - device driver, 64, 128, 150–153, 161, 1029
  - initialization, 85
  - MTU, 85
  - multicast, 400
  - network, 156, 181
  - packet, BPF, 152
  - pseudo-device, 64
  - sockaddr\_d structure, 1070
- loose source and record route, see LSRR
- loopout function, 96, 110, 127, 150–153, 212, 378, 962, 1070
- low-water mark, 477, 479, 496, 530–531, 534, 539, 543
- LSRR (loose source and record route), 249–250, 254–255, 257, 270–271, 283, 398, 404–406, 427, 432–434, 1073
  - multicast tunnel, 427, 430
- Lumley, J., 994, 1126
- Lynch, D. C., 1129
  
- M\_BCAST constant, 3, 39, 103, 125, 221, 234, 245, 325, 697, 943, 1098, 1101
- M\_BUF constant, 716
- M\_COPYALL constant, 53
- M\_COPYFLAGS constant, 39, 61
- M\_DONTWAIT constant, 41–42, 52–53, 88, 585, 763, 781, 874, 884
- M\_EOR constant, 39, 491, 498, 508, 1078
- M\_EXT constant, 31, 33–34, 39, 46, 52, 60–61, 1069
- M\_FREE constant, 40
- M\_FTABLE constant, 40
- M\_HTABLE constant, 40
- M\_IFADDR constant, 40
- M\_LEADINGSPACE macro, 764
- M\_MBUF constant, 40
- M\_BCAST constant, 39, 104, 125, 325, 373, 376, 697, 943, 1098, 1101
- M\_NOWAIT constant, 491
- M\_PCB constant, 40, 715
- M\_PKTDR constant, 16–17, 20, 31–34, 39, 52, 60
- M\_PREPEND macro, 52, 111, 267, 763–764, 793, 1056, 1082
- M\_READABLE constant, 40, 585
- M\_SOCKET constant, 40, 716
- M\_SODOTS constant, 40
- M\_WAIT constant, 41–42, 52, 763
- M\_WAITOK constant, 88, 478
- m\_sact member, 34
- m\_adj function, 53, 771, 912–913, 951, 956, 959, 1088
- m\_cat function, 53, 294, 296–297
- m\_clfree member, 36
- m\_clusters member, 36
- m\_copy function, 53, 56–61, 110, 223, 279, 430, 432, 664, 777, 874, 1055, 1087
- m\_copyback function, 53, 632, 650
- m\_copydata function, 52–53, 650–651, 655, 874
- m\_copym function, 53
- m\_dat member, 31
- m\_data member, 15, 17, 21, 31, 33–34, 47–48, 52, 152, 267–268, 279, 288, 291, 294, 298, 432, 685
- m\_devget function, 44–51, 53, 60, 101, 104, 283, 909
- m\_drain member, 36–37
- m\_drops member, 36–37, 61
- m\_ext structure, 33–34
- m\_flags member, 15–16, 20, 31–32, 34, 39, 774
- m\_free function, 53, 73, 271, 542, 719, 887, 897
- m\_freeen function, 53, 58, 73, 763
- m\_get function, 41, 43, 53, 371, 839, 1024
- m\_getclr function, 53, 554
- m\_gethdr function, 53, 326, 887
- m\_hdr structure, 32, 39, 195
- m\_len member, 15–17, 21, 31–33, 39, 45, 47–48, 52, 259, 267, 291, 298, 427, 432, 519, 546, 987
- m\_nbufs member, 36
- m\_ntypes member, 36–37, 42
- m\_next member, 15–16, 33–34, 39, 42, 52, 54, 776
- m\_nextpkt member, 15, 34, 39, 42, 54, 72, 281, 308, 515, 1081
- m\_pktdot member, 32, 45, 267–268
- m\_pkthdr member, 16, 21, 32–33, 39, 58, 195, 216, 298, 987, 1101
- m\_prepend function, 1082
- m\_pullup function, 44–51, 53, 60–61, 215–216, 283, 312, 391, 769, 909, 925, 927, 1073
- m\_reclaim function, 43, 796, 892
- m\_retry function, 42–43
- m\_space member, 36
- m\_types member, 15, 39, 42, 508
- m\_wait member, 36–37, 61
- machdep.c file, 291, 1069
- main function, 79, 82–83, 85, 93, 571, 584, 756, 796, 1050
- malloc function, 40, 88, 412, 484, 585, 635, 637, 650, 715, 834, 1018, 1087
- MALLOC macro, 40, 42–43, 483, 665, 718, 1087
- Mallory, T., 235, 1129–1130
- management information base, see MIB
- Mann, T., 100, 321, 1127
- Mano, M. M., 235, 1129
- manual pages, Unix, 3

- mapped pages, 33
- mask, lists, 587-591
- mask\_rhhead variable, 568-569, 572, 586, 654
- maskedKey variable, 572, 585, 595-597
- match
  - most specific, 522
  - wildcard, 722
- MAX\_IPOPTLEN constant, 258, 265
- max\_datalen variable, 186, 195
- max\_hdr variable, 186, 195, 498, 763
- max\_keylen variable, 572, 585, 594
- max\_linkhdr variable, 186, 194-195, 267-268, 279, 813, 1077, 1084, 1086
- max\_protohdr variable, 49, 186, 194-195, 813
- max\_rcvwd member, 989
- max\_sndwnd member, 859, 983, 989
- maximum segment lifetime, *see* MSL
- maximum segment size, *see* MSS
- maximum transmission unit, *see* MTU
- MAXTTL constant, 329
- MAXVIFS constant, 406
- MBONE (multicast backbone), 350-351, 1126
- mbstat structure, 36-37, 42
- mbtypes variable, 42
- mbuf, 15-19, 31-61
  - chain, 16, 34
  - cluster, 16, 33
  - external buffer, 33
  - locking, 43
  - packet header, 16, 32
  - queue of, 34, 39
- mbuf structure, 33, 38-40, 47, 267-268, 309
- MBUFLOCK macro, 42-43
- McCanne, S., xxii, 1027, 1129
- McCloghrie, K., 64, 399, 1129
- McGregor, G., 129, 1129
- McKenney, P. E., 750, 791, 994, 1129
- McKusick, M. K., xxii, 5, 24, 33-34, 94, 129, 445, 457, 470, 562, 1097, 1128-1129
- MCLALOC macro, 33, 43, 1069
- MCLBYTES constant, 33, 37, 134, 498, 534, 900-901, 1078, 1082
- MCLFREE macro, 33, 43, 1069
- MCLGET macro, 33, 52, 498
- mclrcfont variable, 1069
- memory leak, 434, 470, 1076
- Mendez, T., 351, 1130
- message boundaries, 508-509
- MF (more fragments flag, IP header), 275-277, 279, 281, 283, 285, 289
- MFREE macro, 43, 52-53
- MGET macro, 41-43, 47, 52-53, 1087
- MGETHER macro, 52-53, 266-267, 432, 764, 874, 1086-1087
- MIL\_ALIGN macro, 52, 326, 328, 498, 685, 1082-1083, 1086
- MHLEN constant, 37, 47, 557, 630, 813, 1069
- MIB (management information base), 64-65, 69, 383, 399, 680, 1129
- MIB-II, SNMP, 65
- microtime function, 264, 1043
- Milliken, W., 351, 1130
- MINCLBYTES constant, 46
- MINCLSIZE constant, 37, 498, 1087
- MLEN constant, 37, 539
- mmap function, 52
- Mogul, J. C., 33, 100, 156, 223, 301, 319, 321, 773, 791, 901, 1071, 1127-1129
- more fragments flag, IP header, *see* MF
- most specific match, 562
- Moy, J., 419, 1129
- MPFail variable, 48
- mrouted program, 339, 363, 391, 397, 401-407, 409, 411-412, 416-421, 424, 428, 433-434, 1059
- MRROUTING constant, 2
- mrt structure, 398, 419-424
  - mrt\_children member, 420
  - mrt\_leaves member, 420
  - mrt\_next member, 419-420
  - mrt\_origin member, 420
  - mrt\_originmask member, 420
  - mrt\_parent member, 420
  - mrtc\_children member, 419
  - mrtc\_leaves member, 419
  - mrtc\_origin member, 419, 423
  - mrtc\_originmask member, 419
  - mrtc\_parent member, 419
  - mrtctl structure, 401, 419-420, 423
  - mrtfind function, 423-424, 428
- MRTHASHMOD macro, 421
- MRTHASHSIZE constant, 419-421
- mrtb\_bad\_tunnel member, 399
- mrtb\_cant\_tunnel member, 399
- mrtb\_grp\_lookups member, 399
- mrtb\_grp\_misses member, 399
- mrtb\_mrt\_lookups member, 399
- mrtb\_mrt\_misses member, 399
- mrtb\_no\_route member, 399
- mrtstat structure, 398-399
- mrtstat variable, 398-399
- mrttable variable, 398, 419-420
- MSG\_CTRUNC constant, 505-506
- MSG\_DONTROUTE constant, 229, 482, 499

- MSD\_DONTWAIT constant, 482, 491, 499, 505, 507, 511, 515, 535
- MSG\_ZOR constant, 190, 482, 490–491, 498, 500, 506, 511, 520, 535, 1078
- MSG\_DON constant, 482, 505–506, 518–519, 986, 1015–1017
- MSG\_PEEK constant, 305, 507, 515, 517, 519–520, 522–523, 1016
- MSG\_TRUNC constant, 506, 523
- MSG\_WAITALL constant, 505, 507, 515, 522
- msg\_control member, 482–483
- msg\_controllen member, 482–483
- msg\_flags member, 482–483, 505–506
- msg\_iov member, 482–483
- msg\_iovlen member, 482–483, 503
- msg\_name member, 482–483
- msg\_nameLen member, 482–483, 1078
- msgndir structure, 482–483, 502–503, 505, 761, 1078
- MSIZE constant, 37, 47
- MSL (maximum segment lifetime), 813–814, 818, 820, 980–981, 991–992, 1087, 1119
- MSS (maximum segment size), 57, 797, 871, 874, 891, 896–905, 944, 1003, 1018, 1023–1024, 1117
  - option, 834, 865, 871–872, 874, 891, 897, 918, 929, 933–934, 1117
- MT\_ALIGN constant, 509–510
- MT\_CONTROL constant, 40, 509–510, 776
- MT\_DATA constant, 21, 34, 40, 509–510
- MT\_FREE constant, 40
- MT\_FTABLE constant, 40
- MT\_HEADER constant, 34, 40
- MT\_IFTABLE constant, 39–40
- MT\_IFADDR constant, 40
- MT\_OOBDATA constant, 40, 510, 519
- MT\_FCB constant, 40
- MT\_RIGHTS constant, 40
- MT\_RTABLE constant, 40
- MT\_SOCKET constant, 40
- MT\_SOMAX constant, 15, 21, 39–40, 453
- MT\_SOOPTS constant, 40
- mtod macro, 44, 46–47, 52–53, 283, 324, 630
- MTU (maximum transmission unit), 85, 92, 99, 147, 153, 232, 234, 276, 278–279, 290, 327, 612, 658, 898, 900–901, 920, 1003, 1018, 1024, 1046
  - discovery, path, 276, 279, 327, 794, 901, 1129
  - Ethernet, 92
  - loopback, 85
  - SLIP, 134, 147
- multibyte options, 248
- multicast
  - address, Ethernet, 100, 103–104, 341–342
  - address, IP, 155–156, 341
  - address mapping, IP to Ethernet, 341–342
  - backbone, see MBONE
  - forwarding, 424–433
  - group, 337
  - group cache, 399, 412, 415, 434
  - group, joining, 355–366
  - group, leaving, 366–371
  - group, local, 339
  - loopback, 400
  - network, example, 416
  - one-behind cache, 398–399, 422, 424, 434
  - output processing, 375–378, 400
  - packet, 99
  - performance, 379
  - routing, algorithms, 416–419
  - routing table, 419–424
  - scope, 348, 351, 428
  - socket option, 348
  - TTL, 348–351
  - tunnel, 398, 404–406, 427, 431–433
  - tunnel, example, 404
  - tunnel LSRR, 427, 430
- multicasting, 1127
  - ARP, 710–711
  - Ethernet, 156
  - hardware, 156, 337
  - III, 155–156, 337–380
  - RFC 1122 compliance, 1110–1111
- multihomed, 100, 155, 219, 245, 329, 337–338, 380, 741, 779–780, 1098, 1100–1101, 1118
- Mutuss, M., 5
- National Optical Astronomy Observatories, *see* NOAO
- MBFFILTER constant, 1034, 1036–1037
- Net/1, 8, 34, 562, 599, 750, 832, 844, 871, 943, 970, 989, 1053
- Net/2, 28, 34, 40, 562, 678, 680, 682, 702, 712, 718, 905, 943, 970, 1053, 1057
- Net/3, 4
- NET\_RT\_DUMP constant, 636–637, 639
- NET\_RT\_FLAGS constant, 636–637, 639
- NET\_RT\_FLIST constant, 636–637, 640
- net\_sysctl function, 202–203, 571, 635, 638, 756
- nethash function, 419–421, 424
- NETISR\_ARP constant, 687
- netstat program, 36–37, 40, 94, 97–98, 207–208, 306, 383, 398, 560–561, 563, 573, 579, 611, 680, 706, 716, 718, 757, 774, 797, 815, 940

- network
  - broadcast address, 162
  - interface, 63–94
  - interface tap, *see* NIT
  - interrupt, 73, 138, 148, 213, 469, 1078
  - loopback, 156, 181
  - mask, 157, 162
  - mask, index of, 576
  - Network File System, *see* NFS
  - Network Time Protocol, *see* NTP
  - next member, 286–289
  - NFDBITS constant, 525, 528
  - NFS (Network File System), 13, 112, 275, 441, 449, 491, 499, 507, 511, 587, 699, 785
  - NIT (network interface tap), 1129
  - NLE constant, 80
  - no operation, *see* NOP
  - NOAO (National Optical Astronomy Observatories), xxii, 28
  - noao.edu networks, 28
  - nonblocking semantics, 459
  - nonlocal source route, 1104
  - NOP (no operation), 249–250, 260, 271, 273, 282, 427, 432, 865, 933–934
  - nselcoll variable, 476, 525, 528, 534–535
  - NTOHS macro, 216
  - NTP (Network Time Protocol), 338
  - nude, IP header, 211
  - null\_sdl variable, 706
  - numvifs variable, 398, 407
- O\_ASYNC constant, 549
- O\_NONBLOCK constant, 549
- Olivier, G., 223, 1129
- O'Malley, S. W., 845, 1126
- open function, 8, 1027, 1034
- open shortest path first, *see* OSPF
- open systems interconnection, *see* OSI
- OPTBIT macro, 244
- options
  - class, IP, 249
  - echo, 866
  - IP, 247–273
    - MSS, 834, 865, 871–872, 874, 891, 897, 918, 929, 933–934, 1117
    - multibyte, 248
    - record route, 252–254
    - single-byte, 248
    - source route, 254–261
    - TCP, 864–866, 1117
    - timestamp, 261–264, 865–871, 874, 933–934, 1128
    - window scale, 865–866, 871, 874, 929, 933–934, 1128
- OPTSET macro, 242
- OPTSIZ constant, 259
- orecv function, 500
- organization, source code, 26–28
- osend function, 480
- OSI (open systems interconnection), 9–10, 23, 27, 39, 66, 69, 75–77, 86, 100–101, 105–106, 121, 123, 150, 158, 162, 185–186, 190, 194, 457, 498, 500, 508, 510, 514, 562, 570, 581, 624, 666
  - address family, 185
  - connectionless transport, 106
  - domain, 75
  - protocol family, 185
- osockaddr structure, 74–76
- OSPF (open shortest path first), 574, 1050, 1129
- out-of-band, data, 40, 505–507, 509–510, 519, 533, 855, 858, 879, 983–988, 1004, 1015–1016, 1117
- output processing
  - ICMP, 324–333
  - IP, 228–234
    - multicast, 375–378, 400
    - TCP, 851–890
    - UDP, 760–768
  - output queue, TTY, 134, 141
  - ovbcopy function, 267, 271
- P\_SELECT constant, 525, 528, 534
- p\_fd member, 13, 446–447
- p\_options structure, 347–348
- packet
  - broadcast, 99
  - header, rmbuf, 16, 32
  - IP, 210
  - multicast, 99
  - unicast, 99
- panic function, 152, 166, 460, 474, 988
- parameter problem, ICMP, 252, 257, 262, 314
- parent interface, 418–419, 1077
- Partridge, C., xxii, 60, 190, 235, 239, 351, 401, 500, 524, 716, 750, 763, 773–774, 791–792, 844, 936, 995, 1068, 1126, 1128–1131
- Pasquale, J., 234, 1128
- passing descriptors, 189, 470
- passive open, 967–969
- path MTU, discovery, 276, 279, 327, 794, 901, 1129
- Patricia tree, 562
- PAWS (protection against wrapped sequence numbers), 798, 868–869, 937, 951–954, 966, 1128
- Paxson, V., xxii, 834, 1130



- PCMTCH constant, 456
- PCB (protocol control block), 13, 347, 569, 578, 602, 618, 650, 713–753, 756–757, 760–761, 763, 768–769, 773–774, 777, 783, 785–786, 788–789, 793, 797–800, 803, 813, 815, 818, 821–822, 833–835, 885, 889, 893, 897–898, 923, 929–930, 932–933, 941, 944–946, 959–960, 966, 981, 994–995, 1009, 1012, 1017–1019, 1021–1022, 1051, 1053, 1055–1056, 1058–1060, 1062, 1090–1091, 1113, 1115
- pdev\_attach member, 78–79, 82–83, 85
- pdevinit structure, 78, 82, 85
- pdevinit variable, 64, 78–79
- performance
  - fragmentation, 291
  - header alignment, 283
  - IP checksum, 236, 239
  - low-water mark, 496
  - multicast, 379
  - send buffer size, 1018
  - SLIP, 147
  - so\_send function, 500
  - TCP, 1126
- Perlman, R., 416, 1130
- persist, timer, 818–821, 827, 835, 855, 858, 861, 878, 882, 939, 976, 1116
- Peterson, L. L., 60, 750, 845, 1126–1127
- PF\_INET constant, 10, 185–186, 197–198, 201, 449, 660, 1072
- PF\_ISO constant, 185
- PF\_LOCAL constant, 185
- PF\_OS1 constant, 185–186, 201
- PF\_ROUTE constant, 185–186, 554, 569–570, 581, 627, 645–647, 660, 662, 664, 666, 671–672
- PF\_UNIX constant, 185–186
- pfctlinput function, 123, 198, 204, 303–304, 323, 617, 743–744, 746–747, 782
- pf\_fasttimo function, 195, 796
- pf\_findproto function, 191, 196–199, 204, 449, 1052, 1072, 1080
- pf\_indtype function, 196–198, 204, 449
- pf\_slowtimo function, 195, 796
- phyint\_send function, 424–425, 429–430
- physical interface, 424
- PIM (protocol independent multicasting), 419
- ping program, 140, 272, 313, 316, 1108
- Pink, S., 239, 500, 524, 750, 763, 773–774, 791–792, 1130
- pipes, Unix, 450
- Piscitello, D. M., 9, 1130
- pkthdr structure, 34
- Plauger, P. J., 1067, 1128
- Plummer, D. C., 100, 675, 1130
- Point-to-Point Protocol, see PPP
- port
  - ephemeral, 21, 715, 719, 725, 729–730, 732, 740, 748, 751–753, 760, 813, 945, 1011–1012, 1081–1083
  - mapper, 1011
  - reserved, 732
  - unreachable, ICMP, 494
  - well-known, 719, 729, 733, 774, 814, 1011
- Portable Operating System Interface, see POSIX
- POSIX (Portable Operating System Interface), 185
- Postel, J. B., 100, 107, 156, 193, 205, 291, 301, 319, 813, 1071, 1091, 1100, 1125–1126, 1129–1130
- PPP (Point-to-Point Protocol), 23, 28, 129, 161–162, 560, 780, 846, 995, 999, 1003–1004, 1094
- frame, 129
- FR\_ADDR constant, 189, 193, 309, 385, 508, 646, 758, 1052
- FR\_ATOMIC constant, 189–190, 193, 309, 385, 490, 493, 499, 508, 515, 534, 646, 758, 1052, 1087
- FR\_CONNRQUIRED constant, 189, 531, 730, 801
- FR\_FASTHZ constant, 386
- FR\_RECV constant, 501
- FR\_RECVOOD constant, 501
- FR\_RIGHTS constant, 189
- FR\_SEND constant, 480
- FR\_SENDOOD constant, 480
- FR\_SHUTDOWN constant, 469
- FR\_SLOWHZ constant, 824, 834
- FR\_WANTRCVD constant, 189, 801
- prctlinput member, 190, 193, 198, 303–304, 309, 315, 385, 646, 744, 758, 782, 801, 1052, 1107–1108
- prctloutput member, 190, 193, 309, 385, 436, 538, 540–542, 646, 758, 801, 1052
- pr\_domain member, 188, 193, 309, 385, 646, 758, 801, 1052
- pr\_drain member, 191, 193, 309, 385, 646, 758, 801, 892, 1052
- pr\_fasttimo member, 190, 193, 196, 309, 385, 646, 758, 801, 821, 1052
- pr\_flags member, 188–189, 193, 309, 385, 531, 646, 758, 801, 1052
- pr\_inif member, 190, 193–194, 309, 385, 646, 758, 801, 1052
- pr\_input member, 190, 193, 220, 286, 309–310, 385, 391, 646, 758, 769, 801, 923, 1052–1053, 1091, 1100, 1102
- pr\_output member, 190, 193, 228, 309, 385, 646, 670, 758, 801, 1052



- DF\_PROTOCOL member, 188, 191, 193, 196, 200, 203, 309, 385, 646, 673, 758, 801, 1052
  - df\_lowtime member, 190, 193, 196, 309, 385, 646, 758, 801, 822, 826, 1052
  - df\_syncall member, 191, 193, 202-203, 309, 385, 646, 758, 801, 1052
  - df\_type member, 188-189, 191, 193, 309, 385, 437, 646, 758, 801, 1052
  - df\_usrreq member, 164, 190, 193, 309, 385, 436, 450, 455, 461, 474, 480, 489-490, 499, 501, 540, 552, 646, 666, 758, 801, 1007, 1052, 1062
- PRC\_HOSTDEAD constant, 316, 744
- PRC\_IFDOWN constant, 123, 316
- PRC\_MSGSIZE constant, 302, 316
- PRC\_PARAMPROB constant, 302, 316
- PRC\_QUEMCH constant, 302, 316, 783, 904
- PRC\_QUEMCH2 constant, 316
- PRC\_REDIRECT\_HOST constant, 302, 316
- PRC\_REDIRECT\_NET constant, 316
- PRC\_REDIRECT\_TOHOST constant, 316
- PRC\_REDIRECT\_TOINET constant, 316
- PRC\_ROUTEDEAD constant, 316
- PRC\_TIMXCEED\_INTRANS constant, 302, 316
- PRC\_TIMXCEED\_REASS constant, 302, 316
- PRC\_UNREACH\_HOST constant, 302, 316
- PRC\_UNREACH\_NET constant, 302, 316
- PRC\_UNREACH\_PORT constant, 302, 316
- PRC\_UNREACH\_PROTOCOL constant, 302, 316
- PRC\_UNREACH\_SRCFAIL constant, 302, 316
- PRCO\_GETOPT constant, 240, 243-244, 401, 546, 1023-1024
- PRCO\_SETOPT constant, 240, 242-243, 401, 540, 1023-1024
- prev member, 286-289
- principle, robustness, 857
- printf function, 7-8, 81, 797
- proc structure, 11, 446, 531
- promiscuous, mode, 101, 104, 125, 1033, 1035, 1070, 1092
- protection against wrapped sequence numbers, see PAWS
- protocol
  - control block, see PCB
  - entry points, 190
  - family, 182, 185, 189-190, 196, 202
  - family, Internet, 185, 202-203, 361
  - family, OSI, 185
  - family, routing, 185
  - family, Unix, 185
  - independent multicasting, see PIM
  - layer, 10
  - protocol structure, 186-196, 198-200, 202, 204, 220, 228, 309, 384-385, 440, 446-447, 449, 500, 647, 744, 795, 801, 1051-1052
  - ICMP, 309-310
  - IGMP, 384-385
  - IP, 186, 191-196
  - raw IP, 1051-1053
  - routing, 646
  - TCP, 801
  - UDP, 758
  - prototypes, ANSI C function, 41
  - proxy ARP, 688, 703-704
  - PRU\_ABORT constant, 450, 471, 669-671, 786, 788-789, 892, 1014, 1059
  - PRU\_ACCEPT constant, 450, 461, 787, 1012-1013
  - PRU\_ATTACH constant, 450, 462-463, 477, 647, 665-667, 671, 717, 785, 822, 833, 932, 966, 1009-1010, 1018, 1053, 1058
  - PRU\_BIND constant, 450, 454-455, 666, 786, 1010, 1059
  - PRU\_CONNECT constant, 450, 465, 468, 666, 787-788, 851, 871, 873, 884, 919, 962, 1059-1060, 1083
  - PRU\_CONNECT2 constant, 450, 668-669, 787, 1012, 1061
  - PRU\_CONTROL constant, 164, 554, 666, 785, 1007
  - PRU\_DETACH constant, 450, 473, 666-667, 669, 671, 719, 786, 788, 1010, 1059
  - PRU\_DISCONNECT constant, 450, 668-669, 671, 787-788, 897, 1010, 1012, 1019, 1059, 1083
  - PRU\_LISTEN constant, 450, 455, 730, 786, 1011
  - PRU\_PERRADDR constant, 450, 556, 670, 741, 789, 1017, 1062
  - PRU\_RCVD constant, 450, 514, 523, 790, 851-852, 863, 1013, 1120
  - PRU\_RCVOOB constant, 450, 513, 790, 985-986, 1015
  - PRU\_SEND constant, 450, 499, 648, 669, 761, 784, 788-789, 851, 884, 920, 1013, 1062, 1087
  - PRU\_SENDOOB constant, 450, 499, 851, 861, 884, 1016
  - PRU\_SENSE constant, 669-670, 789, 1014-1015, 1062
  - PRU\_SHUTDOWN constant, 450, 469-470, 668-669, 788, 851, 884, 1013, 1021, 1061
  - PRU\_SLOWTIME constant, 822, 824, 1017
  - PRU\_SOCKADDR constant, 450, 554, 670, 741, 789, 1017, 1062
  - prune, 418
  - ps program, 456
  - pseudo-device, 78, 83
    - loopback, 64
    - SLIP, 64, 82

- pseudo-header
    - TCP, 880, 885, 926, 995
    - UDP, 764-768
  - PSH (push flag, TCP header), 803, 875, 889, 937, 959, 1000-1001, 1086, 1115-1116
  - pure ACK, 831, 851, 937
  - push flag, TCP header, *see* PSH
  - putc function, 146
  - putmsg function, 8
- Quarterman, J. S., 24, 33-34, 94, 129, 445, 457, 470, 562, 1097, 1129
- queue
  - ARP, input, 97
  - CLNP, input, 97
  - IP, input, 97, 106
- queue of mbufs, 34, 39
- quiet time, 812-814
- 
- R\_malloc macro, 585
- radix node data structures, 573-578
- radix\_mask structure, 577-578, 587, 589-591, 595, 597
- radix\_node structure, 568-569, 573, 575, 577-579, 581, 586-587, 590-591, 595, 597, 600, 602, 609, 612, 641, 1079
- radix\_node\_head structure, 567-569, 573, 575-576, 586-587, 592
- Rago, S. A., xxii, 5, 435, 1130
- Ramsey, N. F., 106, 1127
- RARP (Reverse Address Resolution Protocol), 100, 106, 321, 686
- Ravi, C., xxii
- raw IP, 183, 191, 197, 230, 240, 276, 301, 304-305, 312-313, 384, 391, 440, 451, 477, 541, 1049-1065, 1071, 1102, 1108
  - inetsw variable, 191, 193, 197, 199, 204, 1052, 1072
  - protocols structure, 1051-1053
- raw protocol, default, 191
- raw sockets, ICMP redirect and, 746-748
- raw\_attach function, 667, 671-672
- raw\_ctlinput function, 646
- raw\_detach function, 667, 671-672
- raw\_disconnect function, 669, 671-672
- raw\_init function, 571, 646-648
- raw\_input function, 312, 571, 621, 624-629, 632, 645-646, 650, 659-660, 662-664, 671-672
- raw\_recvspace variable, 572
- raw\_sendspace variable, 572
- raw\_usrreq function, 571, 666-670, 672
- 
- rawcb structure, 440, 647, 665, 672
- rawcb variable, 572, 647, 671
- rawingcb variable, 1051, 1053, 1058-1059
- rch\_raddr member, 647, 664, 669
- rch\_raddr member, 647, 664
- rch\_proto member, 662
- rcmd function, 732
- rcv\_adv member, 808, 860, 863-864, 878, 884, 941, 948
- rcv\_ext member, 808, 830, 860, 863-864, 867, 869, 871, 878, 884, 907-909, 915-916, 937, 940-941, 948, 954, 959, 985, 987-990, 1088-1090
- rcv\_scale member, 866, 949, 969
- rcv\_up member, 951, 985
- rcv\_wnd member, 808, 951, 959
- rcvif member, 16, 20, 33, 53, 152
- RDP (Reliable Datagram Protocol), 189-190, 490, 716, 1130
- read function, 10, 13, 129, 435, 439, 445-447, 475, 500-501, 1025, 1040, 1043
- readv function, 13, 445-446, 475, 500-501
- REASS\_MBUF macro, 913
- reassembly
  - IP, 48-50, 219, 275-277, 283-300
  - TCP, 50-51, 906-916
  - TTL, 298
- receive sequence space, TCP, 808
- record boundaries, 189
- record route option, 252-254
- recovery, fast, 970-974, 1120, 1128
- recv function, 445-446, 500-501, 503, 986
- recvfrom function, 7-8, 10, 19, 21-22, 443, 445-446, 457, 474-475, 500-501, 503, 748, 792, 1078
- recvit function, 501, 503-505, 535, 1078, 1083
- recvmsg function, 21, 435, 443, 445-446, 475, 500-503, 505-506, 511, 523, 776, 781, 1078, 1083, 1114
- redirect, ICMP, 221, 223-228, 252, 321
- reference counts
  - cluster, 56-60
  - control message, 470
  - ether\_multi structure, 343, 346, 364, 369
  - ifaddr structure, 74, 177
  - in\_multi structure, 359-360, 368, 386, 395
  - routing table, 606-607
- refinements
  - TCP implementation, 994-995
  - UDP implementation, 791-792
- reliable
  - protocol buffers, 490
  - protocols, 189

- Reliable Datagram Protocol, *see* RDP
- remote procedure call, *see* RPC
- remote terminal protocol, *see* Telnet
- remque function, 55, 291–292, 298, 671, 683, 710, 719
- reply, ICMP, 234
- Request for Comment, *see* RFC
- request\_r\_scale member, 873
- requested\_s\_scale member, 934
- reserved port, 732
- reset flag, TCP header, *see* RST
- reset segment, TCP, *see* RST
- reset\_d function, 1035, 1039
- resynchronization ACK, 960, 1088
- retransmission
  - ambiguity problem, 976
  - time out, *see* RTO
  - timer, 749, 817, 819–822, 828, 831–833, 835, 841, 843–844, 846, 859, 861, 880, 882, 889, 918, 920, 939, 970–972, 975–976, 978, 1012, 1020, 1086
- retransmit, fast, 908, 970–974, 1120, 1128
- Reverse Address Resolution Protocol, *see* RARP
- reverse path, 417
  - broadcasting, *see* RPB
  - multicasting, *see* RPM
- reversed source route, 240, 255, 258–259, 261, 332, 1103, 1108–1109
- Reynolds, J. K., 100, 107, 193, 1091, 1130
- RFC (Request for Comment), 791, 1130
  - 792, 1130
  - 793, 1130
  - 795, 1130
  - 817, 1126
  - 826, 1130
  - 894, 1127
  - 903, 1127
  - 950, 1129
  - 951, 1126
  - 1009, 1126
  - 1042, 1130
  - 1055, 1130
  - 1071, 1126
  - 1072, 1128
  - 1075, 1131
  - 1112, 1127
  - 1122, 1125
  - 1122 compliance, ARP, 1113
  - 1122 compliance, ICMP, 1105–1110
  - 1122 compliance, IGMP, 1111
  - 1122 compliance, IP, 1098–1105
  - 1122 compliance, link layer, 1097–1098
  - 1122 compliance, multicasting, 1110–1111
  - 1122 compliance, routing, 1111–1113
  - 1122 compliance, TCP, 1115–1123
  - 1122 compliance, UDP, 1113–1115
  - 1123, 1125
  - 1127, 1125
  - 1141, 1129
  - 1144, 1127
  - 1151, 1130
  - 1190, 1131
  - 1191, 1129
  - 1213, 1129
  - 1256, 1127
  - 1323, 1128
  - 1332, 1129
  - 1337, 1126
  - 1349, 1125
  - 1462, 1128
  - 1519, 1127
  - 1541, 1127
  - 1546, 1130
  - 1548, 1130
  - 1624, 1130
  - 1644, 1126
  - 1700, 1130
  - 1716, 1125
  - Gateway Requirements, 1126
  - Host Requirements, 1125
  - how to obtain, 1094
  - Router Requirements, 1125
- Rijisinghani, A., 235, 1129–1130
- RIP (Routing Information Protocol), 291, 574, 1050
  - rip\_output function, 193, 240, 309, 385, 401–402, 412, 538, 541, 1050, 1052, 1063–1065
  - rip\_init function, 193–194, 1050, 1052–1053
  - rip\_input function, 193, 303–304, 310, 312–313, 316, 319, 321, 323, 335, 382, 391, 412, 419, 1050, 1052–1056, 1058–1060, 1065, 1075, 1091, 1105, 1107–1109
  - rip\_output function, 193, 309, 382, 385, 1050, 1052–1053, 1056–1058, 1062, 1065
  - rip\_recvspace variable, 1051
  - rip\_sendspace variable, 1051
  - rip\_usrreq function, 193, 309, 385, 451, 1050, 1052, 1058–1062
  - ripsrc variable, 1051, 1055
- Rlogin, 140, 732, 858, 995, 1002, 1084
- rm\_mask member, 577–578, 597
- rm\_nklist member, 578
- rm\_expire member, 581, 658, 678, 683, 694, 706, 711–712, 894
- rmx\_hopcount member, 658, 894
- rmx\_locks member, 581, 658, 893

- rmx\_sbu member, 658, 894
- rmx\_pssent member, 658, 893
- rmx\_rcvpipe member, 658, 894
- rmx\_rts member, 658, 894, 899, 921
- rmx\_rttvar member, 658, 894, 899
- rmx\_sendpipe member, 658, 894, 896, 901
- rmx\_ssthresh member, 658, 894, 896, 903
- rs\_addmask function, 653
- rs\_addroute function, 571, 575, 611
- rs\_b member, 576-577, 587-588, 591-592, 597, 1079
- rs\_mask member, 576, 587, 592
- rs\_delete function, 571, 575, 608
- rs\_dupedkey member, 577, 587-588, 593, 597
- rs\_flags member, 576, 587
- rs\_init function, 571, 584-587
- rs\_lnichead function, 192, 569, 571, 575, 581, 584-589, 605
- rs\_key member, 577, 588-589, 597, 610, 612-613
- rs\_l member, 577
- rs\_mask member, 577, 588, 597, 654
- rs\_match function, 571, 575, 591-599, 602, 604, 1111
- rs\_mfreeslist variable, 572, 578
- rs\_mnlist member, 576, 578, 587, 595, 597
- rs\_off member, 576-577, 587, 592
- rs\_ones variable, 572, 585, 587, 589
- rs\_p member, 576, 587
- rs\_r member, 577
- rs\_search function, 571, 596-597, 599
- rs\_walktree function, 571, 575, 638, 640
- rs\_xprios variable, 572, 585, 587-589, 619
- RNF\_ACTIVE constant, 576, 603
- RNF\_NORMAL constant, 576
- RNF\_ROOT constant, 576, 587-589, 593, 602, 605
- rnh\_addsadr member, 575, 611
- rnh\_addpkt member, 575
- rnh\_addrsize member, 575
- rnh\_deladdr member, 575, 608
- rnh\_delpkt member, 575
- rnh\_matchsadr member, 575, 591, 602
- rnh\_matchpkt member, 575
- rnh\_nodes member, 587-588
- rnh\_pktsize member, 575
- rnh\_treetop member, 575, 586
- rnh\_walktree member, 575, 640
- rod\_at member, 221, 223, 254, 578, 1090
- rod\_at member, 221, 578-579, 602, 1090-1091
- robustness principle, 857
- Romkey, J. L., 129, 144, 1130
- Rose, M. T., 9, 64, 1129-1130
- ROUTE\_BUFFERS macro, 1033, 1043
- round-trip time, see RTT
- ROUNDUP macro, 612, 632
- route
  - cached, 746-747, 750, 768, 843, 887, 894, 898, 1111
  - characteristics, 893-896
  - default, 181
  - direct, 561, 621, 706
  - held, 606, 659
  - indirect, 561, 569, 580, 608, 615, 706
  - selection, IP, 230-232
- route program, 560-561, 569, 571-572, 601, 606, 612, 650, 893-894, 1112
- route structure, 220-221, 223, 231, 234, 358-359, 568-569, 578-579, 599, 602, 1090-1091
- route\_cb variable, 572, 624, 666
- route\_dst variable, 572, 626, 664
- route\_init function, 571, 581-584, 646
- route\_output function, 571, 606-607, 632-633, 645-646, 648-661, 670, 672-673, 679, 710
- route\_proto variable, 572, 626-627, 660, 664
- route\_src variable, 572, 626, 647, 664, 666, 670
- route\_xrefq function, 571, 646-647, 664-666, 669-670, 672
- routed program, 559, 571-572, 601, 637, 644, 1111
- routedomain variable, 187, 195, 646
- router, discovery, ICMP, 1127
  - solicitation, ICMP, 339
  - vs. host, 157
- Router Requirements RFC, 1125
- routeosw variable, 195, 646, 673
- Routhier, S., 223, 750, 1128
- routing
  - address family, 185
  - control block, 647
  - domain, 67, 75, 437, 539, 554, 569-570, 572, 581, 584, 624, 632, 645-673
  - domain initialization, 646
  - domain structure, 646
  - IP multicast, 397-434
  - mechanism, 559
  - messages, 601-644
  - policy, 559
  - protocol family, 185
  - protoosw structure, 646
  - requests, 601-644
  - RFC 1122 compliance, 1111-1113
  - socket, 569, 645-673
  - structures, 578-581
  - table, 560-569
  - table, ARP, 675-678
  - table internal node, 564



- table leaf node, 564
- table reference counts, 606–607
- Routing Information Protocol, *see* RIP
- RPB (reverse path broadcasting), 417, 419, 434
- RPC (remote procedure call), 729, 1011
- RPM (reverse path multicasting), 418–419, 434
  - example, 418
- rsh program, 732
- RST (reset flag, TCP header), 234, 800, 803, 820, 843, 868, 886–887, 889, 892, 930, 948, 957, 963–966, 992–994, 1014, 1019, 1025, 1086–1088, 1090, 1118–1119, 1126
  - storm, 993
- rt\_addrinfo structure, 621, 623, 625, 627, 629–630, 632, 649–650, 660, 750
- rt\_expire member, 695, 700–701, 1081
- rt\_flags member, 573, 579–580, 609, 641, 677, 681
- rt\_gateway member, 574, 579, 609, 612–614, 619, 621, 629, 636, 677, 681, 692, 695, 698, 702, 706–707, 709, 1080
- rt\_genmask member, 609, 611, 636, 658
- rt\_gwroute member, 580, 608–609, 614–615, 1091
- rt\_ifa member, 580, 609, 617, 619, 656, 677
- rt\_ifmsg function, 125, 571, 627–628
- rt\_ifp member, 224, 580, 656, 677
- rt\_key member, 574, 609, 612, 629, 636, 681
- rt\_llinfo member, 580, 677, 683, 707, 710
- rt\_mask member, 573–574, 629
- rt\_maskedcopy function, 610, 615, 1071
- rt\_metrics member, 661, 683, 711
- rt\_metrics structure, 581, 653, 658, 661, 678, 893–894
- rt\_mismsg function, 571, 621, 624–627, 750
- rt\_msg1 function, 625–627, 629–632
- rt\_msg2 function, 632–635, 638, 640–643, 656
- rt\_msghdr structure, 569–570, 621, 629–630, 632, 650–651, 656, 661, 1080
- rt\_netmask member, 636
- rt\_newaddressg function, 571, 616–617, 628–630
- rt\_nodes member, 579
- rt\_refcnt member, 580, 604, 606, 608, 611, 702
- rt\_rmx member, 658
- rt\_setgate function, 606, 609, 612–615, 621, 656, 706
- rt\_setmetrics function, 656, 658, 661
- rt\_tables variable, 567–568, 572, 586–587, 639
- rt\_use member, 580
- rt\_xadinfo function, 650–651, 660–661
- RTA\_AUTHOR constant, 623
- RTA\_BRD constant, 623
- RTA\_DST constant, 623
- RTA\_GATEWAY constant, 623
- RTA\_GENMASK constant, 623
- RTA\_IFA constant, 623, 656
- RTA\_IFP constant, 623, 656
- RTA\_NETMASK constant, 623
- rtable\_init function, 581–584
- rtalloc function, 223, 232, 254, 358–359, 371–372, 578–579, 591, 601–604, 606, 618–619, 738, 752, 898, 1090
- rtalloc1 function, 571, 591, 593, 601–604, 606–607, 609, 615, 619, 623–624, 654–655, 659, 679, 701–704, 707, 710, 1091
- RTAX\_AUTHOR constant, 623
- RTAX\_BRD constant, 623
- RTAX\_DST constant, 623, 653
- RTAX\_GATEWAY constant, 623
- RTAX\_GENMASK constant, 623
- RTAX\_IFA constant, 623
- RTAX\_IFP constant, 623
- RTAX\_MAX constant, 623, 661
- RTAX\_NETMASK constant, 623
- rtentry structure, 221, 568–569, 575, 578–581, 602–604, 606, 608–609, 611–612, 616, 641, 677, 1079
- RTF\_ANNOUNCE constant, 703, 707
- RTF\_BLACKHOLE constant, 152, 579–580
- RTF\_CLONING constant, 169, 580, 603–604, 609, 612, 615, 653, 703–704, 706–707
- RTF\_DONE constant, 579–580, 619, 650, 659
- RTF\_DYNAMIC constant, 224, 573, 580, 621, 644, 750, 1080
- RTF\_GATEWAY constant, 573, 579–580, 619, 621, 702, 706
- RTF\_HOST constant, 171, 580, 609, 619, 706
- RTF\_LLINFO constant, 580, 637, 641, 680, 702, 707, 710–711
- RTF\_MASK constant, 579–580
- RTF\_MODIFIED constant, 224, 573, 580, 621, 644, 1080
- RTF\_PROTO1 constant, 580
- RTF\_PROTO2 constant, 580
- RTF\_REJECT constant, 109, 152, 580, 692, 695, 699–701
- RTF\_STATIC constant, 580
- RTF\_UP constant, 171, 580, 606–609, 704
- RTF\_XRESOLVE constant, 580, 604
- rtree function, 576, 604–607, 609, 616, 618, 621, 653, 719, 746
- RTFREE macro, 234, 604–608, 614
- rti\_addr member, 623–624, 626, 632, 634, 650–651



- `rtt_info` member, 623-624, 626, 629, 632, 634, 641, 653, 656, 660-661
- `rttlimit` function, 171-172, 371, 601, 606-607, 615-617, 628, 677, 679, 704, 706
- `rttlimit` function, 552, 554
- `RTM_ADD` constant, 570, 606-607, 609, 611, 615, 617, 628-629, 645, 649-650, 653, 657, 661, 704, 706, 711
- `RTM_CHANGE` constant, 570, 645, 654, 656
- `RTM_DELETE` constant, 570, 621, 629
- `RTM_DELETE` constant, 570, 607, 615, 617, 628-629, 645, 653, 656, 704, 706, 709, 750
- `RTM_GET` constant, 570, 632, 636, 641, 645, 654-656, 660, 704, 710
- `RTM_INFO` constant, 570, 621, 636, 643
- `RTM_LOCK` constant, 570, 645, 654, 656, 658
- `RTM_LOCKING` constant, 570, 750, 1112-1113
- `RTM_MISS` constant, 570, 602, 604
- `RTM_NEWADDR` constant, 570, 621, 629, 636, 643
- `RTM_REDIRECT` constant, 570, 1049
- `RTM_RESOLVE` constant, 570, 603-604, 607, 609, 611, 704, 706-707, 709
- `RTM_RTTUNIT` constant, 894
- `rtm_addr` member, 621-622, 626, 650, 660-661
- `rtm_errno` member, 626, 650, 659, 1080
- `rtm_flags` member, 626, 650, 656
- `rtm_inits` member, 658, 661
- `rtm_msglen` member, 637, 651
- `rtm_pid` member, 651
- `rtm_rmx` member, 658, 661
- `rtm_type` member, 569
- `RTO` (retransmission time out), 800, 831-832, 834, 836, 840-841, 843, 900, 1116, 1119, 1121
- `rtrredirect` function, 323, 571, 607, 617-621, 623-624, 1107
- `rtrrequest` function, 371, 601, 603-612, 615-617, 621, 653, 661, 679, 696, 703-704, 706-707, 710-711, 750, 1080-1081
- `rtw_hdrredirect` member, 573
- `rtw_dynamic` member, 573
- `rtw_newgateway` member, 573
- `rtw_unreach` member, 573, 602
- `rtw_wildcard` member, 573
- `statst` structure, 573
- `statst` variable, 572-573
- `RTT` (round-trip time), 612, 658, 797-798, 820, 823, 831-834, 836-837, 839-840, 842-848, 852, 866, 868-871, 889, 894, 899, 917-919, 939, 950, 975-976, 1085, 1089, 1119
- `rtthresh` variable, 572, 605, 608
- `RTV_EXPIRE` constant, 658
- `RTV_MOPCOUNT` constant, 658
- `RTV_MTU` constant, 658
- `RTV_RPIPE` constant, 658
- `RTV_RTF` constant, 658, 899
- `RTV_RTVAR` constant, 658
- `RTV_SPIPE` constant, 658
- `RTV_SSTHRESH` constant, 658
- `rvmt` packet, 103
- `rvmd` program, 571-572
- `s_addr` member, 160, 266, 410, 1071
- `s_host` member, 1071
- `s_inp` member, 1071
- `s_inproc` member, 1071
- `s_lh` member, 1071
- `s_net` member, 1071
- `sa_alen` member, 121
- `sa_data` member, 75-76, 160
- `sa_family` member, 75-76, 94, 109, 153, 160, 212, 602, 608, 627, 684, 686, 1081
- `sa_len` member, 75-76, 94, 117, 160, 453, 661, 686, 1077, 1081
- `sa_nlen` member, 121
- `sa_slten` member, 121
- `SACK` (selective acknowledgment), 866, 1128
- `Salus, P. H.`, xxi, 5, 1130
- `save_rte` function, 257-259, 261, 272, 932
- `SB_ASYNC` constant, 477, 550
- `SB_LOCK` constant, 477
- `SB_NOINTR` constant, 469, 477
- `SB_NOTIFY` constant, 477, 939
- `SB_SEL` constant, 477
- `SB_WAIT` constant, 477
- `SB_WANT` constant, 477
- `sb_loc` member, 476, 478, 490, 515, 530-531, 535, 550, 861, 985, 1078
- `sb_flags` member, 477, 550
- `sb_lowat` member, 463, 476-479, 489-491, 515, 531, 535, 539, 878, 894, 1012, 1078
- `sb_lowat` member, 476-479, 489-490, 508, 515, 530-531, 539
- `sb_max` member, 477, 903
- `sb_nmb` member, 56, 470, 476
- `sb_nmbent` member, 476, 478
- `sb_nmbmax` member, 476-479, 495
- `sb_sel` member, 477, 531
- `sb_timeout` member, 477, 496, 516, 539, 543-544, 548
- `sballot` macro, 478
- `sbandpend` function, 479, 508, 987, 1014, 1016, 1078
- `sbandpendaddr` function, 479, 508, 625, 664, 776-777, 1055, 1078

- sbappendcontrol function, 479, 509, 1078
- sbappendrecord function, 479, 508
- sbcocompress function, 479
- sbdrop function, 479, 978
- sbdroprecord function, 479
- sbflush function, 479, 1020
- sbfree macro, 478
- sbinsertflood function, 479, 509
- sblock macro, 469, 478, 491, 511
- SBLOCKWAIT constant, 491
- sbmax variable, 476
- sbrelease function, 470, 473, 479, 1078
- sbreserve function, 479, 543, 903, 1018
- sbpace macro, 478, 495, 531, 535, 1078
- sbunlock macro, 478
- sbwait function, 478, 496, 515–516, 522
- SC\_AUTOCOMP constant, 83–84
- SC\_COMPRESS constant, 83
- SC\_ERROR constant, 83, 135–137, 153, 1070
- SC\_NOICMP constant, 83, 140
- sc\_ac member, 80
- sc\_addr member, 80–81
- sc\_bpf member, 83, 1030
- sc\_buf member, 83, 131, 133
- sc\_comp member, 83
- sc\_ep member, 83, 131, 133
- sc\_escape member, 83, 136, 138
- sc\_fastq member, 83–84, 131, 140–141, 148
- sc\_it member, 80–81, 83–84, 1040
- sc\_imp member, 83, 131, 133
- sc\_softc structure, 83
- sc\_ttyp member, 83, 131–132
- scatter, 481–482, 486
- Schaller, D., xxii
- schednetisr function, 106–107, 153, 212
- scheduler function, 79
- Schmidt, D. C., xxii
- SCM\_RIGHTS constant, 517
- sd\_lalen member, 87, 91, 677, 692, 695, 698–699, 706, 711, 1070, 1080
- sd\_data member, 87–88, 90, 97
- sd\_family member, 86, 90–91, 706, 1070
- sd\_index member, 86, 90–91, 121, 706, 1070
- sd\_len member, 86, 90–91, 706, 1070
- sd\_nlen member, 87, 90–91, 1070
- sd\_nlen member, 87, 91, 1070
- sd\_slalen member, 87, 91, 1070
- sd\_sltype member, 86, 90–91, 121, 706, 1070
- Sedgewick, R., 562, 1130
- seq\_next member, 834, 906, 908–909, 940
- seq\_prev member, 834, 906, 909
- segments, cached, 972
- sel\_pid member, 534
- select function, 13, 22, 445–447, 463–464, 475–477, 524–525, 528, 531–532, 534–535, 749, 783, 792, 905, 939, 969, 1033, 1043, 1079
- selective acknowledgment, *see* SACK
- self-connect, 956, 960–962, 967, 1005, 1090
- selinfo structure, 477, 531–532, 534, 1033
- selrecord function, 529, 531, 534
- selscan function, 525, 528–529, 534, 1079
- selwait variable, 476, 528, 532, 534
- selwakeop function, 532–535
- send function, 8, 229, 445–446, 480–483, 494, 760–761, 985, 1056
- send sequence space, TCP, 808
- sendit function, 480, 483–485, 487–489, 494, 496, 761
- sendmsg function, 229, 435, 443, 445–446, 453, 475, 480–485, 502, 666, 669, 760–761, 1007, 1056
- sendto function, 7–8, 10, 14–16, 19, 29, 39, 41, 93, 229, 445–446, 453, 475, 480–481, 485, 579, 669, 729, 735, 738, 740–741, 748, 751–752, 760–761, 763, 774, 789, 793, 1056, 1083, 1091, 1115
- SEQ\_GEQ macro, 870
- SEQ\_GT macro, 810, 986
- SEQ\_LEQ macro, 868
- SEQ\_LT macro, 810, 868–869
- sequence numbers, TCP, 807–812
- sequence space
  - TCP receive, 808
  - TCP send, 808
- Sequenced Packet Protocol, *see* SPP
- Serial Line Internet Protocol, *see* SLIP
- setsockopt function, 239–244, 401, 412, 437, 445–446, 463, 537–539, 557, 720, 768, 785, 1007, 1022–1023, 1063, 1065
- shrink, window, 856–857, 878, 884, 1116
- SHRT\_MAX constant, 544
- shutdown function, 445–446, 468–470, 514, 650, 669, 788, 807, 818, 825, 1013, 1061, 1119
- shutdown\_args structure, 469
- SI\_COLL constant, 531–532
- si\_pid member, 532, 534
- STOIO signal, 22, 439–440, 478, 550, 552, 939
- SIGPIPE signal, 489
- SIGURG signal, 550, 552
- silly window syndrome, *see* SWS
- Simple Mail Transfer Protocol, *see* SMTP
- Simple Network Management Protocol, *see* SNMP
- Simpson, W. A., 129, 1130

- simultaneous
  - close, 807, 991, 1005
  - open, 948, 956, 960–962, 967–969
- SIN\_PROXY constant, 694, 701–703, 707
- sin\_addr member, 160, 183, 702, 742, 1082
- sin\_family member, 160, 701
- sin\_len member, 160, 183
- sin\_other member, 701–703, 707
- sin\_port member, 160, 742, 1055, 1082
- sin\_srcaddr member, 701
- sin\_tos member, 701
- sin\_zero member, 160, 732
- single-byte options, 248
- SIODADDRMULTI constant, 344, 356, 361–363, 369, 379, 409, 451
- SIODCAFADDR constant, 162, 170, 174–177
- SIODCATHACK constant, 506–507, 550, 552, 1117
- SIODDELMULTI constant, 344, 356, 361–363, 369, 380, 410, 434, 451, 1077
- SIODGIFADDR constant, 74, 162, 174–177
- SIODGIFADDR constant, 162, 173–174
- SIODGIFBRDADDR constant, 162, 173–174
- SIODGIFCODE constant, 114, 116, 120–121, 637, 1083
- SIODGIFDSTADDR constant, 162, 173–174
- SIODGIFFLAGS constant, 67, 114, 121, 1111
- SIODGIFMETRIC constant, 114, 121
- SIODGIFNETMASK constant, 162, 173–174
- SIODGIFPR constant, 440, 550, 552
- SIODGIFADDR constant, 162, 166–171, 177–180, 683
- SIODSIFBRDADDR constant, 162, 173–174
- SIODSIFDSTADDR constant, 162, 166, 172, 179
- SIODSIFFLAGS constant, 67, 114, 121–125, 178
- SIODSIFMETRIC constant, 114, 121, 123
- SIODSIFNETMASK constant, 162, 166, 170, 172
- SIODSPORT constant, 440, 550, 552
- Stroica, D., 1126
- Skibo, T., 1067
- Skłower, K., xix, 562, 599, 611, 1130
- sl\_bcam function, 130, 138
- sl\_compress\_init function, 133
- sl\_compress\_tcp function, 141, 997
- sl\_flags member, 531
- sl\_gid member, 531
- sl\_softc structure, 77, 83–84, 86, 120, 128, 130–133, 135, 140, 149, 159, 1030, 1069, 1071
- sl\_softc variable, 77, 86, 120, 128, 131, 159, 1030
- sl\_uncompress\_tcp function, 138
- slattach function, 78, 80, 82–84, 86, 94, 129–130, 132–133, 148–149
- slattach program, 84, 130
- SLBUFSIZE constant, 134
- slclose function, 130, 148
- slinit function, 130, 132–133
- slinput function, 130–131, 134–138, 153, 213, 1070
- SLIOCGUNIT constant, 149–150
- slioctl function, 96, 115, 124, 127, 130, 149, 163, 168, 177, 179, 361
- SLIP (Serial Line Internet Protocol), 23, 25, 27–28, 30, 63, 66–69, 71, 76, 78–80, 82–84, 86–87, 94–96, 98, 120–121, 124, 128–150, 158, 212, 219, 246, 283, 290, 337, 341, 361, 380, 451, 561, 995, 999, 1002–1004, 1027, 1030, 1032, 1069, 1127
- BPF, 104
  - cluster, 131
  - device driver, 63–64, 69, 83, 128–150, 161, 179, 1029
  - discarding line noise, 144
  - encapsulation, 128–129
  - END character, 129, 144
  - ESC character, 129, 144
  - frame, 83, 128, 131, 133–137, 143–144, 146–147, 1070
  - initialization, 82–84
  - line discipline, 129–132, 134, 149
  - MTU, 134, 147
  - packet, BPF format, 138
  - performance, 147
  - pseudo-device, 64, 82
  - sockaddr\_sl structure, 1070
  - TOS queueing, 140, 147
- SLIP\_HDRLEN constant, 134, 138
- SLIP\_HIWAT constant, 134, 141, 147–148
- SLIPDIR\_IN constant, 138
- SLIPDIR\_OUT constant, 141
- SLIPDISC constant, 129, 132
- SLMAX constant, 134
- SLMTU constant, 83, 134, 147
- slopen function, 130, 132–133, 149, 451
- sloutput function, 96, 127, 130–131, 139–141
- slow start, 844, 852, 882, 896, 903, 906, 920, 939, 970, 972, 974–975, 977, 1119, 1122, 1127
- slstart function, 130–131, 141–148
- sltioc1 function, 115, 130, 149–150
- SLX\_CBRD constant, 138
- SLX\_DIR constant, 138
- SMTP (Simple Mail Transfer Protocol), 140, 1125
- end\_cwnd member, 835, 844, 854, 903, 939, 977
- snd\_max member, 808, 852, 859, 874, 877–878, 880, 882, 937–939, 948, 968, 970, 975, 979

- and\_nxt member, 808, 844, 857, 859, 871, 874, 877-880, 882, 937, 948, 972-973, 978-979, 1085
- and\_scale member, 866, 949, 969
- and\_sathresh member, 835, 896, 903, 906
- and\_una member, 808, 830, 844, 846, 852, 854, 857, 859, 876, 880, 882, 938-939, 948-949, 968, 970, 975, 978-979, 982, 1088
- and\_up member, 802, 861, 878-879, 1016
- and\_wll member, 951, 969, 981-983
- and\_wll2 member, 951, 981-983
- and\_wnd member, 808, 844, 854, 937, 939, 951, 978, 981, 983
- SNMP (Simple Network Management Protocol), 2-3, 64-65, 69, 97-99, 141, 157, 207, 209, 291, 306, 324, 383, 399, 573, 757, 799, 1051, 1055, 1095
  - group, 65
  - ICMP group, 307
  - interface group, 99
  - IP group, 209, 573
  - MIB-II, 65
  - TCP group, 799
  - UDP group, 757
- SO\_ACCEPTCONN constant, 437, 440, 730, 931
- SO\_BROADCAST socket option, 230, 347, 437, 539, 768, 882, 1118
- SO\_DEBUG socket option, 437, 539, 846, 882, 891, 916, 920, 931, 992, 994, 1010
- SO\_DONTROUTE socket option, 229, 347, 437, 499, 539, 738, 768, 882
- SO\_ERROR socket option, 539
- SO\_ISCONFIRMING constant, 514
- SO\_KEEPACTIVE socket option, 437, 539, 818, 829, 849, 1025, 1085, 1121
- SO\_LINGER socket option, 473, 539, 547, 557, 820, 892, 1009, 1019-1020, 1025, 1090
- SO\_OOBINLINE socket option, 437, 506, 510, 539, 985-986, 1016
- SO\_RCVBUF socket option, 539, 543, 932, 1012
- SO\_RCVLOWAT socket option, 539, 543
- SO\_RCVTIMEO socket option, 477, 516, 539, 543, 548
- SO\_REUSEADDR socket option, 437, 539, 720-721, 723, 725, 730-731, 733-735, 740, 753, 777
- SO\_REUSEPORT socket option, 437, 539, 721, 723, 725, 730-731, 734-735, 777
- SO\_SNDBUF socket option, 491, 539, 543
- SO\_SNDLOWAT socket option, 539, 543
- SO\_SNDTIMEO socket option, 477, 496, 539, 543, 548
- SO\_TYPE socket option, 539
  - SO\_USELOOPBACK socket option, 437, 539, 650, 660, 666
  - so\_dat member, 579
  - so\_error member, 440, 460, 494, 530-531, 539, 548, 783, 905, 948, 1087
  - so\_head member, 440-442, 473
  - so\_linger member, 439, 463, 473, 539, 542, 1010, 1019
  - so\_oobmark member, 440, 522, 531, 985, 1015
  - so\_options member, 347, 437, 463, 539, 547, 1010
  - so\_rcv member, 13, 440, 665-666, 671, 714
  - so\_pgid member, 439-440, 463, 550, 552
  - so\_proto member, 440, 446-447, 450, 531
  - so\_q member, 440-442, 461-464
  - so\_q0 member, 440-442, 461-463
  - so\_q0len member, 440, 463
  - so\_q1en member, 440, 463, 530, 1078
  - so\_q1limit member, 440, 455, 463
  - so\_rcv member, 440, 477, 509-510, 530-531, 539, 550, 878, 985, 1012
  - so\_snd member, 440, 477, 509-510, 531, 539, 861
  - so\_state member, 439, 463, 530-531, 550
  - so\_timeo member, 440
  - so\_tpcb member, 441
  - so\_type member, 13, 437, 446, 539
  - so\_upcall member, 441
  - so\_upcallarg member, 441
  - soabort function, 471
  - soaccept function, 460-461
  - sobind function, 453-455
  - socancelremote function, 442, 470, 990, 1089
  - socancelremote function, 442, 669, 788, 1013, 1061
  - sock program, 712, 740, 846, 1025
  - SOCK\_DGRAM constant, 10, 13, 188-191, 198, 437, 483, 713, 755, 758
  - SOCK\_RAW constant, 188-191, 193, 196-197, 240, 309, 385, 645-646, 1049, 1051-1053, 1072, 1080
  - SOCK\_RDM constant, 188-190, 483, 508
  - SOCK\_SEQPACKET constant, 188-190, 483, 490, 508
  - SOCK\_STREAM constant, 188-191, 196, 198, 437, 449, 483, 490, 508, 713, 801
  - sockaddr structure, 73-76, 89, 94, 117, 120-121, 155, 160, 182, 221, 322, 453, 474, 479, 482-483, 505, 694, 1046
  - sockaddr\_dl structure, 77-78, 86-92, 94, 97, 118, 120-121, 159, 581, 677, 688, 692, 695, 699, 706-707, 709, 1070, 1080



- Ethernet, 91
- loopback, 1070
- SLIP, 1070
- `sockaddr_in` structure, 6, 76, 118, 155, 160–162, 166, 170, 183, 193, 312, 453, 460, 477, 564, 577–578, 581, 585, 588, 593–594, 623, 696, 701, 730, 736, 738, 742, 753, 761, 776–777, 782, 785, 944, 1055, 1059–1060, 1079–1081
- `sockaddr_inarp` structure, 701, 703, 707, 1112
- `sockargs` function, 451–453, 465, 1081
- `sockfd` structure, 56, 470, 476–479
- socket, 14
  - buffers, 476–477
  - descriptor, 6, 445–447
  - I/O, 475–535
  - layer, 9, 435–474
  - pair, 6
  - routing, 569, 645–673
  - TCP, 198
  - UDP, 198
  - utility functions, 477–479
- socket function, 6–8, 10–11, 13–14, 21, 198, 440, 444–448, 476–477, 627, 645, 647, 662, 664–665, 667, 671, 673, 713, 717, 730, 733, 785, 822, 960, 1009–1010, 1013, 1018, 1024–1025, 1051, 1053, 1055–1056, 1058, 1065
- socket option, 537–557
  - `DMRE_ADD_LGSP`, 401, 411–413
  - `DMRE_ADD_MRT`, 401, 419, 421–422
  - `DMRE_ADD_VIP`, 401, 407, 409
  - `DMRE_DEL_LGSP`, 401, 411–412, 414
  - `DMRE_DEL_MRT`, 401, 421
  - `DMRE_DEL_VIP`, 401, 407, 409–410
  - `DMRE_DONE`, 401, 433, 1059
  - `DMRE_INIT`, 401, 403
  - `IP_ADD_MEMBERSHIP`, 348, 356–357, 434, 451, 1076
  - `IP_DROP_MEMBERSHIP`, 348, 366, 451
  - `IP_HDRINCL`, 191, 1053, 1056–1058, 1065, 1091
  - `IP_MULTICAST_IF`, 348, 353–354, 371, 710, 738
  - `IP_MULTICAST_LOOP`, 348, 354–355, 371
  - `IP_MULTICAST_TTL`, 348, 354, 371
  - `IP_OPTIONS`, 230, 240, 242–243, 269–270, 717, 1056, 1065, 1113, 1115, 1122
  - `IP_RECVDSADDR`, 240, 242, 776, 781, 793–794, 1083, 1114
  - `IP_RECVOPTS`, 240, 242–243, 769, 776, 794
  - `IP_RECVTPOPTS`, 240, 242–243, 776
  - `IP_RETPOPTS`, 794
  - `IP_TOS`, 240, 242, 717, 1056, 1099, 1115, 1123
  - `IP_TTL`, 240, 242, 717, 1056, 1100, 1115, 1120
- multicast, 348
- `SO_BROADCAST`, 230, 347, 437, 539, 768, 882, 1118
- `SO_DEBUG`, 437, 539, 846, 882, 891, 916, 920, 931, 992, 994, 1010
- `SO_DONTROUTE`, 229, 347, 437, 409, 539, 738, 768, 882
- `SO_ERROR`, 539
- `SO_KEEPAFLIVE`, 437, 539, 818, 829, 849, 1025, 1085, 1121
- `SO_LINGER`, 473, 539, 547, 557, 820, 892, 1009, 1019–1020, 1025, 1090
- `SO_OOBINLINE`, 437, 506, 510, 539, 985–986, 1016
- `SO_RCVBUF`, 539, 543, 932, 1012
- `SO_RCVLOWAT`, 539, 543
- `SO_RCVTIMEO`, 477, 516, 539, 543, 548
- `SO_REUSEADDR`, 437, 539, 720–721, 723, 725, 730–731, 733–735, 740, 753, 777
- `SO_REUSEPORT`, 437, 539, 721, 723, 725, 730–731, 734–735, 777
- `SO_SNDBUF`, 491, 539, 543
- `SO_SNDLOWAT`, 539, 543
- `SO_SNDTIMEO`, 477, 496, 539, 543, 548
- `SO_TYPE`, 539
- `SO_USELOOPBACK`, 437, 539, 650, 660, 666
- `TCP_MAXSEG`, 1022, 1024
- `TCP_NODELAY`, 858, 1022, 1024, 1120
- socket structure, 11, 13–14, 21, 56, 347, 398, 437–442, 446–447, 449–450, 453, 461, 463, 471, 476, 509–510, 664–666, 671–672, 713–714, 716–719, 746, 777, 1009, 1018, 1058, 1059, 1085
- `socket_args` structure, 444, 447
- socketops variable, 437, 446–448
- socketpair function, 669, 787, 1012
- sockets API, 5–6
- sockmod streams module, 8
- `sockproto` structure, 626, 647, 664
- `socksize` function, 471–473, 1010, 1014, 1019, 1090
- `sockconnect` function, 464–465, 467–468, 788, 962, 1059, 1083
- `sockcreate` function, 166, 447–451
- `sockdisconnect` function, 468, 473, 897, 1083
- `sockfree` function, 473, 719, 1018
- software interrupt, 106–107, 138, 153, 205, 212, 436
- `socketopt` function, 240, 538, 545–548
- `socketoutofband` function, 533, 552, 985
- `sockconnected` function, 461–465, 787, 949, 962, 969, 1089
- `sockconnecting` function, 442, 464–465, 1012
- `sockdisconnected` function, 442, 669, 1021



- soisdisconnecting function, 442, 1020
- SOL\_SOCKET constant, 240, 539–540, 1022
- Solaris, 85, 721
- socket function, 455
- solutions to exercises, 1069–1092
- SOMAXCONN constant, 440, 455
- someconn function, 459, 461–464, 931, 944, 963, 1009–1010, 1018
- soo\_close function, 446, 448, 471
- soo\_ioctl function, 164, 446, 448, 549, 552–554
- soo\_read function, 446, 448, 501
- soo\_select function, 446, 448, 528–532
- soo\_write function, 446, 448, 480, 761
- sooqsque function, 442, 461
- sooqsque function, 442, 460–461, 463–464, 474, 1078
- soqueue variable, 463
- soreadable macro, 442, 529–530
- soreceive function, 475, 501, 503, 505–524, 534–535, 792, 1014–1015, 1078
- soreserve function, 671, 785
- sorflush function, 469–470, 473–474
- sorwakep macro, 463–464, 478, 533, 776–777, 779, 916, 940
- sosend function, 59, 111, 475, 480, 483, 485, 489–500, 506, 515–516, 522, 524, 534–535, 648, 650, 669, 761, 763–764, 789, 792–793, 1016, 1078, 1082, 1086–1087
  - performance, 500
- sosendallatonce macro, 442, 493, 522
- so\_setopt function, 240, 412, 538–544, 546
- soshutdown function, 469, 1013
- source address
  - Ethernet, 99
  - IP, 232
- source code
  - conventions, 1–3
  - copyright, xxi–xxii
  - organization, 26–28
- source quench, ICMP, 226, 314
- source route
  - example, 255
  - failure, ICMP, 257
  - nonlocal, 1104
  - option, 254–261
  - reversed, 240, 255, 258–259, 261, 332, 1103, 1108–1109
- sowakeup function, 478, 533, 552
- sowriteable macro, 442, 529, 531
- sowakeup macro, 464, 478, 533, 939, 978
- sp\_family member, 647, 660, 662
- sp\_protocol member, 627, 647, 662
- spanning tree, 416–418, 1077
  - example, 416
- sp10 function, 24
- sp1bio function, 24
- sp1clock function, 24
- sp1high function, 24
- sp1ing function, 23–26, 30, 43, 73, 94, 112, 138, 213, 469, 1069, 1078
- sp1net function, 23–25, 212, 298, 434, 436, 467, 499, 763, 785, 1007, 1022, 1069, 1083
- sp1softclock function, 24
- sp1tty function, 24–25, 138, 1069
- sp1x function, 24–26, 43, 73, 94, 148, 434, 436, 470, 496, 499, 1077–1078
- SPP (Sequenced Packet Protocol), 189–190, 490
- sprint\_d function, 88
- SS\_ACCEPTCONN constant, 455
- SS\_ASYNC constant, 439–440, 552
- SS\_CANTBCVMORE constant, 439, 442, 530, 1078
- SS\_CANTSENDMORE constant, 439, 442, 531
- SS\_ISCONFIRMING constant, 439, 463, 500
- SS\_ISCONNECTED constant, 439, 442, 531
- SS\_ISCONNECTING constant, 439, 442, 467
- SS\_ISDISCONNECTING constant, 439, 442
- SS\_NBIO constant, 439–440, 550
- SS\_NOFDREF constant, 439, 741, 1018
- SS\_PRIV constant, 166, 439, 450–451
- SS\_RCVATMARK constant, 439, 522, 531, 550, 1015
- SSRR (strict source and record route), 249–250, 254–255, 257, 270–271
- st\_bikeize member, 789, 1015
- Stallings, W., 100, 106, 1131
- stat structure, 789, 1015
- state transition diagram, TCP, 805–807
- statistics, kernel, 37
- Stevens, D. A., xxii
- Stevens, D. L., 457, 1126
- Stevens, E. M., xxii
- Stevens, S. H., xxii
- Stevens, W. R., 5–7, 9, 11, 186, 435, 440, 470, 524, 732, 1077, 1131
- Stevens, W. R., xxii
- strcpy function, 8
- streams module, sockmod, 8
- streams subsystem, SVR4, 5, 8, 749
- strict source and record route, see SSRR
- strong end system model, 219, 780
- subnet
  - address, 1129
  - mask, 162
- SUBNETSARELOCAL constant, 901
- subnetsarelocal variable, 181, 901, 1105

- subnetting, 170
  - IP, 156, 170, 181, 1071
- superman, building leaping ability, 1102
- supermetting, 170
  - IP, 170, 1071
- superuser privileges, 451
- suser function, 451
- SVR4, 4
  - streams subsystem, 5, 8, 749
- SWS (silly window syndrome), 858, 878, 1120
- sy\_call member, 443
- sy\_narg member, 443
- SYN (synchronize sequence numbers flag, TCP header), 441, 461, 463, 803, 805, 828, 871
- synchronize sequence numbers flag, TCP header, *see* SYN
- syncall function, 441, 443–444, 454, 456, 489, 1045
- \_\_sysctl function, 202
- sysctl function, 67, 201–202, 239, 244–245, 334, 571–572, 601, 632–635, 637–638, 640, 644, 672, 679, 756, 790, 1083, 1100
- sysctl names, 201
- sysctl program, 191, 201, 319, 334, 680, 790, 1114
- sysctl\_dumpentry function, 632, 638, 640–642
- sysctl\_iflist function, 632, 638, 640, 642–643
- sysctl\_int function, 245, 334, 790
- sysctl\_rtable function, 203, 571, 635–642, 646, 679
- sysent structure, 437, 441, 443
- sysent variable, 437, 443
- system call
  - accept, 457–461
  - bind, 453–455
  - close, 471–473
  - connect, 464–468
  - fcntl, 548–552
  - getpeername, 554–556
  - getsockname, 554
  - ioctl, 548–550, 552–554
  - listen, 455
  - read, 500–501
  - readv, 500–501
  - recvfrom, 500–501
  - recvmsg, 500–503
  - select, 524–528
  - sendmsg, 480–484
  - sendto, 480–483
  - shutdown, 468–470
  - socket, 447–450
  - write, 480–483
  - writv, 480–483
- system calls, 7, 441–445
  - and library functions, 7–8
- system, vs. router and host, 157
- t\_dupacks member, 844, 970, 972–973
- t\_flags member, 805, 819, 871–872
- t\_force member, 827, 855, 862, 874, 882, 1017
- t\_idle member, 822, 826–828, 830–831, 849, 933, 1085
- t\_inpcb member, 714, 834
- t\_labc member, 985–987, 1016
- t\_maxseg member, 834, 844, 852, 896, 901, 903, 934, 1023–1024
- t\_newtcpcb member, 841, 843
- t\_oobflags member, 1016
- t\_oproc member, 131, 141
- t\_ospeed member, 132
- t\_outq member, 131, 141
- t\_rcvuderr function, 749
- t\_rtseq member, 837, 880
- t\_stt member, 823, 837, 840, 844, 880, 899, 917–918, 950, 972, 976
- t\_rttmin member, 832, 834, 841, 898–899
- t\_rttvar member, 832–835, 839, 843–844, 846, 848, 894, 898–899
- t\_rxtcur member, 832, 834, 840, 843, 846, 848, 882, 900, 917, 939, 976, 1119
- t\_rxtshift member, 832, 835–836, 840–843, 846, 855
- t\_sc member, 131–132, 135, 149
- t\_softerror member, 841, 843, 905, 1122
- t\_srtt member, 832–834, 837, 843–844, 846, 848, 894, 898–899
- t\_ssthresh member, 844
- t\_state member, 800, 805, 807
- t\_template member, 876, 880, 884, 888, 892
- t\_timer member, 819, 843, 1017
- TA\_DROP constant, 917
- TA\_INPUT constant, 917
- TA\_OUTPUT constant, 917
- TA\_USER constant, 917
- TAC, 982
- TAI (International Atomic Time), 1127
- Tanenbaum, A. S., 416, 1131
- Taylor, I. L., xxii
- tcp variable, 715–716, 718, 732, 744, 797, 813, 824
- TCP (Transmission Control Protocol), 65, 189, 191, 228, 240, 267, 440, 477, 541, 795–1025, 1130
  - checksum, 800

- control block, 713, 718, 800, 803–805, 808, 818–819, 821–822, 832–835, 837, 846, 866–867, 871–872, 884, 887–888, 893, 897, 906–907, 909, 916, 930, 932, 944–946, 949, 959–960, 966, 969, 981, 986, 989, 1009–1010, 1018–1019, 1021, 1023, 1084
- header, 801–803
- implementation refinements, 994–995
- inetaw variable, 198–199, 801
- input processing, 923–1005
- one-behind cache, 231, 798, 897, 929, 941
- options, 864–866, 1117
- output processing, 851–890
- performance, 1126
- protoaw structure, 801
- pseudo-header, 880, 885, 926, 995
- reassembly, 50–51, 906–916
- receive sequence space, 808
- RFC 1122 compliance, 1115–1123
- segments, demultiplexing, 721–723, 728
- send sequence space, 808
- sequence numbers, 807–812
- socket, 198
- state transition diagram, 805–807
- three-way handshake, 440, 465, 556, 722, 915, 917, 969, 1014
- timer, 817–849
- transactions, 866, 1089, 1126
- TCP\_COMPAT\_42 constant, 887
- TCP\_ISSINCR constant, 824, 945, 959, 1012
- TCP\_LINGERTIME constant, 820, 1009–1010, 1025
- TCP\_MAXRXTSHIFT constant, 820, 832, 836, 842
- TCP\_MAXSEG socket option, 1022, 1024
- TCP\_MAXWIN constant, 835, 864, 932, 1012
- TCP\_MAX\_WINSHIFT constant, 835, 932
- TCP\_NODELAY socket option, 858, 1022, 1024, 1120
- TCP\_PAWS\_IDLE constant, 953, 966
- TCP\_REASS Macro, 891, 906–916, 940, 987–989, 1004, 1089
- TCP\_RECVFVAL macro, 840, 843
- TCP\_RTT\_SCALE constant, 833, 899
- TCP\_RTT\_SHIFT constant, 833
- TCP\_RTTVAR\_SCALE constant, 833
- TCP\_RTTVAR\_SHIFT constant, 833
- tcp\_attach function, 1009, 1018
- tcp\_backoff variable, 832, 836, 1121
- tcp\_cancel timers: function, 821
- tcp\_close function, 825–826, 891–897, 930, 959, 981, 1010, 1019, 1021, 1087
- tcp\_ctlinput function, 198, 617, 743–744, 796, 801, 891, 904, 906, 1107, 1122
- tcp\_ctloutput function, 240, 538, 541, 796, 801, 1007, 1022–1025
- tcp\_debug structure, 916
- tcp\_debug variable, 917
- tcp\_debugx variable, 917
- tcp\_disconnect function, 851, 892, 1010, 1012, 1019–1021, 1025, 1087
- tcp\_dooptions function, 928, 933–934, 944, 951, 1067, 1088
- tcp\_dooptions variable, 797, 834, 865, 872, 889
- tcp\_drain function, 796, 801, 891–892
- tcp\_drop function, 828–830, 841, 851, 886, 891–892, 905, 948, 1014, 1019
- tcp\_fasttime function, 796, 801, 821–822, 848, 851, 1084, 1120
- tcp\_init function, 194, 796, 801, 812–815, 822, 824, 849
- tcp\_input function, 2, 461–462, 735, 743, 796, 801–802, 821–822, 825, 828, 830, 837, 843, 851, 859–860, 863–864, 867–868, 871, 873–875, 880, 884, 886–887, 892–893, 897, 901, 903, 906–907, 909, 917, 923–1005, 1009, 1018, 1023, 1088–1089, 1120, 1122
- tcp\_iss variable, 797, 812–814, 824, 871, 945–947, 1012, 1118
- tcp\_keepidle variable, 797, 819, 828, 830–831, 1121
- tcp\_keepintvl variable, 797, 819, 826, 830, 1084
- tcp\_last\_inpcb variable, 797, 897, 929
- tcp\_maxidle variable, 797, 819, 822, 825–826, 830, 849, 1084
- tcp\_mss function, 572, 834–835, 872, 891, 894, 897–903, 934, 944, 1023–1024
- tcp\_mssdiff variable, 797, 834, 898, 901
- tcp\_newtcb function, 832–835, 837, 846, 865, 871, 884, 899, 901, 903, 949, 1018, 1023, 1119–1120
- tcp\_notify function, 303–304, 743–744, 807, 843, 891, 904–905, 1122
- tcp\_now variable, 797, 824, 836–837, 867–868, 874, 890, 934, 937, 953, 966, 975
- tcp\_outflags variable, 797, 805, 808, 854, 892, 946, 961, 1012, 1021, 1090
- tcp\_output function, 56–58, 764, 795–796, 802, 805, 808, 821–823, 827, 836, 841, 843–844, 851–850, 892, 897, 901, 903, 906, 916–917, 923, 932, 939, 946–948, 961–962, 967, 972–973, 975, 978, 992, 994–995, 1004, 1012–1014, 1017, 1020–1021, 1084, 1086, 1088–1091, 1120
- tcp\_pullpcb function, 985–988, 1016
- tcp\_quench function, 2, 303–304, 743, 882, 891, 904, 906, 1122
- tcp\_recvqinit macro, 948

- tcp\_rwnd function, 51, 891, 906–916, 949, 969, 987–989, 1089, 1120
- tcp\_sndwspace variable, 797, 932, 1012, 1018
- tcp\_respond function, 45, 830, 885–888, 994, 1086
- tcp\_rstofid variable, 797, 834, 849, 1085
- tcp\_saver structure, 931
- tcp\_sndasginit macro, 947–948, 968, 1012
- tcp\_sndwspace variable, 797, 1018
- tcp\_seq data type, 810
- tcp\_setpersist function, 827, 835–836
- tcp\_slowtimo function, 796, 801, 822–824, 826, 830, 837, 848–849, 1017, 1084
- tcp\_template function, 876, 884–885, 888, 944, 966, 1012, 1088
- tcp\_timer function, 571, 822, 824–831, 841–846, 851, 886, 893, 906, 1017
- tcp\_trace function, 882, 891, 916–920, 931, 992, 994, 1009, 1017
- tcp\_upto\_closed function, 1013, 1020–1021
- tcp\_usrreq function, 461, 465, 796, 801, 822, 826, 851, 884, 917, 962, 1007–1018, 1025
- tcp\_xmit\_timer function, 834, 836–841, 843, 846, 882, 905, 939, 950, 975
- tcpActiveOpens variable, 800
- tcpAttemptFail variable, 800
- tcpb structure, 440, 714, 716–717, 803, 832, 916
- tcpConnLocalAddress variable, 800
- tcpConnLocalPort variable, 800
- tcpConnRemAddress variable, 800
- tcpConnRemPort variable, 800
- tcpConnState variable, 800
- tcpconnchug variable, 916
- tcpCurrEstab variable, 799–800
- TCPDEBUG constant, 916
- tcpdump program, 101, 917, 919–920, 1002, 1025, 1027, 1034, 1043
- tcpEstablishes variable, 800
- tcpHdr structure, 801, 906, 913
- tcpInErrs variable, 800
- tcpInSegs variable, 800
- tcpiphdr structure, 802, 906–907, 916
- tcpMaxConn variable, 800
- TCPLEN\_MAXSEG constant, 934
- TCPLEN\_TSTAMP constant, 934
- TCPLEN\_TSTAMP\_APPA constant, 874, 928
- TCPLEN\_WINDOW constant, 934
- TCPPOB\_ADDDATA constant, 985, 1016
- TCPPOB\_HAVEADATA constant, 985, 987, 1016
- TCPPOB\_MAXSEG constant, 872
- TCPPOB\_TSTAMP\_RDR constant, 874, 928
- tcpPortRes variable, 800
- tcpOutSegs variable, 800
- tcpPassiveOpens variable, 800
- tcpRetransSegs variable, 800
- tcpresmtthresh variable, 797, 970, 972
- tcpSrtAlgorithm variable, 800
- knprtoMax variable, 800
- knprtoMin variable, 800
- TCP\_S\_CLOSED constant, 807
- TCP\_S\_CLOSE\_WAIT constant, 807
- TCP\_S\_CLOSING constant, 807
- TCP\_S\_ESTABLISHED constant, 807
- TCP\_S\_FIN\_WAIT\_1 constant, 807
- TCP\_S\_FIN\_WAIT\_2 constant, 807
- TCP\_S\_HAVERCVDFIN macro, 807, 983
- TCP\_S\_HAVERCVDSYN macro, 807
- TCP\_S\_LAST\_ACK constant, 807
- TCP\_S\_LISTEN constant, 807
- TCP\_S\_SYN\_RECEIVED constant, 807
- TCP\_S\_SYN\_SENT constant, 807
- TCP\_S\_TIME\_WAIT constant, 807
- tcp\_s\_accepts member, 798–800
- tcp\_s\_closed member, 798
- tcp\_s\_connattempt member, 798–800
- tcp\_s\_conndrops member, 798–800
- tcp\_s\_connects member, 798–799
- tcp\_s\_delay member, 798–799
- tcp\_s\_drops member, 798–800
- tcp\_s\_keepprobes member, 798–799, 831
- tcp\_s\_keepprobe member, 798–799
- tcp\_s\_keeptimeo member, 798–799
- tcp\_s\_pawdrop member, 798–799, 954
- tcp\_s\_policyschemiss member, 798–799
- tcp\_s\_persisttimeo member, 798–799
- tcp\_s\_prndack member, 798–799
- tcp\_s\_prnddat member, 798–799
- tcp\_s\_rcvackbyte member, 797–799
- tcp\_s\_rcvackpack member, 798–799
- tcp\_s\_rcvacktoomuch member, 798–799
- tcp\_s\_rcvafterclose member, 798–799
- tcp\_s\_rcvbadoff member, 798–800
- tcp\_s\_rcvbadsum member, 798–800
- tcp\_s\_rcvbyte member, 798–799
- tcp\_s\_rcvbyteafterwin member, 798–799
- tcp\_s\_rcvdupack member, 798–799, 972
- tcp\_s\_rcvdupbyte member, 798–799
- tcp\_s\_rcvduppack member, 798–799, 954
- tcp\_s\_rcvordbyte member, 798–799
- tcp\_s\_rcvordpack member, 798–799
- tcp\_s\_rcvpack member, 798–799
- tcp\_s\_rcvpackafterwin member, 798–799
- tcp\_s\_rcvpartdupbyte member, 798–799
- tcp\_s\_rcvpartduppack member, 798–799



- tcps\_rcvshort member, 798-800, 927
- tcps\_rcvtotal member, 798-800
- tcps\_rcvwindup member, 799
- tcps\_rcvwinprobe member, 798-799
- tcps\_rcvwinupd member, 798, 983
- tcps\_rexmttimeo member, 798-799
- tcps\_rttupdated member, 798-799
- tcps\_segstimed member, 798-799
- tcps\_sndacks member, 798
- tcps\_sndbyte member, 797-799
- tcps\_sndctrl member, 798-799
- tcps\_sndpack member, 797-799
- tcps\_sndprobe member, 798-799
- tcps\_sndrexmitbyte member, 798-799
- tcps\_sndrexmitpack member, 798-800
- tcps\_sndtotal member, 798-800
- tcps\_sndurg member, 798-799
- tcps\_sndwinup member, 798-799
- tcps\_timeoutdrop member, 798-799
- tcpstat structure, 797-799
- tcpstat variable, 797
- TCPT\_2MSL constant, 819, 825, 849
- TCPT\_KEEP constant, 819, 828, 831
- TCPT\_NTIMERS constant, 819
- TCPT\_PERSIST constant, 819, 827
- TCPT\_RANGESET macro, 820, 834-836, 900
- TCPT\_REXMT constant, 819
- tcpTable variable, 799
- TCPTV\_KEEPCNT constant, 820, 830
- TCPTV\_KEEP\_IDLE constant, 820
- TCPTV\_KEEP\_INIT constant, 819-820, 828, 946
- TCPTV\_KEEPIVTL constant, 820
- TCPTV\_MIN constant, 820, 834
- TCPTV\_MSL constant, 813, 819-820, 825
- TCPTV\_PERSMAX constant, 820
- TCPTV\_PERSMIN constant, 820
- TCPTV\_SRTTMAX constant, 820, 834, 841
- TCPTV\_SRTTBASE constant, 820, 834
- TCPTV\_SRTTDELT constant, 820, 835
- td\_act member, 917
- Telnet (remote terminal protocol), 4, 140, 157, 272, 721, 753, 858, 982, 995, 1084, 1125
- telnet: program, 183
- test network, 28
- TF\_ACKNOW constant, 805, 821, 852, 861, 876, 884, 946, 948, 956, 961, 990, 1089
- TF\_DELACK constant, 805, 819, 821, 884, 1120
- TF\_NODELAY constant, 805, 858, 1023-1024
- TF\_NOOPT constant, 805, 871
- TF\_RCVD\_SCALE constant, 805
- TF\_RCVD\_TSTMP constant, 805, 874
- TF\_REQ\_SCALE constant, 805, 834, 872
- TF\_REQ\_TSTMP constant, 805, 834, 872
- TF\_SENTFIN constant, 805, 861, 877, 880
- TFTP (Trivial File Transfer Protocol), 140, 291, 776, 1125
- TH\_ACK constant, 803, 854, 872, 994
- TH\_FIN constant, 803, 854, 908, 916, 1090
- TH\_PUSH constant, 803, 854
- TH\_RST constant, 803, 854, 994
- TH\_SYN constant, 803, 854
- TH\_URG constant, 803, 854
- th\_ack member, 803, 1000-1001
- th\_flags member, 802-803
- th\_off member, 802, 878
- th\_seq member, 802, 1000-1001
- th\_urg member, 1001
- th\_urg member, 802-803
- th\_win member, 1000-1001
- Thimer, M., 100, 321, 1127
- three-way handshake, TCP, 440, 465, 556, 722, 915, 917, 969, 1014
- ti\_ack member, 837, 867, 938-939, 975
- ti\_sport member, 906, 909
- ti\_len member, 809-810, 868, 884, 906-907, 909, 927, 938-939, 951, 956, 959, 963, 987, 994, 1088
- ti\_next member, 885, 906, 909
- ti\_off member, 809, 885, 926
- ti\_pr member, 884
- ti\_prev member, 885, 906, 909
- ti\_seq member, 809, 871, 906-907, 909, 915, 937, 951, 954, 956, 961, 981-982, 989, 1088
- ti\_sport member, 906, 908
- ti\_t member, 913
- ti\_urg member, 879, 954, 983
- ti\_win member, 878
- tick variable, 544, 548
- time exceeded, ICMP, 223, 292-293, 300, 314
- time variable, 105, 699, 1081
- TIME\_WAIT, assassination, 964, 1089, 1126
- timeout function, 94, 195-196, 706, 1010
- timer
  - 2MSL, 818-819, 821-822, 825-827, 893, 967
  - connection-establishment, 817, 819, 828-831, 892, 946, 948, 1012, 1121
  - delayed ACK, 817-818, 821, 861, 864
  - example, 846-848
  - FIN\_WAIT\_2, 818-819, 821-822, 825-827, 980, 991, 1085
  - functions, ARP, 694-696
  - keepalive, 818-819, 821, 828-831, 842, 887, 892, 918, 933
  - persist, 818-821, 827, 835, 855, 858, 861, 878, 882, 939, 976, 1116



- retransmission, 749, 817, 819–822, 828, 831–833, 835, 841, 843–844, 846, 859, 861, 880, 882, 889, 918, 920, 939, 970–972, 975–976, 978, 1012, 1020, 1086
- TCP, 817–849
- timestamp option, 261–264, 865–871, 874, 933, 934, 1126
- timestamp reply, ICMP, 318
- timestamp request, ICMP, 318
- time-to-live, *see* TTL
- timeval structure, 105–106, 264, 525, 539, 543–544, 548, 1035
- timevaladd function, 525
- TIOCPCMSG constant, 552
- tk\_nin variable, 128, 135
- TLI (Transport Layer Interface), 5, 749
- Topolcic, C., 215, 1131
- Torek, C., 730, 1131
- TOS (type of service), 140–141, 147–148, 153, 226, 230, 240, 285, 302–303, 316, 328, 717, 768, 785, 882, 998, 1002, 1056, 1112, 1115, 1123
  - queuing, SLIP, 140, 147
- TP4, 189–190, 457, 463–464, 490, 494, 498, 508, 510, 514, 555–556
  - links program, 30
  - traceroute program, 140, 191, 272, 313
  - transactions, TCP, 866, 1089, 1126
  - Transmission Control Protocol, *see* TCP
  - Transport Layer Interface, *see* TLI
  - Trivial File Transfer Protocol, *see* TFTP
  - TRPB (truncated reverse path broadcast), 401, 416–419, 434
    - routing, example, 417
  - trpt program, 846, 891, 916–917, 1017
  - truncated reverse path broadcast, *see* TRPB
  - truss program, 30
  - ts\_eof variable, 837, 867–868, 928, 934, 975–976
  - ts\_present variable, 837, 868, 928, 934, 951
  - ts\_recent member, 867–871, 874, 934, 937, 951, 953, 963
  - ts\_recent\_age member, 867–868, 934, 937, 953, 966
  - ts\_val variable, 867, 870, 928, 934, 937, 951
  - tsleep function, 441, 456–457, 459, 461–465, 467, 473, 478, 528, 532, 534, 544, 1010, 1046, 1079
  - TSTMF\_GRQ macro, 870, 937
  - TSTMF\_LT macro, 953
  - TTL (time-to-live), 209, 216, 221, 223, 230, 239, 244–245, 292, 329, 339, 348, 351, 354–355, 371, 377, 379, 389, 428–430, 432, 717, 768, 785, 835, 882, 889, 1053, 1056, 1115, 1120
    - default, 207
    - multicast, 348–351
    - reassembly, 298
  - TTY, 141
    - device driver, 129–130, 134–135, 141, 148–149
    - output queue, 134, 141
  - tty structure, 131–132, 135, 141, 149
  - TTY\_CHARMAX constant, 135
  - ttyflush function, 132
  - tunnel\_send function, 424–425, 429, 431–433
  - tv\_sec member, 544, 1081
  - tv\_usec member, 544, 548, 669
  - type, ICMP, 250, 302–303
  - type of service, *see* TOS
  - typographical conventions, 3
  - u\_char data type, 250, 340, 348
  - u\_int data type, 1035
  - u\_long data type, 160–162
  - udb variable, 14, 21, 715–716, 718, 732, 744, 756, 760, 786, 1081
  - UDP (User Datagram Protocol), 65, 73, 189, 191, 228, 240, 440, 477, 541, 755–794
    - cache hiding, 791
    - checksum, 758, 764–768, 792
    - datagrams, demultiplexing, 723–724
    - header, 759–760
    - implementation refinements, 791–792
    - inetsw variable, 203, 758
    - input processing, 769–780
    - one-behind cache, 231, 757, 773–774, 786, 791, 794
    - output processing, 760–768
    - protocol structure, 758
    - pseudo-header, 764–768
    - RFC 1122 compliance, 1113–1115
    - socket, 198
    - socket, connected, 721, 755, 779–780
    - socket, unconnected, 721, 755
    - sockets, ICMP errors and, 748–749
  - udp\_ctlinput function, 198, 617, 743–744, 756, 758, 782–784, 793, 904, 1107
  - udp\_detach function, 786
  - udp\_in variable, 756, 776–777, 794
  - udp\_init function, 194, 756, 758, 760
  - udp\_input function, 743, 756, 758–759, 769–781, 791, 793–794, 929, 1084, 1113–1115
  - udp\_last\_inpcb variable, 756, 773–774
  - udp\_notify function, 303–304, 743–744, 783–784
  - udp\_output function, 741, 756, 758–768, 772, 789, 792–793, 882, 1082, 1084, 1115
  - udp\_recvspace variable, 756, 785

- udp\_saveopt function, 776, 781, 793
- udp\_sendspace variable, 756, 785
- udp\_soreceive function, 792
- udp\_sosend function, 792
- udp\_sysctl function, 203, 756, 758, 790-791
- udp\_usrreq function, 116, 163-164, 327, 465, 756, 758, 761, 784-790, 793, 1083
- udpccksum variable, 756, 768, 772, 790, 1114
- UDPCTL\_CHECKSUM constant, 202
- udphdr structure, 759
  - udpInDatagrams variable, 758
  - udpInErrors variable, 758
  - udpiphdr structure, 759, 765, 767-768, 885
  - udpLocalAddress variable, 758
  - udpLocalPort variable, 758
  - udpNoPorts variable, 758
  - udpOutDatagrams variable, 758
- udpsa\_pcbcachealias member, 757, 774
- udpsa\_badlen member, 757-758
- udpsa\_badsum member, 757-758
- udpsa\_fullsock member, 757
- udpsa\_hdrops member, 757-758
- udpsa\_ipackets member, 757-758, 774
- udpsa\_noport member, 757-758, 774
- udpsa\_noportbcst member, 757-758, 774
- udpsa\_opackets member, 757
- udpstat structure, 757, 774
- udpstat variable, 756
- udpstable variable, 757
- uh\_dport member, 759
- uh\_sport member, 759
- uh\_sum member, 759
- uh\_ulen member, 759, 771-772
- ui\_dst member, 765
- ui\_len member, 765, 768
- ui\_next member, 765
- ui\_pr member, 765
- ui\_prev member, 765
- ui\_src member, 765
- ui\_ulen member, 768
- ui\_xl member, 765
- uio structure, 476, 485-487, 489, 491, 503, 510, 1078
  - UIO\_MAXIOV constant, 481, 483, 500
  - UIO\_READ constant, 486
  - UIO\_SMALLIOV constant, 483
  - UIO\_SYSPACE constant, 486
  - UIO\_USERSPACE constant, 486
  - UIO\_USERSPACE constant, 486-487
  - UIO\_WRITE constant, 486-487
  - uio\_iov member, 485-487
  - uio\_iovcnt member, 486-487
- uio\_offset member, 485-487
- uio\_procp member, 486-487
- uio\_resid member, 485-487, 489, 503, 505, 511, 515, 519, 1078
- uio\_rw member, 485-487
- uio\_segflg member, 485-487
- uio\_move function, 485-487, 498, 519-520, 792, 1046, 1079
- unconnected UDP socket, 721, 755
- unicast, 155
  - address, Ethernet, 100
  - address, IP, 155-156, 182
  - one-behind cache, 223, 253
  - packet, 99
- uniform resource locator, *see* URL
- Unix
  - address family, 185
  - domain, 75, 189, 450, 460, 470, 510, 518, 1077
  - domain protocol, 9-10, 40, 581, 718, 787
  - Epoch, 105, 683, 695
  - manual pages, 3
  - protocol family, 185
- unixdomain variable, 187, 193, 195
- unixsw variable, 195
- unreachable, ICMP, 314
- unreliable protocol buffers, 490-491
- unsocial behavior, 951
- URG (urgent pointer flag, TCP header), 803, 878, 956, 983-986, 1000-1002, 1004, 1017, 1086
- urgent
  - offset, 802-803, 878-880, 889, 954, 956, 983, 985, 987-988, 1000-1002, 1016-1017, 1086
  - pointer, 802, 861, 876, 878-880, 951, 984-987, 1004, 1016, 1116-1117
  - pointer flag, TCP header, *see* URG
- URL (uniform resource locator), 1093, 1125
- uselookback variable, 680, 709
- User Datagram Protocol, *see* UDP
- UTC (Coordinated Universal Time), 105-106, 261, 264, 318, 1127
- utility functions
  - ifnet, 182
  - IP, 181
  - socket, 477-479
- v\_cached\_group member, 406, 412
- v\_cached\_result member, 406, 412
- v\_flags member, 406
- v\_ifp member, 406-407
- v\_lcl\_addr member, 406, 410
- v\_lcl\_grps member, 406-407, 411-412
- v\_lcl\_grps\_max member, 406-407, 412

- v\_lcl\_grps\_0 member, 406-407
- v\_rmt\_addr member, 406
- v\_threshold member, 406
- Varadhan, K., 170, 1127
- Vardhana, G. N. A., xxd
- VAX, 24, 60, 1068, 1117
- vif structure, 398, 406-410, 412
- vifc\_flags member, 408
- vifc\_lcl\_addr member, 408-409
- vifc\_rmt\_addr member, 408
- vifc\_threshold member, 408
- vifc\_vifi member, 408-409
- vifctl structure, 401, 407, 409
- VIFF\_TUNNEL constant, 406
- vifc\_r data type, 398, 401, 406
- viftable variable, 398, 406-408, 410, 418, 429, 434
- villain, 1067
- virtual
  - interface, 404-411
  - interface table, 406-410
- vmstat program, 40, 716
- vnode structure, 13
- vsprintf function, 6
  
- w\_arg member, 639
- w\_given member, 639-640
- w\_needed member, 634-635, 639-640
- w\_op member, 639
- w\_tmem member, 635, 642
- w\_tmemsize member, 635
- w\_where member, 635, 639-640, 642-643
- Wait, J. W., xcd
- Waitzman, D., 401, 1131
- Wakeman, L., 1126
- wakeup function, 441, 456-457, 461, 463-465, 467, 477, 532
- walkarg structure, 632, 634, 639-641
- WAN (wide area network), 901, 940-941
- wandering duplicate, 813
- Wang, Z., 1126
- Watson, G., 994, 1126
- weak end system model, 219, 741, 780, 1072, 1100-1101
- Wei, L., 419, 1127
- well-known
  - IP multicast groups, 338-339
  - multicast groups, 338
  - port, 719, 729, 733, 774, 814, 1011
- wide area network, *see* WAN
- wildcard match, 722
  
- window
  - scale option, 865-866, 871, 874, 929, 933-934, 1128
  - shrink, 856-857, 878, 884, 1116
  - update, 859-861, 863-864, 876, 981, 983, 1014
- Wolff, R., xcd
- Wolf, S., xcd
- write function, 8, 10, 13, 56, 129, 435, 439, 445-447, 475, 478, 480-481, 534, 650, 752, 760-761, 874, 890, 920, 1025, 1046, 1056, 1080, 1092
- writew function, 13, 445-446, 475, 480-481, 760-761, 890, 1056, 1087
  
- X.25, 27, 580
- Xerox Network Systems, *see* XNS
- XNS (Xerox Network Systems), 9-10, 23, 27, 39, 189, 562, 581, 624
- X/Open, 5
  - Transport Layer Interface, *see* XTI
- XTI (X/Open Transport Layer Interface), 5
- XXX comment, 70, 91, 141, 655, 763, 882, 913, 944, 1018
  
- Yu, J. Y., 170, 1127
  
- zero\_in\_addr variable, 715



## Function and Macro Definitions

|                        |      |                   |     |                  |     |
|------------------------|------|-------------------|-----|------------------|-----|
| accept                 | 458  | ifa_ifwithnet     | 182 | ip_init          | 200 |
| add_lgrp               | 413  | ifa_ifwithroute   | 182 | ip_insertoptions | 266 |
| add_mrt                | 422  | ifaof_ifpforaddr  | 182 | ipintr           | 213 |
| add_vif                | 408  | if_attach         | 88  | ip_mforward      | 426 |
| arpintr                | 687  | ifconf            | 118 | ip_mloopback     | 378 |
| arplookup              | 702  | IF_DEQUEUE        | 72  | ip_mroutercmd    | 402 |
| arprequest             | 685  | if_down           | 123 | ip_mrouterdone   | 433 |
| arpresolve             | 697  | IF_DROP           | 72  | ip_mroutercmd    | 402 |
| arp_rtrrequest         | 705  | IF_ENQUEUE        | 72  | ip_mroutercmd    | 402 |
| arptfree               | 696  | ifinit            | 93  | ip_optcopy       | 282 |
| arptimer               | 695  | ifioctl           | 116 | ip_output        | 229 |
| arpwhoas               | 683  | IF_PREPEND        | 72  | ip_pcbopts       | 269 |
|                        |      | if_qflush         | 72  | ip_reass         | 290 |
| bind                   | 454  | IF_QFULL          | 72  | ip_rtaddr        | 254 |
| bpfattach              | 1031 | if_slowtimo       | 93  | ip_setmoptions   | 352 |
| bpf_attach             | 1040 | ifunit            | 182 | ip_slowtimo      | 299 |
| bpfioc1                | 1035 | ifup              | 123 | ip_srcroute      | 260 |
| bpfopen                | 1034 | igmp_fasttimo     | 389 | ip_sysctl        | 244 |
| bpfread                | 1044 | igmp_input        | 392 | iptime           | 264 |
| bpf_setif              | 1038 | igmp_joingroup    | 386 |                  |     |
| bpf_tap                | 1041 | igmp_leavegroup   | 395 | leattach         | 82  |
| bpfwrite               | 1047 | IGMP_RANDOM_DELAY | 387 | leioctl          | 124 |
|                        |      | igmp_sendreport   | 390 | leread           | 102 |
| catchpacket            | 1042 | in_addmulti       | 359 | lestart          | 113 |
| connect                | 466  | in_arpinput       | 689 | listen           | 455 |
|                        |      | in_broadcast      | 181 | loioctl          | 180 |
| del_lgrp               | 414  | in_canforward     | 181 | loopattach       | 85  |
| del_mrt                | 421  | in_cksum          | 237 | looutput         | 150 |
| del_vif                | 410  | in_control        | 165 |                  |     |
| domaininit             | 194  | in_delmulti       | 368 | m_adj            | 53  |
| dtom                   | 46   | IN_FIRST_MULTI    | 388 | main             | 79  |
|                        |      | in_ifinit         | 169 | m_cat            | 53  |
| ether_addmulti         | 364  | in_localaddr      | 181 | MCLGET           | 52  |
| ether_delmulti         | 370  | IN_LOOKUP_MULTI   | 347 | m_copy           | 53  |
| ether_ifattach         | 92   | in_losing         | 749 | m_copyback       | 53  |
| ether_input            | 104  | in_netof          | 181 | m_copydata       | 53  |
| ETHER_LOOKUP_MULTI     | 344  | IN_NEXT_MULTI     | 388 | m_copym          | 53  |
| ETHER_MAP_IP_MULTICAST | 342  | in_pcballoc       | 718 | m_devget         | 53  |
| ether_output           | 108  | in_pcbbind        | 729 | MFREE            | 52  |
|                        |      | in_pcbconnect     | 735 | m_free           | 53  |
| fcntl                  | 550  | in_pcbdetach      | 719 | m_freem          | 53  |
|                        |      | in_pcbdisconnect  | 741 | m_get            | 41  |
| getpeername            | 556  | in_pcblookup      | 726 | MGET             | 42  |
| getsock                | 452  | in_pcbnotify      | 745 | m_getclr         | 53  |
| getsockname            | 555  | in_rtchange       | 746 | MGETHDR          | 52  |
| getsockopt             | 545  | in_setpeeraddr    | 742 | m_gethdr         | 53  |
| grplst_member          | 415  | in_setsockaddr    | 742 | MH_ALIGN         | 52  |
|                        |      | insque            | 292 | M_LEADINGSPACE   | 764 |
| icmp_error             | 325  | ip_ctloutput      | 241 | M_PREPEND        | 52  |
| icmp_input             | 311  | ip_deq            | 292 | m_pullup         | 53  |
| icmp_reflect           | 330  | ip_dooptions      | 251 | m_retry          | 43  |
| icmp_send              | 333  | ip_drain          | 299 | mrtfind          | 423 |
| icmp_sysctl            | 334  | ip_enq            | 292 | mtod             | 46  |
| ifa_ifwithaddr         | 182  | ip_forward        | 222 |                  |     |
| ifa_ifwithaf           | 182  | ip_freef          | 299 | nethash          | 420 |
| ifa_ifwithdstaddr      | 182  | ip_getmoptions    | 372 | net_sysctl       | 203 |



## Function and Macro Definitions

|     |                 |      |                   |     |                   |      |
|-----|-----------------|------|-------------------|-----|-------------------|------|
| 200 | pfctlinput      | 198  | sblock            | 478 | sosetopt          | 541  |
| 266 | pffasttimo      | 196  | sbrelease         | 479 | soshutdown        | 469  |
| 213 | pffindproto     | 197  | sbreserve         | 479 | sowakeup          | 478  |
| 426 | pffindtype      | 197  | sbspace           | 478 | sowriteable       | 531  |
| 378 | pfslowtimo      | 196  | sbunlock          | 478 | sowakeup          | 478  |
| 402 | phyint_send     | 430  | sbwait            | 478 | sysctl_dumpentry  | 641  |
| 433 |                 |      | select            | 526 | sysctl_iflist     | 642  |
| 404 | raw_attach      | 671  | selrecord         | 532 | sysctl_rtable     | 638  |
| 282 | raw_detach      | 672  | selscan           | 529 |                   |      |
| 229 | raw_disconnect  | 672  | selwakeup         | 533 | tcp_attach        | 1019 |
| 269 | raw_init        | 648  | sendit            | 488 | tcp_canceltimers  | 821  |
| 290 | raw_input       | 662  | sendmsg           | 484 | tcp_close         | 895  |
| 254 | raw_usrreq      | 667  | SEQ_GEQ           | 810 | tcp_ctlinput      | 904  |
| 352 | recvit          | 503  | SEQ_GT            | 810 | tcp_ctloutput     | 1022 |
| 299 | recvmsg         | 502  | SEQ_LEQ           | 810 | tcp_disconnect    | 1020 |
| 260 | remque          | 292  | SEQ_LT            | 810 | tcp_dooptions     | 933  |
| 244 | rip_ctloutput   | 1064 | setsockopt        | 540 | tcp_drop          | 893  |
| 264 | rip_init        | 1053 | shutdown          | 468 | tcp_fasttimo      | 821  |
|     | rip_input       | 1054 | slattach          | 84  | tcp_init          | 812  |
| 82  | rip_output      | 1057 | slclose           | 148 | tcp_input         | 926  |
| 124 | rip_usrreq      | 1058 | slinit            | 133 | tcp_mss           | 898  |
| 102 | rn_init         | 584  | slinput           | 134 | tcp_newtcpcb      | 833  |
| 113 | rn_match        | 591  | slioctl           | 179 | tcp_notify        | 905  |
| 455 | rn_search       | 599  | slopen            | 132 | tcp_output        | 853  |
| 180 | route_init      | 584  | sloutput          | 139 | tcp_pulloutofband | 986  |
| 85  | route_output    | 652  | slstart           | 142 | tcp_quench        | 906  |
| 150 | route_usrreq    | 664  | sltioctl          | 149 | tcp_rcvseqinit    | 946  |
|     | rtable_init     | 584  | soaccept          | 460 | TCP_REASS         | 908  |
| 53  | rtalloc         | 602  | sobind            | 454 | tcp_reass         | 911  |
| 79  | rtalloc1        | 603  | socantrcvmore     | 442 | tcp_respond       | 886  |
| 53  | RTFREE          | 605  | socantsendmore    | 442 | tcp_sendseqinit   | 946  |
| 52  | rtfree          | 605  | sockargs          | 452 | tcp_setpersist    | 835  |
| 53  | rt_ifmsg        | 627  | socket            | 448 | tcp_slowtimo      | 823  |
| 53  | rtinit          | 616  | soclose           | 472 | tcp_template      | 885  |
| 53  | rt_missmsg      | 625  | soconnect         | 467 | tcp_timers        | 824  |
| 53  | rt_msg1         | 631  | socreate          | 449 | tcp_trace         | 918  |
| 53  | rt_msg2         | 633  | sofree            | 473 | TCPT_RANGESET     | 820  |
| 52  | rt_newaddrmsg   | 628  | sogetopt          | 546 | tcp_usrclosed     | 1021 |
| 53  | rtredirect      | 618  | soisconnected     | 464 | tcp_usrreq        | 1008 |
| 53  | rtrequest       | 607  | soisconnecting    | 442 | tcp_xmit_timer    | 838  |
| 41  | rt_setgate      | 614  | soisdisconnected  | 442 | tunnel_send       | 431  |
| 42  | rt_setmetrics   | 662  | soisdisconnecting | 442 |                   |      |
| 53  | rt_xaddrs       | 660  | solisten          | 456 | udp_ctlinput      | 783  |
| 52  |                 |      | sonewconn         | 462 | udp_detach        | 786  |
| 53  | save_rte        | 259  | soo_close         | 471 | udp_init          | 760  |
| 52  | sballloc        | 478  | soo_ioctl         | 553 | udp_input         | 770  |
| 764 | sbappend        | 479  | soo_select        | 530 | udp_notify        | 784  |
| 52  | sbappendaddr    | 479  | soqinsque         | 442 | udp_output        | 762  |
| 53  | sbappendcontrol | 479  | soqremque         | 442 | udp_saveopt       | 781  |
| 43  | sbappendrecord  | 479  | soreadable        | 530 | udp_sysctl        | 790  |
| 423 | sbcompress      | 479  | soreceive         | 512 | udp_usrreq        | 784  |
| 46  | sbdrop          | 479  | soreserve         | 479 |                   |      |
|     | sbdroprecord    | 479  | sorflush          | 470 |                   |      |
| 420 | sbflush         | 479  | sorwakeup         | 478 |                   |      |
| 203 | sbfree          | 478  | sosend            | 492 |                   |      |
|     | sbinsertoob     | 479  | sosendallatonce   | 442 |                   |      |

Networking, TCP/IP

# TCP/IP Illustrated, Volume 2

"Teaching Cisco internetworking classes around the country, I get to hear about every new data communications book of note because my students bring them to class. *TCP/IP Illustrated, Volume 1* arrived like a bomb going off. I've just read *TCP/IP Illustrated, Volume 2*. If *Volume 2* is the Big Bang!

— Greg Gentile, Director of Training  
Protocol Interface

LIBRARY OF CONGRESS



00023560690

"After reading Stevens' *TCP/IP Illustrated, Volume 1*, I thought it would be hard to come up with another book as useful, but *Volume 2* is just that. Some of you might ask how this book applies to you if you aren't a full-time network programmer. Since my earliest days of using UNIX®, the standard answer to 'How does this really work?' has been 'Look at the source.' With this book, you can not only look at the source, but also have one of the clearest explanations as to how it all fits together."

— Greg Merrell, Systems Engineer  
International Network Services

*TCP/IP Illustrated*, an ongoing series covering the many facets of TCP/IP, brings a highly-effective visual approach to learning about this networking protocol suite.

*TCP/IP Illustrated, Volume 2* contains a thorough explanation of how TCP/IP protocols are implemented. There isn't a more practical or up-to-date book — this volume is the only one to cover the de facto standard implementation from the 4.4BSD-Lite release, the foundation for TCP/IP implementations run daily on hundreds of thousands of systems worldwide.

Combining 500 illustrations with 15,000 lines of real, working code, *TCP/IP Illustrated, Volume 2* uses a teach-by-example approach to help you master TCP/IP implementation. You will learn about such topics as the relationship between the sockets API and the protocol suite, and the differences between a host implementation and a router. In addition, the book covers the newest features of the 4.4BSD-Lite release, including multicasting, long fat pipe support, window scale, timestamp options, and protection against wrapped sequence numbers, and many other topics.

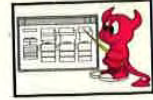
Comprehensive in scope, based on a working standard, and thoroughly illustrated, this book is an indispensable resource for anyone working with TCP/IP.

**Gary R. Wright** has worked with TCP/IP for more than eight years. He is President of Connix, a Connecticut-based company providing Internet access and consulting services. **W. Richard Stevens** is the highly-respected author of three best-selling books, *TCP/IP Illustrated, Volume 1* (Addison-Wesley, 1994), *Advanced Programming in the UNIX Environment* (Addison-Wesley, 1992), and *UNIX Network Programming* (Prentice-Hall). He is also a popular tutorials instructor and consultant.

Cover art by Fine Line Design  
♻️ Text printed on recycled paper  
Corporate & Professional Publishing Group  
♣️ Addison-Wesley Publishing Company



ISBN 0-201-63354-X



TCP/IP Illustrated, Volume 2

Wright  
Stevens

TK 5105  
.55  
.S74  
1994  
v. 2  
Copy 1

  
Addison  
Wesley  
63354