

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.

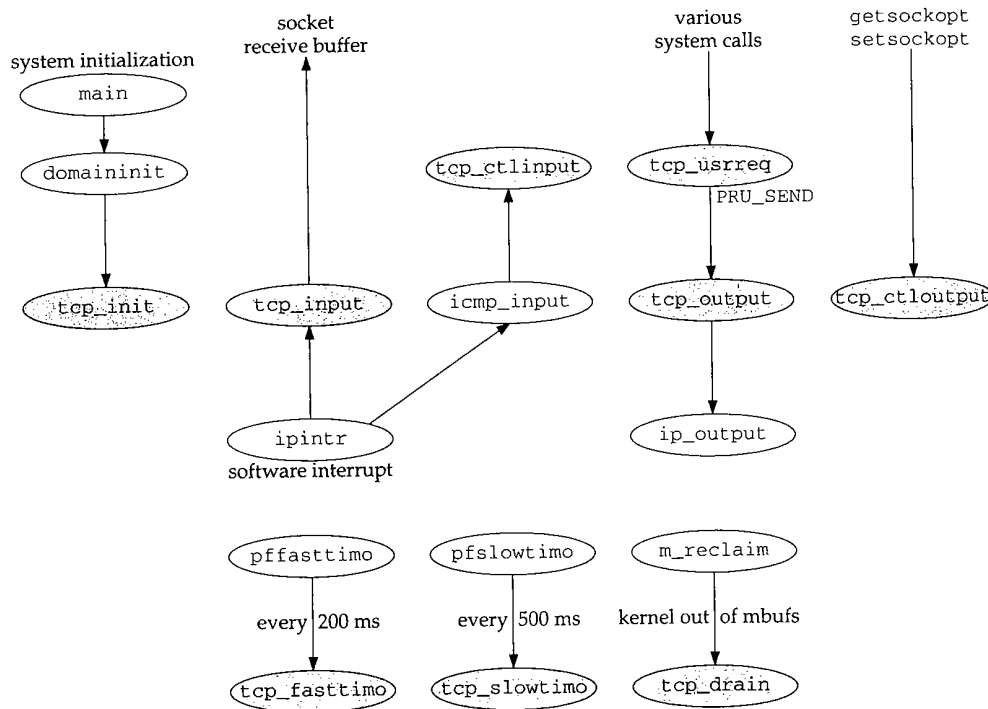


Figure 24.2 Relationship of TCP functions to rest of the kernel.

Global V

Statist

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_rexmtthresh</code> | <code>int</code> | number of duplicate ACKs to trigger fast retransmit (3) |
| <code>tcp_mssdflt</code> | <code>int</code> | default MSS (512 bytes) |
| <code>tcp_rttdeflt</code> | <code>int</code> | default RTT if no data (3 seconds) |
| <code>tcp_do_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 connect) |
| <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 | |
|----------------------|---|-----------------|---------|
| tcps_accepts | #SYNs received in LISTEN state | • | 10,655, |
| tcps_closed | #connections closed (includes drops) | • | 9,17 |
| tcps_connattempt | #connections initiated (calls to connect) | • | 257, |
| tcps_conndrops | #embryonic connections dropped (before SYN received) | • | 862, |
| tcps_connects | #connections established actively or passively | • | 229 |
| tcps_delack | #delayed ACKs sent | • | 3,45 |
| tcps_drops | #connections dropped (after SYN received) | • | 74,9 |
| tcps_keepprobes | #connections dropped in keepalive (established or awaiting SYN) | • | 279, |
| tcps_keeptimeo | #keepalive probes sent | • | 8,801,9 |
| tcps_pawdrop | #times keepalive timer or connection-establishment timer expire | • | 6,61 |
| tcps_pcbcachemiss | #segments dropped due to PAWS | • | 235, |
| tcps_persisttimeo | #times PCB cache comparison fails | • | 0 ac |
| tcps_predack | #times persist timer expires | • | 4,67 |
| tcps_preddat | #times header prediction correct for ACKs | • | 46,9 |
| tcps_rcvackbyte | #times header prediction correct for data packets | • | 22 c |
| tcps_rcvackpack | #bytes ACKed by received ACKs | • | 3,44 |
| tcps_rcvacktoomuch | #received ACK packets | • | 77,1 |
| tcps_rcvafterclose | #received ACKs for unsent data | • | 1,88 |
| tcps_rcvbadoff | #packets received after connection closed | • | 1,75 |
| tcps_rcvbadsum | #packets received with invalid header length | • | 175, |
| tcps_rcvbyte | #packets received with checksum errors | • | 1,01 |
| tcps_rcvbyteafterwin | #bytes received in sequence | • | 60,3 |
| tcps_rcvdupack | #bytes received beyond advertised window | • | 279 |
| tcps_rcvdupbyte | #duplicate ACKs received | • | 0 d: |
| tcps_rcvduppack | #bytes received in completely duplicate packets | • | 144,020 |
| tcps_rcvoobyte | #packets received with completely duplicate bytes | • | 92,595 |
| tcps_rcvoopack | #out-of-order bytes received | • | 126,820 |
| tcps_rcvpack | #out-of-order packets received | • | 237,740 |
| tcps_rcvpackafterwin | #packets received in sequence | • | 110,010 |
| tcps_rcvpartdupbyte | #packets with some data beyond advertised window | • | 6,363,1 |
| tcps_rcvpartduppack | #duplicate bytes in part-duplicate packets | • | 114,790 |
| tcps_rcvshort | #packets with some duplicate data | • | 86 |
| tcps_rcvtotal | #packets received too short | • | 1,173 |
| tcps_rcvwinprobe | total #packets received | • | 16,419 |
| tcps_rcvwinupd | #window probe packets received | • | 6,8 |
| tcps_rexmttimeo | #received window update packets | • | 3,2 |
| tcps_rttupdated | #retransmit timeouts | • | 733,13 |
| tcps_segstimed | #times RTT estimators updated | • | 1,266, |
| tcps_sndacks | #segments for which TCP tried to measure RTT | • | 1,851, |
| tcps_sndbyte | #ACK-only packets sent (data length = 0) | • | |
| tcps_sndctrl | #data bytes sent | • | |
| tcps_sndpack | #control (SYN, FIN, RST) packets sent (data length = 0) | • | |
| tcps_sndprobe | #data packets sent (data length > 0) | • | |
| tcps_sndrexmitbyte | #window probes sent (1 byte of data forced by persist timer) | • | |
| tcps_sndrexmitpack | #data bytes retransmitted | • | |
| tcps_sndtotal | #data packets retransmitted | • | |
| tcps_sndurg | total #packets sent | • | |
| tcps_sndwinup | #packets sent with URG-only (data length = 0) | • | |
| tcps_timeoutdrop | #window update-only packets sent (data length = 0) | • | |
| | #connections dropped in retransmission timeout | • | |

Figure 24.4 TCP statistics maintained in the tcpstat structure.

SNMP

| netstat -s output | tcpstat members |
|---|--|
| 10,655,999 packets sent 9,177,823 data packets (-22,194,928 bytes) 257,295 data packets (81,075,086 bytes) retransmitted 862,900 ack-only packets (531,285 delayed) 229 URG-only packets 3,453 window probe packets 74,925 window update packets 279,387 control packets | tcps_sndtotal tcps_snd(pack,byte) tcps_sndrexit(pack,byte) tcps_sndacks,tcps_delack tcps_sndurg tcps_sndprobe tcps_sndwinup tcps_sndctrl |
| 8,801,953 packets received 6,617,079 acks (for -21,264,360 bytes) 235,311 duplicate acks 0 acks for unsent data 4,670,615 packets (324,965,351 bytes) rcvd in-sequence 46,953 completely duplicate packets (1,549,785 bytes) 22 old duplicate packets 3,442 packets with some dup. data (54,483 bytes duped) 77,114 out-of-order packets (13,938,456 bytes) 1,892 packets (1,755 bytes) of data after window 1,755 window probes 175,476 window update packets 1,017 packets received after close 60,370 discarded for bad checksums 279 discarded for bad header offset fields 0 discarded because packet too short | tcps_rcvtotal tcps_rcvack(pack,byte) tcps_rcvdupack tcps_rcvacktoomuch tcps_rcv(pack,byte) tcps_rcvdup(pack,byte) tcps_pawsdrop tcps_rcvpartdup(pack,byte) tcps_rcvoo(pack,byte) tcps_rcv(pack,byte)afterwin tcps_rcwinprobe tcps_rcwindup tcps_rcvafterclose tcps_rcvbadsum tcps_rcvbadoff tcps_rcvshort |
| 144,020 connection requests 92,595 connection accepts 126,820 connections established (including accepts) 237,743 connections closed (including 1,061 drops) 110,016 embryonic connections dropped | tcps_connattempt tcps_accepts tcps_connects tcps_closed,tcps_drops tcps_conndrops |
| 6,363,546 segments updated rtt (of 6,444,667 attempts) 114,797 retransmit timeouts 86 connection dropped by rexmit timeout 1,173 persist timeouts 16,419 keepalive timeouts 6,899 keepalive probes sent 3,219 connections dropped by keepalive | tcps_{rttupdated,segstimed} tcps_rexmttimeo tcps_timeoutdrop tcps_persisttimeo tcps_keeptimeo tcps_keepprobe tcps_keepdrops |
| 733,130 correct ACK header predictions 1,266,889 correct data packet header predictions 1,851,557 cache misses | tcps_predack tcps_preddat 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.

Explore Litigation Insights

Docket Alarm provides insights to develop a more informed litigation strategy and the peace of mind of knowing you're on top of things.

Real-Time Litigation Alerts



Keep your litigation team up-to-date with **real-time alerts** and advanced team management tools built for the enterprise, all while greatly reducing PACER spend.

Our comprehensive service means we can handle Federal, State, and Administrative courts across the country.

Advanced Docket Research



With over 230 million records, Docket Alarm's cloud-native docket research platform finds what other services can't. Coverage includes Federal, State, plus PTAB, TTAB, ITC and NLRB decisions, all in one place.

Identify arguments that have been successful in the past with full text, pinpoint searching. Link to case law cited within any court document via Fastcase.

Analytics At Your Fingertips



Learn what happened the last time a particular judge, opposing counsel or company faced cases similar to yours.

Advanced out-of-the-box PTAB and TTAB analytics are always at your fingertips.

API

Docket Alarm offers a powerful API (application programming interface) to developers that want to integrate case filings into their apps.

LAW FIRMS

Build custom dashboards for your attorneys and clients with live data direct from the court.

Automate many repetitive legal tasks like conflict checks, document management, and marketing.

FINANCIAL INSTITUTIONS

Litigation and bankruptcy checks for companies and debtors.

E-DISCOVERY AND LEGAL VENDORS

Sync your system to PACER to automate legal marketing.