
19581 Mon May 26 17:49:34 2014

new/usr/src/man/man7p/tcp.7p

4774 Typos in tcp(7P) manpage

```

1 \" te
2.\" Copyright (c) 2006, Sun Microsystems, Inc. All Rights Reserved.
3.\" Copyright (c) 2011 Nexenta Systems, Inc. All rights reserved.
4.\" Copyright 1989 AT&T
5.\" The contents of this file are subject to the terms of the Common Development
6.\" You can obtain a copy of the license at usr/src/OPENSOLARIS.LICENSE or http:
7.\" When distributing Covered Code, include this CDDL HEADER in each file and in
8.TH TCP 7P "Jun 30, 2006"
9.SH NAME
10 tcp, TCP \- Internet Transmission Control Protocol
11.SH SYNOPSIS
12.LP
13.nf
14 \fB#include <sys/socket.h>\fR
15.fi

17.LP
18.nf
19 \fB#include <netinet/in.h>\fR
20.fi

22.LP
23.nf
24 \fBs = socket(AF_INET, SOCK_STREAM, 0);\fR
25.fi

27.LP
28.nf
29 \fBs = socket(AF_INET6, SOCK_STREAM, 0);\fR
30.fi

32.LP
33.nf
34 \fBt = t_open("/dev/tcp", O_RDWR);\fR
35.fi

37.LP
38.nf
39 \fBt = t_open("/dev/tcp6", O_RDWR);\fR
40.fi

42.SH DESCRIPTION
43.sp
44.LP
45 \fBTCP\fR is the virtual circuit protocol of the Internet protocol family. It
46 provides reliable, flow-controlled, in order, two-way transmission of data. It
47 is a byte-stream protocol layered above the Internet Protocol (\fBIP\fR), or
48 the Internet Protocol Version 6 (\fBIPv6\fR), the Internet protocol family's
49 internetwork datagram delivery protocol.
50.sp
51.LP
52 Programs can access \fBTCP\fR using the socket interface as a \fB SOCK_STREAM\fR
53 socket type, or using the Transport Level Interface (\fBTLI\fR) where it
54 supports the connection-oriented (\fB COTS_ORD\fR) service type.
55.sp
56.LP
57 \fBTCP\fR uses \fBIP\fR's host-level addressing and adds its own per-host
58 collection of "port addresses." The endpoints of a \fBTCP\fR connection are
59 identified by the combination of an \fBIP\fR or IPv6 address and a \fBTCP\fR
60 port number. Although other protocols, such as the User Datagram Protocol
61 (UDP), may use the same host and port address format, the port space of these

```

```

62 protocols is distinct. See \fBinet\fR(7P) and \fBinet6\fR(7P) for details on
63 the common aspects of addressing in the Internet protocol family.
64.sp
65.LP
66 Sockets utilizing \fBTCP\fR are either "active" or "passive." Active sockets
67 initiate connections to passive sockets. Both types of sockets must have their
68 local \fBIP\fR or IPv6 address and \fBTCP\fR port number bound with the
69 \fBbind\fR(3SOCKET) system call after the socket is created. By default,
70 \fBTCP\fR sockets are active. A passive socket is created by calling the
71 \fBlisten\fR(3SOCKET) system call after binding the socket with \fBbind()\fR.
72 This establishes a queueing parameter for the passive socket. After this,
73 connections to the passive socket can be received with the
74 \fBaccept\fR(3SOCKET) system call. Active sockets use the
75 \fBconnect\fR(3SOCKET) call after binding to initiate connections.
76.sp
77.LP
78 By using the special value \fBINADDR_ANY\fR with \fBIP\fR, or the unspecified
79 address (all zeroes) with IPv6, the local \fBIP\fR address can be left
80 unspecified in the \fBbind()\fR call by either active or passive \fBTCP\fR
81 sockets. This feature is usually used if the local address is either unknown or
82 irrelevant. If left unspecified, the local \fBIP\fR or IPv6 address will be
83 bound at connection time to the address of the network interface used to
84 service the connection.
85.sp
86.LP
87 Note that no two TCP sockets can be bound to the same port unless the bound IP
88 addresses are different. IPv4 \fBINADDR_ANY\fR and IPv6 unspecified addresses
89 compare as equal to any IPv4 or IPv6 address. For example, if a socket is bound
90 to \fBINADDR_ANY\fR or unspecified address and port X, no other socket can bind
91 to port X, regardless of the binding address. This special consideration of
92 \fBINADDR_ANY\fR and unspecified address can be changed using the socket option
93 \fB SO_REUSEADDR\fR. If \fB SO_REUSEADDR\fR is set on a socket doing a bind, IPv4
94 \fBINADDR_ANY\fR and IPv6 unspecified address do not compare as equal to any IP
95 address. This means that as long as the two sockets are not both bound to
96 \fBINADDR_ANY\fR/unspecified address or the same IP address, the two sockets
97 can be bound to the same port.
98.sp
99.LP
100 If an application does not want to allow another socket using the
101 \fB SO_REUSEADDR\fR option to bind to a port its socket is bound to, the
102 application can set the socket level option \fB SO_EXCLBIND\fR on a socket. The
103 option values of 0 and 1 mean enabling and disabling the option respectively.
104 Once this option is enabled on a socket, no other socket can be bound to the
105 same port.
106.sp
107.LP
108 Once a connection has been established, data can be exchanged using the
109 \fBread\fR(2) and \fBwrite\fR(2) system calls.
110.sp
111.LP
112 Under most circumstances, \fBTCP\fR sends data when it is presented. When
113 outstanding data has not yet been acknowledged, \fBTCP\fR gathers small amounts
114 of output to be sent in a single packet once an acknowledgement has been
115 received. For a small number of clients, such as window systems that send a
116 stream of mouse events which receive no replies, this packetization may cause
117 significant delays. To circumvent this problem, \fBTCP\fR provides a
118 socket-level boolean option, \fB TCP_NODELAY\fR. \fB TCP_NODELAY\fR is defined in
119 \fB<netinet/tcp.h>\fR, and is set with \fBsetsockopt\fR(3SOCKET) and tested
120 with \fBgetsockopt\fR(3SOCKET). The option level for the \fBsetsockopt()\fR
121 call is the protocol number for \fBTCP\fR available from
122 \fBgetprotobyname\fR(3SOCKET).
123.sp
124.LP
125 For some applications, it may be desirable for TCP not to send out data unless
126 a full TCP segment can be sent. To enable this behavior, an application can use
127 the \fB TCP_CORK\fR socket option. When \fB TCP_CORK\fR is set with a non-zero

```

128 value, TCP sends out a full TCP segment only. When `\fBTCP_CORK`\fR is set to
 129 zero after it has been enabled, all buffered data is sent out (as permitted by
 130 the peer's receive window and the current congestion window). `\fBTCP_CORK`\fR is
 131 defined in `<\fBnetinet/tcp.h`\fR, and is set with `\fBsetsockopt`\fR(3SOCKET)
 132 and tested with `\fBgetsockopt`\fR(3SOCKET). The option level for the
 133 `\fBsetsockopt`\fR() call is the protocol number for TCP, available from
 134 `\fBgetprotobyname`\fR(3SOCKET).

135 .sp
 136 .LP
 137 The `SO_RCVBUF` socket level option can be used to control the window that TCP
 138 advertises to the peer. IP level options may also be used with TCP. See
 139 `\fBip`\fR(7P) and `\fBip6`\fR(7P).

140 .sp
 141 .LP
 142 Another socket level option, `\fBSO_RCVBUF`\fR, can be used to control the window
 143 that `\fBTCP`\fR advertises to the peer. `\fBIP`\fR level options may also be used
 144 with `\fBTCP`\fR. See `\fBip`\fR(7P) and `\fBip6`\fR(7P).

145 .sp
 146 .LP
 147 `\fBTCP`\fR provides an urgent data mechanism, which may be invoked using the
 148 out-of-band provisions of `\fBsend`\fR(3SOCKET). The caller may mark one byte as
 149 "urgent" with the `\fBMSG_OOB`\fR flag to `\fBsend`\fR(3SOCKET). This sets an
 150 "urgent pointer" pointing to this byte in the `\fBTCP`\fR stream. The receiver on
 151 the other side of the stream is notified of the urgent data by a `\fBSIGURG`\fR
 152 signal. The `\fBSIOCATMARK`\fR `\fBioctl`\fR(2) request returns a value indicating
 153 whether the stream is at the urgent mark. Because the system never returns data
 154 across the urgent mark in a single `\fBread`\fR(2) call, it is possible to
 155 advance to the urgent data in a simple loop which reads data, testing the
 156 socket with the `\fBSIOCATMARK`\fR `\fBioctl`\fR() request, until it reaches the
 157 mark.

158 .sp
 159 .LP
 160 Incoming connection requests that include an `\fBIP`\fR source route option are
 161 noted, and the reverse source route is used in responding.

162 .sp
 163 .LP
 164 A checksum over all data helps `\fBTCP`\fR implement reliability. Using a
 165 window-based flow control mechanism that makes use of positive
 166 acknowledgements, sequence numbers, and a retransmission strategy, `\fBTCP`\fR
 167 can usually recover when datagrams are damaged, delayed, duplicated or
 168 delivered out of order by the underlying communication medium.

169 .sp
 170 .LP
 171 If the local `\fBTCP`\fR receives no acknowledgements from its peer for a period
 172 of time, (for example, if the remote machine crashes), the connection is closed
 173 and an error is returned.

174 .sp
 175 .LP
 176 TCP follows the congestion control algorithm described in `\fIRFC 2581`\fR, and
 177 also supports the initial congestion window (`wnd`) changes in `\fIRFC 3390`\fR.
 178 The initial `wnd` calculation can be overridden by the socket option
 179 `TCP_INIT_CWND`. An application can use this option to set the initial `wnd` to a
 180 specified number of TCP segments. This applies to the cases when the connection
 181 first starts and restarts after an idle period. The process must have the
 182 `PRIV_SYS_NET_CONFIG` privilege if it wants to specify a number greater than that
 183 calculated by `\fIRFC 3390`\fR.

184 .sp
 185 .LP
 186 SunOS supports `\fBTCP`\fR Extensions for High Performance (`\fIRFC 1323`\fR) which
 187 **includes the window scale and timestamp options, and Protection Against Wrap**
 188 *includes the window scale and time stamp options, and Protection Against Wrap*
 189 Around Sequence Numbers (PAWS). SunOS also supports Selective Acknowledgment
 189 (SACK) capabilities (RFC 2018) and Explicit Congestion Notification (ECN)
 190 mechanism (`\fIRFC 3168`\fR).

191 .sp
 192 .LP

193 Turn on the window scale option in one of the following ways:

194 .RS +4
 195 .TP
 196 .ie t \(\bu
 197 .el o
 198 An application can set `\fBSO_SNDBUF`\fR or `\fBSO_RCVBUF`\fR size in the
 199 `\fBsetsockopt`\fR() call option to be larger than 64K. This must be done `\fBbefore`\fR
 200 the program calls `\fBlisten`\fR() or `\fBconnect`\fR(), because the window scale
 201 option is negotiated when the connection is established. Once the connection
 202 has been made, it is too late to increase the send or receive window beyond the
 203 default `\fBTCP`\fR limit of 64K.

204 .RE
 205 .RS +4
 206 .TP
 207 .ie t \(\bu
 208 .el o
 209 For all applications, use `\fBndd`\fR(1M) to modify the configuration parameter
 210 `\fBtcp_wscale_always`\fR. If `\fBtcp_wscale_always`\fR is set to `\fB1`\fR, the
 211 window scale option will always be set when connecting to a remote system. If
 212 `\fBtcp_wscale_always`\fR is `\fB0`\fR, the window scale option will be set only if
 213 the user has requested a send or receive window larger than 64K. The default
 214 value of `\fBtcp_wscale_always`\fR is `\fB1`\fR.

215 .RE
 216 .RS +4
 217 .TP
 218 .ie t \(\bu
 219 .el o
 220 Regardless of the value of `\fBtcp_wscale_always`\fR, the window scale option
 221 will always be included in a connect acknowledgement if the connecting system
 222 has used the option.

223 .RE
 224 .sp
 225 .LP
 226 Turn on `\fBSACK`\fR capabilities in the following way:

227 .RS +4
 228 .TP
 229 .ie t \(\bu
 230 .el o
 231 Use `\fBndd`\fR to modify the configuration parameter `\fBtcp_sack_permitted`\fR.
 232 If `\fBtcp_sack_permitted`\fR is set to `\fB0`\fR, `\fBTCP`\fR will not accept
 233 `\fBSACK`\fR or send out `\fBSACK`\fR information. If `\fBtcp_sack_permitted`\fR is
 234 set to `\fB1`\fR, `\fBTCP`\fR will not initiate a connection with `\fBSACK`\fR
 235 permitted option in the `\fBSYN`\fR segment, but will respond with `\fBSACK`\fR
 236 permitted option in the `\fBSYN`\fR segment if an incoming connection request
 237 has the `\fBSACK`\fR permitted option. This means that `\fBTCP`\fR will only
 238 accept `\fBSACK`\fR information if the other side of the connection also accepts
 239 `\fBSACK`\fR information. If `\fBtcp_sack_permitted`\fR is set to `\fB2`\fR, it will
 240 both initiate and accept connections with `\fBSACK`\fR information. The default
 241 for `\fBtcp_sack_permitted`\fR is `\fB2`\fR (active enabled).

242 .RE
 243 .sp
 244 .LP
 245 Turn on `\fBTCP`\fR ECN mechanism in the following way:

246 .RS +4
 247 .TP
 248 .ie t \(\bu
 249 .el o
 250 Use `\fBndd`\fR to modify the configuration parameter `\fBtcp_ecn_permitted`\fR. If
 251 `\fBtcp_ecn_permitted`\fR is set to `\fB0`\fR, `\fBTCP`\fR will not negotiate with a
 252 peer that supports `\fBECN`\fR mechanism. If `\fBtcp_ecn_permitted`\fR is set to
 253 `\fB1`\fR when initiating a connection, TCP will not tell a peer that it supports
 254 ECN mechanism. However, it will tell a peer that it supports `\fBECN`\fR
 255 mechanism when accepting a new incoming connection request if the peer
 256 indicates that it supports `\fBECN`\fR mechanism in the SYN segment. If
 257 `tcp_ecn_permitted` is set to 2, in addition to negotiating with a peer on ECN
 258 mechanism when accepting connections, TCP will indicate in the outgoing SYN

259 segment that it supports \fBECN\fR mechanism when \fBTCP\fR makes active
 260 outgoing connections. The default for \fBtcp_ecn_permitted\fR is 1.

261 .RE
 262 .sp
 263 .LP

264 **Turn on the timestamp option in the following way:**
 264 *Turn on the time stamp option in the following way:*

265 .RS +4
 266 .TP
 267 .ie t \(\bu
 268 .el o
 269 Use \fBndd\fR to modify the configuration parameter \fBtcp_timestamp_always\fR. If
 270 \fBtcp_timestamp_always\fR is \fB1\fR, the timestamp option will always be set
 271 when connecting to a remote machine. If \fBtcp_timestamp_always\fR is \fB0\fR, the
 272 timestamp option will not be set when connecting to a remote system. The
 273 default for \fBtcp_timestamp_always\fR is \fB0\fR.

274 .RE
 275 .RS +4
 276 .TP
 277 .ie t \(\bu
 278 .el o
 279 **Regardless of the value of \fBtcp_timestamp_always\fR, the timestamp option will**
 279 *Regardless of the value of \fBtcp_timestamp_always\fR, the time stamp option will*
 280 *always be included in a connect acknowledgement (and all succeeding packets) if*
 281 **the connecting system has used the timestamp option.**
 281 *the connecting system has used the time stamp option.*

282 .RE
 283 .sp
 284 .LP

285 **Use the following procedure to turn on the timestamp option only when the**
 285 *Use the following procedure to turn on the time stamp option only when the*
 286 *window scale option is in effect:*

287 .RS +4
 288 .TP
 289 .ie t \(\bu
 290 .el o
 291 Use \fBndd\fR to modify the configuration parameter \fBtcp_timestamp_if_wscale\fR.
 292 **Setting \fBtcp_timestamp_if_wscale\fR to \fB1\fR will cause the timestamp option**
 292 *Setting \fBtcp_timestamp_if_wscale\fR to \fB1\fR will cause the time stamp option*
 293 *to be set when connecting to a remote system, if the window scale option has*
 294 **been set. If \fBtcp_timestamp_if_wscale\fR is \fB0\fR, the timestamp option will**
 294 *been set. If \fBtcp_timestamp_if_wscale\fR is \fB0\fR, the time stamp option will*
 295 *not be set when connecting to a remote system. The default for*
 296 *\fBtcp_timestamp_if_wscale\fR is \fB1\fR.*

297 .RE
 298 .sp
 299 .LP

300 Protection Against Wrap Around Sequence Numbers (PAWS) is always used when the
 301 **timestamp option is set.**
 301 *time stamp option is set.*

302 .sp
 303 .LP

304 SunOS also supports multiple methods of generating initial sequence numbers.
 305 One of these methods is the improved technique suggested in \fBIRFC\fR 1948. We
 306 \fBHIGHLY\fR recommend that you set sequence number generation parameters as
 307 close to boot time as possible. This prevents sequence number problems on
 308 connections that use the same connection-ID as ones that used a different
 309 sequence number generation. The \fBsvc:/network/initial:default\fR service
 310 configures the initial sequence number generation. The service reads the value
 311 contained in the configuration file \fB/etc/default/inetinit\fR to determine
 312 which method to use.

313 .sp
 314 .LP

315 The \fB/etc/default/inetinit\fR file is an unstable interface, and may change
 316 in future releases.

317 .sp
 318 .LP

319 \fBTCP\fR may be configured to report some information on connections that
 320 terminate by means of an \fBRST\fR packet. By default, no logging is done. If
 321 the \fBndd\fR(1M) parameter \fBtcp_trace\fR is set to 1, then trace data is
 322 collected for all new connections established after that time.

323 .sp
 324 .LP

325 The trace data consists of the \fBTCP\fR headers and \fBIP\fR source and
 326 destination addresses of the last few packets sent in each direction before RST
 327 occurred. Those packets are logged in a series of \fBstrlog\fR(9F) calls. This
 328 trace facility has a very low overhead, and so is superior to such utilities as
 329 \fBsnop\fR(1M) for non-intrusive debugging for connections terminating by
 330 means of an \fBRST\fR.

331 .sp
 332 .LP

333 SunOS supports the keep-alive mechanism described in \fBIRFC 1122\fR. It is
 334 enabled using the socket option SO_KEEPALIVE. When enabled, the first
 335 keep-alive probe is sent out after a TCP is idle for two hours. If the peer does
 336 not respond to the probe within eight minutes, the TCP connection is aborted.
 337 You can alter the interval for sending out the first probe using the socket
 338 option TCP_KEEPALIVE_THRESHOLD. The option value is an unsigned integer in
 339 milliseconds. The system default is controlled by the TCP ndd parameter
 340 tcp_keepalive_interval. The minimum value is ten seconds. The maximum is ten
 341 days, while the default is two hours. If you receive no response to the probe,
 342 you can use the TCP_KEEPALIVE_ABORT_THRESHOLD socket option to change the time
 343 threshold for aborting a TCP connection. The option value is an unsigned
 344 integer in milliseconds. The value zero indicates that TCP should never time
 345 out and abort the connection when probing. The system default is controlled by
 346 the TCP ndd parameter tcp_keepalive_abort_interval. The default is eight
 347 minutes.

348 .sp
 349 .LP

350 socket options TCP_KEEPIIDLE, TCP_KEEPCNT and TCP_KEEPIIDLE are also supported
 351 for compatibility with other Unix Flavors. TCP_KEEPIIDLE option specifies the
 352 interval in seconds for sending out the first keep-alive probe. TCP_KEEPCNT
 353 specifies the number of keep-alive probes to be sent before aborting the
 354 connection in the event of no response from peer. TCP_KEEPIIDLE specifies the
 355 interval in seconds between successive keep-alive probes.

356 .SH SEE ALSO

357 .sp
 358 .LP

359 \fBsvcs\fR(1), \fBndd\fR(1M), \fBbioctl\fR(2), \fBbradvd\fR(2), \fBsvcadm\fR(1M),
 360 \fBwrite\fR(2), \fBaccept\fR(3SOCKET), \fBbind\fR(3SOCKET),
 361 \fBconnect\fR(3SOCKET), \fBgetprotobyname\fR(3SOCKET),
 362 \fBgetsockopt\fR(3SOCKET), \fBlisten\fR(3SOCKET), \fBsend\fR(3SOCKET),
 363 \fBsmf\fR(5), \fBinet\fR(7P), \fBinet6\fR(7P), \fBip\fR(7P), \fBip6\fR(7P)

364 .sp
 365 .LP

366 Ramakrishnan, K., Floyd, S., Black, D., RFC 3168, \fBThe Addition of Explicit
 367 Congestion Notification (ECN) to IP\fR, September 2001.

368 .sp
 369 .LP

370 Mathias, M. and Hahdavi, J. Pittsburgh Supercomputing Center; Ford, S. Lawrence
 371 Berkeley National Laboratory; Romanow, A. Sun Microsystems, Inc. \fBIRFC 2018,
 372 TCP Selective Acknowledgement Options\fR, October 1996.

373 .sp
 374 .LP

375 Bellovin, S., \fBIRFC 1948, Defending Against Sequence Number Attacks\fR, May
 376 1996.

377 .sp
 378 .LP

379 Jacobson, V., Braden, R., and Borman, D., \fBIRFC 1323, TCP Extensions for High
 380 Performance\fR, May 1992.

381 .sp
 382 .LP

```

383 Postel, Jon, \fIRFC 793, Transmission Control Protocol - DARPA Internet Program
384 Protocol Specification\fR, Network Information Center, SRI International, Menlo
385 Park, CA., September 1981.
386 .SH DIAGNOSTICS
387 .sp
388 .LP
389 A socket operation may fail if:
390 .sp
391 .ne 2
392 .na
393 \fB\fBEISCONN\fR\fR
394 .ad
395 .RS 17n
396 A \fBconnect()\fR operation was attempted on a socket on which a
397 \fBconnect()\fR operation had already been performed.
398 .RE

400 .sp
401 .ne 2
402 .na
403 \fB\fBETIMEDOUT\fR\fR
404 .ad
405 .RS 17n
406 A connection was dropped due to excessive retransmissions.
407 .RE

409 .sp
410 .ne 2
411 .na
412 \fB\fBECONNRESET\fR\fR
413 .ad
414 .RS 17n
415 The remote peer forced the connection to be closed (usually because the remote
416 machine has lost state information about the connection due to a crash).
417 .RE

419 .sp
420 .ne 2
421 .na
422 \fB\fBECONNREFUSED\fR\fR
423 .ad
424 .RS 17n
425 The remote peer actively refused connection establishment (usually because no
426 process is listening to the port).
427 .RE

429 .sp
430 .ne 2
431 .na
432 \fB\fBEADDRINUSE\fR\fR
433 .ad
434 .RS 17n
435 A \fBbind()\fR operation was attempted on a socket with a network address/port
436 pair that has already been bound to another socket.
437 .RE

439 .sp
440 .ne 2
441 .na
442 \fB\fBEADDRNOTAVAIL\fR\fR
443 .ad
444 .RS 17n
445 A \fBbind()\fR operation was attempted on a socket with a network address for
446 which no network interface exists.
447 .RE

```

```

449 .sp
450 .ne 2
451 .na
452 \fB\fBEACCES\fR\fR
453 .ad
454 .RS 17n
455 A \fBbind()\fR operation was attempted with a "reserved" port number and the
456 effective user \fBID\fR of the process was not the privileged user.
457 .RE

459 .sp
460 .ne 2
461 .na
462 \fB\fBENOBUFS\fR\fR
463 .ad
464 .RS 17n
465 The system ran out of memory for internal data structures.
466 .RE

468 .SH NOTES
469 .sp
470 .LP
471 The \fBtcp\fR service is managed by the service management facility,
472 \fBsmf(5)\fR, under the service identifier:
473 .sp
474 .in +2
475 .nf
476 svc:/network/initial:default
477 .fi
478 .in -2
479 .sp

481 .sp
482 .LP
483 Administrative actions on this service, such as enabling, disabling, or
484 requesting restart, can be performed using \fBsvcadm\fR(1M). The service's
485 status can be queried using the \fBsvcs\fR(1) command.

```