# Linux Cross Reference

## Free Electrons

## Embedded Linux Experts

• *source navigation* • diff markup • identifier search • freetext search •

Version: 2.0.40 2.2.26 2.4.37 3.1 3.2 3.3 3.4 3.5 3.6 3.7 3.8 3.9 3.10 3.11 3.12 3.13 3.14 3.15 3.16 *3.17*

# Linux/net/ipv4/tcp_minisocks.c

```
1  /*
2   * INET         An implementation of the TCP/IP protocol suite for the LINUX
3   *              operating system.  INET is implemented using the  BSD Socket
4   *              interface as the means of communication with the user level.
5   *
6   *              Implementation of the Transmission Control Protocol(TCP).
7   *
8   * Authors:     Ross Biro
9   *              Fred N. van Kempen, <waltje@uWalt.NL.Mugnet.ORG>
10  *              Mark Evans, <evansmp@uhura.aston.ac.uk>
11  *              Corey Minyard <wf-rch!minyard@relay.EU.net>
12  *              Florian La Roche, <flla@stud.uni-sb.de>
13  *              Charles Hedrick, <hedrick@klinzhai.rutgers.edu>
14  *              Linus Torvalds, <torvalds@cs.helsinki.fi>
15  *              Alan Cox, <gw4pts@gw4pts.ampr.org>
16  *              Matthew Dillon, <dillon@apollo.west.oic.com>
17  *              Arnt Gulbrandsen, <agulbra@nvg.unit.no>
18  *              Jorge Cwik, <jorge@laser.satlink.net>
19  */
20
21  #include <linux/mm.h>
22  #include <linux/module.h>
23  #include <linux/slab.h>
24  #include <linux/sysctl.h>
25  #include <linux/workqueue.h>
26  #include <net/tcp.h>
27  #include <net/inet_common.h>
28  #include <net/xfrm.h>
29
30  int sysctl_tcp_syncookies __read_mostly = 1;
31  EXPORT_SYMBOL(sysctl_tcp_syncookies);
32
33  int sysctl_tcp_abort_on_overflow __read_mostly;
34
35  struct inet_timewait_death_row tcp_death_row = {
36          .sysctl_max_tw_buckets = NR_FILE * 2,
37          .period         = TCP_TIMEWAIT_LEN / INET_TWDR_TWKILL_SLOTS,
38          .death_lock     = __SPIN_LOCK_UNLOCKED(tcp_death_row.death_lock),
39          .hashinfo       = &tcp_hashinfo,
40          .tw_timer       = TIMER_INITIALIZER(inet_twdr_hangman, 0,
41                                              (unsigned long)&tcp_death_row),
42          .twkill_work    = __WORK_INITIALIZER(tcp_death_row.twkill_work,
43                                              inet_twdr_twkill_work),
44  /* Short-time timewait calendar */
45
46          .twcal_hand     = -1,
47          .twcal_timer    = TIMER_INITIALIZER(inet_twdr_twcal_tick, 0,
48                                              (unsigned long)&tcp_death_row),
49  };
50  EXPORT_SYMBOL_GPL(tcp_death_row);
51
52  static bool tcp_in_window(u32 seq, u32 end_seq, u32 s_win, u32 e_win)
53  {
54          if (seq == s_win)
55                  return true;
56          if (after(end_seq, s_win) && before(seq, e_win))
57                  return true;
58          return seq == e_win && seq == end_seq;
```

```
 59 }
 60
 61 /*
 62  * * Main purpose of TIME-WAIT state is to close connection gracefully,
 63  *   when one of ends sits in LAST-ACK or CLOSING retransmitting FIN
 64  *   (and, probably, tail of data) and one or more our ACKs are lost.
 65  * * What is TIME-WAIT timeout? It is associated with maximal packet
 66  *   lifetime in the internet, which results in wrong conclusion, that
 67  *   it is set to catch "old duplicate segments" wandering out of their path.
 68  *   It is not quite correct. This timeout is calculated so that it exceeds
 69  *   maximal retransmission timeout enough to allow to lose one (or more)
 70  *   segments sent by peer and our ACKs. This time may be calculated from RTO.
 71  * * When TIME-WAIT socket receives RST, it means that another end
 72  *   finally closed and we are allowed to kill TIME-WAIT too.
 73  * * Second purpose of TIME-WAIT is catching old duplicate segments.
 74  *   Well, certainly it is pure paranoia, but if we load TIME-WAIT
 75  *   with this semantics, we MUST NOT kill TIME-WAIT state with RSTs.
 76  * * If we invented some more clever way to catch duplicates
 77  *   (f.e. based on PAWS), we could truncate TIME-WAIT to several RTOs.
 78  *
 79  * The algorithm below is based on FORMAL INTERPRETATION of RFCs.
 80  * When you compare it to RFCs, please, read section SEGMENT ARRIVES
 81  * from the very beginning.
 82  *
 83  * NOTE. With recycling (and later with fin-wait-2) TW bucket
 84  * is _not_ stateless. It means, that strictly speaking we must
 85  * spinlock it. I do not want! Well, probability of misbehaviour
 86  * is ridiculously low and, seems, we could use some mb() tricks
 87  * to avoid misread sequence numbers, states etc.  --ANK
 88  *
 89  * We don't need to initialize tmp_out.sack_ok as we don't use the results
 90  */
 91 enum tcp_tw_status
 92 tcp_timewait_state_process(struct inet_timewait_sock *tw, struct sk_buff *skb,
 93                            const struct tcphdr *th)
 94 {
 95         struct tcp_options_received tmp_opt;
 96         struct tcp_timewait_sock *tcptw = tcp_twsk((struct sock *)tw);
 97         bool paws_reject = false;
 98
 99         tmp_opt.saw_tstamp = 0;
100         if (th->doff > (sizeof(*th) >> 2) && tcptw->tw_ts_recent_stamp) {
101                 tcp_parse_options(skb, &tmp_opt, 0, NULL);
102
103                 if (tmp_opt.saw_tstamp) {
104                         tmp_opt.rcv_tsecr       -= tcptw->tw_ts_offset;
105                         tmp_opt.ts_recent       = tcptw->tw_ts_recent;
106                         tmp_opt.ts_recent_stamp = tcptw->tw_ts_recent_stamp;
107                         paws_reject = tcp_paws_reject(&tmp_opt, th->rst);
108                 }
109         }
110
111         if (tw->tw_substate == TCP_FIN_WAIT2) {
112                 /* Just repeat all the checks of tcp_rcv_state_process() */
113
114                 /* Out of window, send ACK */
115                 if (paws_reject ||
116                     !tcp_in_window(TCP_SKB_CB(skb)->seq, TCP_SKB_CB(skb)->end_seq,
117                                    tcptw->tw_rcv_nxt,
118                                    tcptw->tw_rcv_nxt + tcptw->tw_rcv_wnd))
119                         return TCP_TW_ACK;
120
121                 if (th->rst)
122                         goto kill;
123
124                 if (th->syn && !before(TCP_SKB_CB(skb)->seq, tcptw->tw_rcv_nxt))
125                         goto kill_with_rst;
126
127                 /* Dup ACK? */
128                 if (!th->ack ||
129                     !after(TCP_SKB_CB(skb)->end_seq, tcptw->tw_rcv_nxt) ||
130                     TCP_SKB_CB(skb)->end_seq == TCP_SKB_CB(skb)->seq) {
131                         inet_twsk_put(tw);
132                         return TCP_TW_SUCCESS;
133                 }
134
135                 /* New data or FIN. If new data arrive after half-duplex close,
136                  * reset.
```

```
137                        */
138                    if (!th->fin ||
139                        TCP_SKB_CB(skb)->end_seq != tcptw->tw_rcv_nxt + 1) {
140 kill_with_rst:
141                            inet_twsk_deschedule(tw, &tcp_death_row);
142                            inet_twsk_put(tw);
143                            return TCP_TW_RST;
144                    }
145
146                    /* FIN arrived, enter true time-wait state. */
147                    tw->tw_substate   = TCP_TIME_WAIT;
148                    tcptw->tw_rcv_nxt = TCP_SKB_CB(skb)->end_seq;
149                    if (tmp_opt.saw_tstamp) {
150                            tcptw->tw_ts_recent_stamp = get_seconds();
151                            tcptw->tw_ts_recent       = tmp_opt.rcv_tsval;
152                    }
153
154                    if (tcp_death_row.sysctl_tw_recycle &&
155                        tcptw->tw_ts_recent_stamp &&
156                        tcp_tw_remember_stamp(tw))
157                            inet_twsk_schedule(tw, &tcp_death_row, tw->tw_timeout,
158                                            TCP_TIMEWAIT_LEN);
159                    else
160                            inet_twsk_schedule(tw, &tcp_death_row, TCP_TIMEWAIT_LEN,
161                                            TCP_TIMEWAIT_LEN);
162                    return TCP_TW_ACK;
163            }
164
165        /*
166         *       Now real TIME-WAIT state.
167         *
168         *       RFC 1122:
169         *       "When a connection is [...] on TIME-WAIT state [...]
170         *       [a TCP] MAY accept a new SYN from the remote TCP to
171         *       reopen the connection directly, if it:
172         *
173         *       (1)  assigns its initial sequence number for the new
174         *       connection to be larger than the largest sequence
175         *       number it used on the previous connection incarnation,
176         *       and
177         *
178         *       (2)  returns to TIME-WAIT state if the SYN turns out
179         *       to be an old duplicate".
180         */
181
182        if (!paws_reject &&
183            (TCP_SKB_CB(skb)->seq == tcptw->tw_rcv_nxt &&
184             (TCP_SKB_CB(skb)->seq == TCP_SKB_CB(skb)->end_seq || th->rst))) {
185                    /* In window segment, it may be only reset or bare ack. */
186
187                    if (th->rst) {
188                            /* This is TIME_WAIT assassination, in two flavors.
189                             * Oh well... nobody has a sufficient solution to this
190                             * protocol bug yet.
191                             */
192                            if (sysctl_tcp_rfc1337 == 0) {
193 kill:
194                                    inet_twsk_deschedule(tw, &tcp_death_row);
195                                    inet_twsk_put(tw);
196                                    return TCP_TW_SUCCESS;
197                            }
198                    }
199                    inet_twsk_schedule(tw, &tcp_death_row, TCP_TIMEWAIT_LEN,
200                                    TCP_TIMEWAIT_LEN);
201
202                    if (tmp_opt.saw_tstamp) {
203                            tcptw->tw_ts_recent       = tmp_opt.rcv_tsval;
204                            tcptw->tw_ts_recent_stamp = get_seconds();
205                    }
206
207                    inet_twsk_put(tw);
208                    return TCP_TW_SUCCESS;
209            }
210
211        /* Out of window segment.
212
213            All the segments are ACKed immediately.
214
```

```
215                  The only exception is new SYN. We accept it, if it is
216                  not old duplicate and we are not in danger to be killed
217                  by delayed old duplicates. RFC check is that it has
218                  newer sequence number works at rates <40Mbit/sec.
219                  However, if paws works, it is reliable AND even more,
220                  we even may relax silly seq space cutoff.
221
222                  RED-PEN: we violate main RFC requirement, if this SYN will appear
223                  old duplicate (i.e. we receive RST in reply to SYN-ACK),
224                  we must return socket to time-wait state. It is not good,
225                  but not fatal yet.
226           */
227
228          if (th->syn && !th->rst && !th->ack && !paws_reject &&
229              (after(TCP_SKB_CB(skb)->seq, tcptw->tw_rcv_nxt) ||
230               (tmp_opt.saw_tstamp &&
231                (s32)(tcptw->tw_ts_recent - tmp_opt.rcv_tsval) < 0))) {
232                  u32 isn = tcptw->tw_snd_nxt + 65535 + 2;
233                  if (isn == 0)
234                          isn++;
235                  TCP_SKB_CB(skb)->when = isn;
236                  return TCP_TW_SYN;
237          }
238
239          if (paws_reject)
240                  NET_INC_STATS_BH(twsk_net(tw), LINUX_MIB_PAWSESTABREJECTED);
241
242          if (!th->rst) {
243                  /* In this case we must reset the TIMEWAIT timer.
244                   *
245                   * If it is ACKless SYN it may be both old duplicate
246                   * and new good SYN with random sequence number <rcv_nxt.
247                   * Do not reschedule in the last case.
248                   */
249                  if (paws_reject || th->ack)
250                          inet_twsk_schedule(tw, &tcp_death_row, TCP_TIMEWAIT_LEN,
251                                             TCP_TIMEWAIT_LEN);
252
253                  /* Send ACK. Note, we do not put the bucket,
254                   * it will be released by caller.
255                   */
256                  return TCP_TW_ACK;
257          }
258          inet_twsk_put(tw);
259          return TCP_TW_SUCCESS;
260 }
261 EXPORT_SYMBOL(tcp_timewait_state_process);
262
263 /*
264  * Move a socket to time-wait or dead fin-wait-2 state.
265  */
266 void tcp_time_wait(struct sock *sk, int state, int timeo)
267 {
268          struct inet_timewait_sock *tw = NULL;
269          const struct inet_connection_sock *icsk = inet_csk(sk);
270          const struct tcp_sock *tp = tcp_sk(sk);
271          bool recycle_ok = false;
272
273          if (tcp_death_row.sysctl_tw_recycle && tp->rx_opt.ts_recent_stamp)
274                  recycle_ok = tcp_remember_stamp(sk);
275
276          if (tcp_death_row.tw_count < tcp_death_row.sysctl_max_tw_buckets)
277                  tw = inet_twsk_alloc(sk, state);
278
279          if (tw != NULL) {
280                  struct tcp_timewait_sock *tcptw = tcp_twsk((struct sock *)tw);
281                  const int rto = (icsk->icsk_rto << 2) - (icsk->icsk_rto >> 1);
282                  struct inet_sock *inet = inet_sk(sk);
283
284                  tw->tw_transparent      = inet->transparent;
285                  tw->tw_rcv_wscale       = tp->rx_opt.rcv_wscale;
286                  tcptw->tw_rcv_nxt       = tp->rcv_nxt;
287                  tcptw->tw_snd_nxt       = tp->snd_nxt;
288                  tcptw->tw_rcv_wnd       = tcp_receive_window(tp);
289                  tcptw->tw_ts_recent     = tp->rx_opt.ts_recent;
290                  tcptw->tw_ts_recent_stamp = tp->rx_opt.ts_recent_stamp;
291                  tcptw->tw_ts_offset     = tp->tsoffset;
292
```

```
293 #if IS_ENABLED(CONFIG_IPV6)
294             if (tw->tw_family == PF_INET6) {
295                     struct ipv6_pinfo *np = inet6_sk(sk);
296
297                     tw->tw_v6_daddr = sk->sk_v6_daddr;
298                     tw->tw_v6_rcv_saddr = sk->sk_v6_rcv_saddr;
299                     tw->tw_tclass = np->tclass;
300                     tw->tw_flowlabel = np->flow_label >> 12;
301                     tw->tw_ipv6only = sk->sk_ipv6only;
302             }
303 #endif
304
305 #ifdef CONFIG_TCP_MD5SIG
306             /*
307              * The timewait bucket does not have the key DB from the
308              * sock structure. We just make a quick copy of the
309              * md5 key being used (if indeed we are using one)
310              * so the timewait ack generating code has the key.
311              */
312             do {
313                     struct tcp_md5sig_key *key;
314                     tcptw->tw_md5_key = NULL;
315                     key = tp->af_specific->md5_lookup(sk, sk);
316                     if (key != NULL) {
317                             tcptw->tw_md5_key = kmemdup(key, sizeof(*key), GFP_ATOMIC);
318                             if (tcptw->tw_md5_key && !tcp_alloc_md5sig_pool())
319                                     BUG();
320                     }
321             } while (0);
322 #endif
323
324             /* Linkage updates. */
325             __inet_twsk_hashdance(tw, sk, &tcp_hashinfo);
326
327             /* Get the TIME_WAIT timeout firing. */
328             if (timeo < rto)
329                     timeo = rto;
330
331             if (recycle_ok) {
332                     tw->tw_timeout = rto;
333             } else {
334                     tw->tw_timeout = TCP_TIMEWAIT_LEN;
335                     if (state == TCP_TIME_WAIT)
336                             timeo = TCP_TIMEWAIT_LEN;
337             }
338
339             inet_twsk_schedule(tw, &tcp_death_row, timeo,
340                             TCP_TIMEWAIT_LEN);
341             inet_twsk_put(tw);
342     } else {
343             /* Sorry, if we're out of memory, just CLOSE this
344              * socket up.  We've got bigger problems than
345              * non-graceful socket closings.
346              */
347             NET_INC_STATS_BH(sock_net(sk), LINUX_MIB_TCPTIMEWAITOVERFLOW);
348     }
349
350     tcp_update_metrics(sk);
351     tcp_done(sk);
352 }
353
354 void tcp_twsk_destructor(struct sock *sk)
355 {
356 #ifdef CONFIG_TCP_MD5SIG
357     struct tcp_timewait_sock *twsk = tcp_twsk(sk);
358
359     if (twsk->tw_md5_key)
360             kfree_rcu(twsk->tw_md5_key, rcu);
361 #endif
362 }
363 EXPORT_SYMBOL_GPL(tcp_twsk_destructor);
364
365 void tcp_openreq_init_rwin(struct request_sock *req,
366                     struct sock *sk, struct dst_entry *dst)
367 {
368     struct inet_request_sock *ireq = inet_rsk(req);
369     struct tcp_sock *tp = tcp_sk(sk);
370     __u8 rcv_wscale;
```

```
371                int mss = dst_metric_advmss(dst);
372
373                if (tp->rx_opt.user_mss && tp->rx_opt.user_mss < mss)
374                        mss = tp->rx_opt.user_mss;
375
376        /* Set this up on the first call only */
377        req->window_clamp = tp->window_clamp ? : dst_metric(dst, RTAX_WINDOW);
378
379        /* limit the window selection if the user enforce a smaller rx buffer */
380        if (sk->sk_userlocks & SOCK_RCVBUF_LOCK &&
381            (req->window_clamp > tcp_full_space(sk) || req->window_clamp == 0))
382                req->window_clamp = tcp_full_space(sk);
383
384        /* tcp_full_space because it is guaranteed to be the first packet */
385        tcp_select_initial_window(tcp_full_space(sk),
386                mss - (ireq->tstamp_ok ? TCPOLEN_TSTAMP_ALIGNED : 0),
387                &req->rcv_wnd,
388                &req->window_clamp,
389                ireq->wscale_ok,
390                &rcv_wscale,
391                dst_metric(dst, RTAX_INITRWND));
392        ireq->rcv_wscale = rcv_wscale;
393 }
394 EXPORT_SYMBOL(tcp_openreq_init_rwin);
395
396 static inline void TCP_ECN_openreq_child(struct tcp_sock *tp,
397                                          struct request_sock *req)
398 {
399        tp->ecn_flags = inet_rsk(req)->ecn_ok ? TCP_ECN_OK : 0;
400 }
401
402 /* This is not only more efficient than what we used to do, it eliminates
403  * a lot of code duplication between IPv4/IPv6 SYN recv processing. -DaveM
404  *
405  * Actually, we could lots of memory writes here. tp of listening
406  * socket contains all necessary default parameters.
407  */
408 struct sock *tcp_create_openreq_child(struct sock *sk, struct request_sock *req, struct sk_buff *skb)
409 {
410        struct sock *newsk = inet_csk_clone_lock(sk, req, GFP_ATOMIC);
411
412        if (newsk != NULL) {
413                const struct inet_request_sock *ireq = inet_rsk(req);
414                struct tcp_request_sock *treq = tcp_rsk(req);
415                struct inet_connection_sock *newicsk = inet_csk(newsk);
416                struct tcp_sock *newtp = tcp_sk(newsk);
417
418                /* Now setup tcp_sock */
419                newtp->pred_flags = 0;
420
421                newtp->rcv_wup = newtp->copied_seq =
422                newtp->rcv_nxt = treq->rcv_isn + 1;
423
424                newtp->snd_sml = newtp->snd_una =
425                newtp->snd_nxt = newtp->snd_up = treq->snt_isn + 1;
426
427                tcp_prequeue_init(newtp);
428                INIT_LIST_HEAD(&newtp->tsq_node);
429
430                tcp_init_wl(newtp, treq->rcv_isn);
431
432                newtp->srtt_us = 0;
433                newtp->mdev_us = jiffies_to_usecs(TCP_TIMEOUT_INIT);
434                newicsk->icsk_rto = TCP_TIMEOUT_INIT;
435
436                newtp->packets_out = 0;
437                newtp->retrans_out = 0;
438                newtp->sacked_out = 0;
439                newtp->fackets_out = 0;
440                newtp->snd_ssthresh = TCP_INFINITE_SSTHRESH;
441                tcp_enable_early_retrans(newtp);
442                newtp->tlp_high_seq = 0;
443                newtp->lsndtime = treq->snt_synack;
444                newtp->total_retrans = req->num_retrans;
445
446                /* So many TCP implementations out there (incorrectly) count the
447                 * initial SYN frame in their delayed-ACK and congestion control
448                 * algorithms that we must have the following bandaid to talk
```

```
449                 * efficiently to them.  -DaveM
450                 */
451                newtp->snd_cwnd = TCP_INIT_CWND;
452                newtp->snd_cwnd_cnt = 0;
453
454                if (newicsk->icsk_ca_ops != &tcp_init_congestion_ops &&
455                    !try_module_get(newicsk->icsk_ca_ops->owner))
456                        newicsk->icsk_ca_ops = &tcp_init_congestion_ops;
457
458                tcp_set_ca_state(newsk, TCP_CA_Open);
459                tcp_init_xmit_timers(newsk);
460                __skb_queue_head_init(&newtp->out_of_order_queue);
461                newtp->write_seq = newtp->pushed_seq = treq->snt_isn + 1;
462
463                newtp->rx_opt.saw_tstamp = 0;
464
465                newtp->rx_opt.dsack = 0;
466                newtp->rx_opt.num_sacks = 0;
467
468                newtp->urg_data = 0;
469
470                if (sock_flag(newsk, SOCK_KEEPOPEN))
471                        inet_csk_reset_keepalive_timer(newsk,
472                                                keepalive_time_when(newtp));
473
474                newtp->rx_opt.tstamp_ok = ireq->tstamp_ok;
475                if ((newtp->rx_opt.sack_ok = ireq->sack_ok) != 0) {
476                        if (sysctl_tcp_fack)
477                                tcp_enable_fack(newtp);
478                }
479                newtp->window_clamp = req->window_clamp;
480                newtp->rcv_ssthresh = req->rcv_wnd;
481                newtp->rcv_wnd = req->rcv_wnd;
482                newtp->rx_opt.wscale_ok = ireq->wscale_ok;
483                if (newtp->rx_opt.wscale_ok) {
484                        newtp->rx_opt.snd_wscale = ireq->snd_wscale;
485                        newtp->rx_opt.rcv_wscale = ireq->rcv_wscale;
486                } else {
487                        newtp->rx_opt.snd_wscale = newtp->rx_opt.rcv_wscale = 0;
488                        newtp->window_clamp = min(newtp->window_clamp, 65535U);
489                }
490                newtp->snd_wnd = (ntohs(tcp_hdr(skb)->window) <<
491                                        newtp->rx_opt.snd_wscale);
492                newtp->max_window = newtp->snd_wnd;
493
494                if (newtp->rx_opt.tstamp_ok) {
495                        newtp->rx_opt.ts_recent = req->ts_recent;
496                        newtp->rx_opt.ts_recent_stamp = get_seconds();
497                        newtp->tcp_header_len = sizeof(struct tcphdr) + TCPOLEN_TSTAMP_ALIGNED;
498                } else {
499                        newtp->rx_opt.ts_recent_stamp = 0;
500                        newtp->tcp_header_len = sizeof(struct tcphdr);
501                }
502                newtp->tsoffset = 0;
503 #ifdef CONFIG_TCP_MD5SIG
504                newtp->md5sig_info = NULL;      /*XXX*/
505                if (newtp->af_specific->md5_lookup(sk, newsk))
506                        newtp->tcp_header_len += TCPOLEN_MD5SIG_ALIGNED;
507 #endif
508                if (skb->len >= TCP_MSS_DEFAULT + newtp->tcp_header_len)
509                        newicsk->icsk_ack.last_seg_size = skb->len - newtp->tcp_header_len;
510                newtp->rx_opt.mss_clamp = req->mss;
511                TCP_ECN_openreq_child(newtp, req);
512                newtp->fastopen_rsk = NULL;
513                newtp->syn_data_acked = 0;
514
515                TCP_INC_STATS_BH(sock_net(sk), TCP_MIB_PASSIVEOPENS);
516        }
517        return newsk;
518 }
519 EXPORT_SYMBOL(tcp_create_openreq_child);
520
521 /*
522  * Process an incoming packet for SYN_RECV sockets represented as a
523  * request_sock. Normally sk is the listener socket but for TFO it
524  * points to the child socket.
525  *
526  * XXX (TFO) - The current impl contains a special check for ack
```

```
527     * validation and inside tcp_v4_reqsk_send_ack(). Can we do better?
528     *
529     * We don't need to initialize tmp_opt.sack_ok as we don't use the results
530     */
531
532    struct sock *tcp_check_req(struct sock *sk, struct sk_buff *skb,
533                               struct request_sock *req,
534                               struct request_sock **prev,
535                               bool fastopen)
536    {
537            struct tcp_options_received tmp_opt;
538            struct sock *child;
539            const struct tcphdr *th = tcp_hdr(skb);
540            __be32 flg = tcp_flag_word(th) & (TCP_FLAG_RST|TCP_FLAG_SYN|TCP_FLAG_ACK);
541            bool paws_reject = false;
542
543            BUG_ON(fastopen == (sk->sk_state == TCP_LISTEN));
544
545            tmp_opt.saw_tstamp = 0;
546            if (th->doff > (sizeof(struct tcphdr)>>2)) {
547                    tcp_parse_options(skb, &tmp_opt, 0, NULL);
548
549                    if (tmp_opt.saw_tstamp) {
550                            tmp_opt.ts_recent = req->ts_recent;
551                            /* We do not store true stamp, but it is not required,
552                             * it can be estimated (approximately)
553                             * from another data.
554                             */
555                            tmp_opt.ts_recent_stamp = get_seconds() - ((TCP_TIMEOUT_INIT/HZ)<<req->num_timeout);
556                            paws_reject = tcp_paws_reject(&tmp_opt, th->rst);
557                    }
558            }
559
560            /* Check for pure retransmitted SYN. */
561            if (TCP_SKB_CB(skb)->seq == tcp_rsk(req)->rcv_isn &&
562                flg == TCP_FLAG_SYN &&
563                !paws_reject) {
564                    /*
565                     * RFC793 draws (Incorrectly! It was fixed in RFC1122)
566                     * this case on figure 6 and figure 8, but formal
567                     * protocol description says NOTHING.
568                     * To be more exact, it says that we should send ACK,
569                     * because this segment (at least, if it has no data)
570                     * is out of window.
571                     *
572                     *  CONCLUSION: RFC793 (even with RFC1122) DOES NOT
573                     *  describe SYN-RECV state. All the description
574                     *  is wrong, we cannot believe to it and should
575                     *  rely only on common sense and implementation
576                     *  experience.
577                     *
578                     * Enforce "SYN-ACK" according to figure 8, figure 6
579                     * of RFC793, fixed by RFC1122.
580                     *
581                     * Note that even if there is new data in the SYN packet
582                     * they will be thrown away too.
583                     *
584                     * Reset timer after retransmitting SYNACK, similar to
585                     * the idea of fast retransmit in recovery.
586                     */
587                    if (!inet_rtx_syn_ack(sk, req))
588                            req->expires = min(TCP_TIMEOUT_INIT << req->num_timeout,
589                                               TCP_RTO_MAX) + jiffies;
590                    return NULL;
591            }
592
593            /* Further reproduces section "SEGMENT ARRIVES"
594               for state SYN-RECEIVED of RFC793.
595               It is broken, however, it does not work only
596               when SYNs are crossed.
597
598               You would think that SYN crossing is impossible here, since
599               we should have a SYN_SENT socket (from connect()) on our end,
600               but this is not true if the crossed SYNs were sent to both
601               ends by a malicious third party.  We must defend against this,
602               and to do that we first verify the ACK (as per RFC793, page
603               36) and reset if it is invalid.  Is this a true full defense?
604               To convince ourselves, let us consider a way in which the ACK
```

```
605                 test can still pass in this 'malicious crossed SYNs' case.
606                 Malicious sender sends identical SYNs (and thus identical sequence
607                 numbers) to both A and B:
608
609                     A: gets SYN, seq=7
610                     B: gets SYN, seq=7
611
612                 By our good fortune, both A and B select the same initial
613                 send sequence number of seven :-)
614
615                     A: sends SYN|ACK, seq=7, ack_seq=8
616                     B: sends SYN|ACK, seq=7, ack_seq=8
617
618                 So we are now A eating this SYN|ACK, ACK test passes.  So
619                 does sequence test, SYN is truncated, and thus we consider
620                 it a bare ACK.
621
622                 If icsk->icsk_accept_queue.rskq_defer_accept, we silently drop this
623                 bare ACK.  Otherwise, we create an established connection.  Both
624                 ends (listening sockets) accept the new incoming connection and try
625                 to talk to each other. 8-)
626
627                 Note: This case is both harmless, and rare.  Possibility is about the
628                 same as us discovering intelligent life on another plant tomorrow.
629
630                 But generally, we should (RFC lies!) to accept ACK
631                 from SYNACK both here and in tcp_rcv_state_process().
632                 tcp_rcv_state_process() does not, hence, we do not too.
633
634                 Note that the case is absolutely generic:
635                 we cannot optimize anything here without
636                 violating protocol. All the checks must be made
637                 before attempt to create socket.
638             */
639
640         /* RFC793 page 36: "If the connection is in any non-synchronized state ...
641          *                 and the incoming segment acknowledges something not yet
642          *                 sent (the segment carries an unacceptable ACK) ...
643          *                 a reset is sent."
644          *
645          * Invalid ACK: reset will be sent by listening socket.
646          * Note that the ACK validity check for a Fast Open socket is done
647          * elsewhere and is checked directly against the child socket rather
648          * than req because user data may have been sent out.
649          */
650         if ((flg & TCP_FLAG_ACK) && !fastopen &&
651             (TCP_SKB_CB(skb)->ack_seq !=
652              tcp_rsk(req)->snt_isn + 1))
653                 return sk;
654
655         /* Also, it would be not so bad idea to check rcv_tsecr, which
656          * is essentially ACK extension and too early or too late values
657          * should cause reset in unsynchronized states.
658          */
659
660         /* RFC793: "first check sequence number". */
661
662         if (paws_reject || !tcp_in_window(TCP_SKB_CB(skb)->seq, TCP_SKB_CB(skb)->end_seq,
663                                  tcp_rsk(req)->rcv_nxt, tcp_rsk(req)->rcv_nxt + req->rcv_wnd)) {
664                 /* Out of window: send ACK and drop. */
665                 if (!(flg & TCP_FLAG_RST))
666                         req->rsk_ops->send_ack(sk, skb, req);
667                 if (paws_reject)
668                         NET_INC_STATS_BH(sock_net(sk), LINUX_MIB_PAWSESTABREJECTED);
669                 return NULL;
670         }
671
672         /* In sequence, PAWS is OK. */
673
674         if (tmp_opt.saw_tstamp && !after(TCP_SKB_CB(skb)->seq, tcp_rsk(req)->rcv_nxt))
675                 req->ts_recent = tmp_opt.rcv_tsval;
676
677         if (TCP_SKB_CB(skb)->seq == tcp_rsk(req)->rcv_isn) {
678                 /* Truncate SYN, it is out of window starting
679                    at tcp_rsk(req)->rcv_isn + 1. */
680                 flg &= ~TCP_FLAG_SYN;
681         }
682
```

```
683              /* RFC793: "second check the RST bit" and
684               *          "fourth, check the SYN bit"
685               */
686              if (flg & (TCP_FLAG_RST|TCP_FLAG_SYN)) {
687                      TCP_INC_STATS_BH(sock_net(sk), TCP_MIB_ATTEMPTFAILS);
688                      goto embryonic_reset;
689              }
690
691              /* ACK sequence verified above, just make sure ACK is
692               * set.  If ACK not set, just silently drop the packet.
693               *
694               * XXX (TFO) - if we ever allow "data after SYN", the
695               * following check needs to be removed.
696               */
697              if (!(flg & TCP_FLAG_ACK))
698                      return NULL;
699
700              /* For Fast Open no more processing is needed (sk is the
701               * child socket).
702               */
703              if (fastopen)
704                      return sk;
705
706              /* While TCP_DEFER_ACCEPT is active, drop bare ACK. */
707              if (req->num_timeout < inet_csk(sk)->icsk_accept_queue.rskq_defer_accept &&
708                  TCP_SKB_CB(skb)->end_seq == tcp_rsk(req)->rcv_isn + 1) {
709                      inet_rsk(req)->acked = 1;
710                      NET_INC_STATS_BH(sock_net(sk), LINUX_MIB_TCPDEFERACCEPTDROP);
711                      return NULL;
712              }
713
714              /* OK, ACK is valid, create big socket and
715               * feed this segment to it. It will repeat all
716               * the tests. THIS SEGMENT MUST MOVE SOCKET TO
717               * ESTABLISHED STATE. If it will be dropped after
718               * socket is created, wait for troubles.
719               */
720              child = inet_csk(sk)->icsk_af_ops->syn_recv_sock(sk, skb, req, NULL);
721              if (child == NULL)
722                      goto listen_overflow;
723
724              inet_csk_reqsk_queue_unlink(sk, req, prev);
725              inet_csk_reqsk_queue_removed(sk, req);
726
727              inet_csk_reqsk_queue_add(sk, req, child);
728              return child;
729
730      listen_overflow:
731              if (!sysctl_tcp_abort_on_overflow) {
732                      inet_rsk(req)->acked = 1;
733                      return NULL;
734              }
735
736      embryonic_reset:
737              if (!(flg & TCP_FLAG_RST)) {
738                      /* Received a bad SYN pkt - for TFO We try not to reset
739                       * the local connection unless it's really necessary to
740                       * avoid becoming vulnerable to outside attack aiming at
741                       * resetting legit local connections.
742                       */
743                      req->rsk_ops->send_reset(sk, skb);
744              } else if (fastopen) { /* received a valid RST pkt */
745                      reqsk_fastopen_remove(sk, req, true);
746                      tcp_reset(sk);
747              }
748              if (!fastopen) {
749                      inet_csk_reqsk_queue_drop(sk, req, prev);
750                      NET_INC_STATS_BH(sock_net(sk), LINUX_MIB_EMBRYONICRSTS);
751              }
752              return NULL;
753      }
754      EXPORT_SYMBOL(tcp_check_req);
755
756      /*
757       * Queue segment on the new socket if the new socket is active,
758       * otherwise we just shortcircuit this and continue with
759       * the new socket.
760       *
```

```
761   * For the vast majority of cases child->sk_state will be TCP_SYN_RECV
762   * when entering. But other states are possible due to a race condition
763   * where after __inet_lookup_established() fails but before the listener
764   * locked is obtained, other packets cause the same connection to
765   * be created.
766   */
767
768  int tcp_child_process(struct sock *parent, struct sock *child,
769                        struct sk_buff *skb)
770  {
771          int ret = 0;
772          int state = child->sk_state;
773
774          if (!sock_owned_by_user(child)) {
775                  ret = tcp_rcv_state_process(child, skb, tcp_hdr(skb),
776                                              skb->len);
777                  /* Wakeup parent, send SIGIO */
778                  if (state == TCP_SYN_RECV && child->sk_state != state)
779                          parent->sk_data_ready(parent);
780          } else {
781                  /* Alas, it is possible again, because we do lookup
782                   * in main socket hash table and lock on listening
783                   * socket does not protect us more.
784                   */
785                  __sk_add_backlog(child, skb);
786          }
787
788          bh_unlock_sock(child);
789          sock_put(child);
790          return ret;
791  }
792  EXPORT_SYMBOL(tcp_child_process);
793
```

This page was automatically generated by [LXR](#) 0.3.1 ([source](#)). • Linux is a registered trademark of Linus Torvalds •
Contact us

- Home
- Development
- Services
- Training
- Docs
- Community
- Company
- Blog