# 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_lp.c](#)

```
1  /*
2   * TCP Low Priority (TCP-LP)
3   *
4   * TCP Low Priority is a distributed algorithm whose goal is to utilize only
5   *   the excess network bandwidth as compared to the ``fair share`` of
6   *   bandwidth as targeted by TCP.
7   *
8   * As of 2.6.13, Linux supports pluggable congestion control algorithms.
9   * Due to the limitation of the API, we take the following changes from
10  * the original TCP-LP implementation:
11  *   o We use newReno in most core CA handling. Only add some checking
12  *     within cong_avoid.
13  *   o Error correcting in remote HZ, therefore remote HZ will be kept
14  *     on checking and updating.
15  *   o Handling calculation of One-Way-Delay (OWD) within rtt_sample, since
16  *     OWD have a similar meaning as RTT. Also correct the buggy formular.
17  *   o Handle reaction for Early Congestion Indication (ECI) within
18  *     pkts_acked, as mentioned within pseudo code.
19  *   o OWD is handled in relative format, where local time stamp will in
20  *     tcp_time_stamp format.
21  *
22  * Original Author:
23  *   Aleksandar Kuzmanovic <akuzma@northwestern.edu>
24  * Available from:
25  *     http://www.ece.rice.edu/~akuzma/Doc/akuzma/TCP-LP.pdf
26  * Original implementation for 2.4.19:
27  *     http://www-ece.rice.edu/networks/TCP-LP/
28  *
29  * 2.6.x module Authors:
30  *   Wong Hoi Sing, Edison <hswong3i@gmail.com>
31  *   Hung Hing Lun, Mike <hlhung3i@gmail.com>
32  * SourceForge project page:
33  *     http://tcp-lp-mod.sourceforge.net/
34  */
35
36 #include <linux/module.h>
37 #include <net/tcp.h>
38
39 /* resolution of owd */
40 #define LP_RESOL        1000
41
42 /**
43  * enum tcp_lp_state
```

```
 44    * @LP_VALID_RHZ: is remote HZ valid?
 45    * @LP_VALID_OWD: is OWD valid?
 46    * @LP_WITHIN_THR: are we within threshold?
 47    * @LP_WITHIN_INF: are we within inference?
 48    *
 49    * TCP-LP's state flags.
 50    * We create this set of state flag mainly for debugging.
 51    */
 52  enum tcp_lp_state {
 53          LP_VALID_RHZ = (1 << 0),
 54          LP_VALID_OWD = (1 << 1),
 55          LP_WITHIN_THR = (1 << 3),
 56          LP_WITHIN_INF = (1 << 4),
 57  };
 58
 59  /**
 60   * struct lp
 61   * @flag: TCP-LP state flag
 62   * @sowd: smoothed OWD << 3
 63   * @owd_min: min OWD
 64   * @owd_max: max OWD
 65   * @owd_max_rsv: resrved max owd
 66   * @remote_hz: estimated remote HZ
 67   * @remote_ref_time: remote reference time
 68   * @local_ref_time: local reference time
 69   * @last_drop: time for last active drop
 70   * @inference: current inference
 71   *
 72   * TCP-LP's private struct.
 73   * We get the idea from original TCP-LP implementation where only left those we
 74   * found are really useful.
 75   */
 76  struct lp {
 77          u32 flag;
 78          u32 sowd;
 79          u32 owd_min;
 80          u32 owd_max;
 81          u32 owd_max_rsv;
 82          u32 remote_hz;
 83          u32 remote_ref_time;
 84          u32 local_ref_time;
 85          u32 last_drop;
 86          u32 inference;
 87  };
 88
 89  /**
 90   * tcp_lp_init
 91   *
 92   * Init all required variables.
 93   * Clone the handling from Vegas module implementation.
 94   */
 95  static void tcp_lp_init(struct sock *sk)
 96  {
 97          struct lp *lp = inet_csk_ca(sk);
 98
 99          lp->flag = 0;
100          lp->sowd = 0;
101          lp->owd_min = 0xffffffff;
102          lp->owd_max = 0;
103          lp->owd_max_rsv = 0;
104          lp->remote_hz = 0;
105          lp->remote_ref_time = 0;
106          lp->local_ref_time = 0;
107          lp->last_drop = 0;
108          lp->inference = 0;
```

```
109  }
110
111  /**
112   * tcp_lp_cong_avoid
113   *
114   * Implementation of cong_avoid.
115   * Will only call newReno CA when away from inference.
116   * From TCP-LP's paper, this will be handled in additive increasement.
117   */
118  static void tcp_lp_cong_avoid(struct sock *sk, u32 ack, u32 acked)
119  {
120          struct lp *lp = inet_csk_ca(sk);
121
122          if (!(lp->flag & LP_WITHIN_INF))
123                  tcp_reno_cong_avoid(sk, ack, acked);
124  }
125
126  /**
127   * tcp_lp_remote_hz_estimator
128   *
129   * Estimate remote HZ.
130   * We keep on updating the estimated value, where original TCP-LP
131   * implementation only guest it for once and use forever.
132   */
133  static u32 tcp_lp_remote_hz_estimator(struct sock *sk)
134  {
135          struct tcp_sock *tp = tcp_sk(sk);
136          struct lp *lp = inet_csk_ca(sk);
137          s64 rhz = lp->remote_hz << 6;    /* remote HZ << 6 */
138          s64 m = 0;
139
140          /* not yet record reference time
141           * go away!! record it before come back!! */
142          if (lp->remote_ref_time == 0 || lp->local_ref_time == 0)
143                  goto out;
144
145          /* we can't calc remote HZ with no different!! */
146          if (tp->rx_opt.rcv_tsval == lp->remote_ref_time ||
147              tp->rx_opt.rcv_tsecr == lp->local_ref_time)
148                  goto out;
149
150          m = HZ * (tp->rx_opt.rcv_tsval -
151                      lp->remote_ref_time) / (tp->rx_opt.rcv_tsecr -
152                                              lp->local_ref_time);
153          if (m < 0)
154                  m = -m;
155
156          if (rhz > 0) {
157                  m -= rhz >> 6;  /* m is now error in remote HZ est */
158                  rhz += m;       /* 63/64 old + 1/64 new */
159          } else
160                  rhz = m << 6;
161
162    out:
163          /* record time for successful remote HZ calc */
164          if ((rhz >> 6) > 0)
165                  lp->flag |= LP_VALID_RHZ;
166          else
167                  lp->flag &= ~LP_VALID_RHZ;
168
169          /* record reference time stamp */
170          lp->remote_ref_time = tp->rx_opt.rcv_tsval;
171          lp->local_ref_time = tp->rx_opt.rcv_tsecr;
172
173          return rhz >> 6;
```

```
174  }
175
176  /**
177   * tcp_lp_owd_calculator
178   *
179   * Calculate one way delay (in relative format).
180   * Original implement OWD as minus of remote time difference to local time
181   * difference directly. As this time difference just simply equal to RTT, when
182   * the network status is stable, remote RTT will equal to local RTT, and result
183   * OWD into zero.
184   * It seems to be a bug and so we fixed it.
185   */
186  static u32 tcp_lp_owd_calculator(struct sock *sk)
187  {
188          struct tcp_sock *tp = tcp_sk(sk);
189          struct lp *lp = inet_csk_ca(sk);
190          s64 owd = 0;
191
192          lp->remote_hz = tcp_lp_remote_hz_estimator(sk);
193
194          if (lp->flag & LP_VALID_RHZ) {
195                  owd =
196                      tp->rx_opt.rcv_tsval * (LP_RESOL / lp->remote_hz) -
197                      tp->rx_opt.rcv_tsecr * (LP_RESOL / HZ);
198                  if (owd < 0)
199                          owd = -owd;
200          }
201
202          if (owd > 0)
203                  lp->flag |= LP_VALID_OWD;
204          else
205                  lp->flag &= ~LP_VALID_OWD;
206
207          return owd;
208  }
209
210  /**
211   * tcp_lp_rtt_sample
212   *
213   * Implementation or rtt_sample.
214   * Will take the following action,
215   *   1. calc OWD,
216   *   2. record the min/max OWD,
217   *   3. calc smoothed OWD (SOWD).
218   * Most ideas come from the original TCP-LP implementation.
219   */
220  static void tcp_lp_rtt_sample(struct sock *sk, u32 rtt)
221  {
222          struct lp *lp = inet_csk_ca(sk);
223          s64 mowd = tcp_lp_owd_calculator(sk);
224
225          /* sorry that we don't have valid data */
226          if (!(lp->flag & LP_VALID_RHZ) || !(lp->flag & LP_VALID_OWD))
227                  return;
228
229          /* record the next min owd */
230          if (mowd < lp->owd_min)
231                  lp->owd_min = mowd;
232
233          /* always forget the max of the max
234           * we just set owd_max as one below it */
235          if (mowd > lp->owd_max) {
236                  if (mowd > lp->owd_max_rsv) {
237                          if (lp->owd_max_rsv == 0)
238                                  lp->owd_max = mowd;
```

```
239                            else
240                                    lp->owd_max = lp->owd_max_rsv;
241                            lp->owd_max_rsv = mowd;
242                    } else
243                            lp->owd_max = mowd;
244            }
245
246            /* calc for smoothed owd */
247            if (lp->sowd != 0) {
248                    mowd -= lp->sowd >> 3;    /* m is now error in owd est */
249                    lp->sowd += mowd;          /* owd = 7/8 owd + 1/8 new */
250            } else
251                    lp->sowd = mowd << 3;    /* take the measured time be owd */
252  }
253
254  /**
255   * tcp_lp_pkts_acked
256   *
257   * Implementation of pkts_acked.
258   * Deal with active drop under Early Congestion Indication.
259   * Only drop to half and 1 will be handle, because we hope to use back
260   * newReno in increase case.
261   * We work it out by following the idea from TCP-LP's paper directly
262   */
263  static void tcp_lp_pkts_acked(struct sock *sk, u32 num_acked, s32 rtt_us)
264  {
265          struct tcp_sock *tp = tcp_sk(sk);
266          struct lp *lp = inet_csk_ca(sk);
267
268          if (rtt_us > 0)
269                  tcp_lp_rtt_sample(sk, rtt_us);
270
271          /* calc inference */
272          if (tcp_time_stamp > tp->rx_opt.rcv_tsecr)
273                  lp->inference = 3 * (tcp_time_stamp - tp->rx_opt.rcv_tsecr);
274
275          /* test if within inference */
276          if (lp->last_drop && (tcp_time_stamp - lp->last_drop < lp->inference))
277                  lp->flag |= LP_WITHIN_INF;
278          else
279                  lp->flag &= ~LP_WITHIN_INF;
280
281          /* test if within threshold */
282          if (lp->sowd >> 3 <
283              lp->owd_min + 15 * (lp->owd_max - lp->owd_min) / 100)
284                  lp->flag |= LP_WITHIN_THR;
285          else
286                  lp->flag &= ~LP_WITHIN_THR;
287
288          pr_debug("TCP-LP: %05o|%5u|%5u|%15u|%15u|%15u\n", lp->flag,
289                  tp->snd_cwnd, lp->remote_hz, lp->owd_min, lp->owd_max,
290                  lp->sowd >> 3);
291
292          if (lp->flag & LP_WITHIN_THR)
293                  return;
294
295          /* FIXME: try to reset owd_min and owd_max here
296           * so decrease the chance the min/max is no longer suitable
297           * and will usually within threshold when whithin inference */
298          lp->owd_min = lp->sowd >> 3;
299          lp->owd_max = lp->sowd >> 2;
300          lp->owd_max_rsv = lp->sowd >> 2;
301
302          /* happened within inference
303           * drop snd_cwnd into 1 */
```

```
304              if (lp->flag & LP_WITHIN_INF)
305                      tp->snd_cwnd = 1U;
306
307              /* happened after inference
308               * cut snd_cwnd into half */
309              else
310                      tp->snd_cwnd = max(tp->snd_cwnd >> 1U, 1U);
311
312              /* record this drop time */
313              lp->last_drop = tcp_time_stamp;
314 }
315
316 static struct tcp_congestion_ops tcp_lp __read_mostly = {
317         .init = tcp_lp_init,
318         .ssthresh = tcp_reno_ssthresh,
319         .cong_avoid = tcp_lp_cong_avoid,
320         .pkts_acked = tcp_lp_pkts_acked,
321
322         .owner = THIS_MODULE,
323         .name = "lp"
324 };
325
326 static int __init tcp_lp_register(void)
327 {
328         BUILD_BUG_ON(sizeof(struct lp) > ICSK_CA_PRIV_SIZE);
329         return tcp_register_congestion_control(&tcp_lp);
330 }
331
332 static void __exit tcp_lp_unregister(void)
333 {
334         tcp_unregister_congestion_control(&tcp_lp);
335 }
336
337 module_init(tcp_lp_register);
338 module_exit(tcp_lp_unregister);
339
340 MODULE_AUTHOR("Wong Hoi Sing Edison, Hung Hing Lun Mike");
341 MODULE_LICENSE("GPL");
342 MODULE_DESCRIPTION("TCP Low Priority");
343
```

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](#)