# 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_probe.c

```
1  /*
2   * tcpprobe - Observe the TCP flow with kprobes.
3   *
4   * The idea for this came from Werner Almesberger's umlsim
5   * Copyright (C) 2004, Stephen Hemminger <shemminger@osdl.org>
6   *
7   * This program is free software; you can redistribute it and/or modify
8   * it under the terms of the GNU General Public License as published by
9   * the Free Software Foundation; either version 2 of the License.
10  *
11  * This program is distributed in the hope that it will be useful,
12  * but WITHOUT ANY WARRANTY; without even the implied warranty of
13  * MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE.  See the
14  * GNU General Public License for more details.
15  *
16  * You should have received a copy of the GNU General Public License
17  * along with this program; if not, write to the Free Software
18  * Foundation, Inc., 675 Mass Ave, Cambridge, MA 02139, USA.
19  */
20
21  #define pr_fmt(fmt) KBUILD_MODNAME ": " fmt
22
23  #include <linux/kernel.h>
24  #include <linux/kprobes.h>
25  #include <linux/socket.h>
26  #include <linux/tcp.h>
27  #include <linux/slab.h>
28  #include <linux/proc_fs.h>
29  #include <linux/module.h>
30  #include <linux/ktime.h>
31  #include <linux/time.h>
32  #include <net/net_namespace.h>
33
34  #include <net/tcp.h>
35
36  MODULE_AUTHOR("Stephen Hemminger <shemminger@linux-foundation.org>");
37  MODULE_DESCRIPTION("TCP cwnd snooper");
38  MODULE_LICENSE("GPL");
39  MODULE_VERSION("1.1");
40
41  static int port __read_mostly;
42  MODULE_PARM_DESC(port, "Port to match (0=all)");
43  module_param(port, int, 0);
```

```
44
45  static unsigned int buffsize __read_mostly = 4096;
46  MODULE_PARM_DESC(bufsize, "Log buffer size in packets (4096)");
47  module_param(bufsize, uint, 0);
48
49  static unsigned int fwmark __read_mostly;
50  MODULE_PARM_DESC(fwmark, "skb mark to match (0=no mark)");
51  module_param(fwmark, uint, 0);
52
53  static int full __read_mostly;
54  MODULE_PARM_DESC(full, "Full log (1=every ack packet received,  0=only cwnd changes)");
55  module_param(full, int, 0);
56
57  static const char procname[] = "tcpprobe";
58
59  struct tcp_log {
60          ktime_t tstamp;
61          union {
62                  struct sockaddr         raw;
63                  struct sockaddr_in      v4;
64                  struct sockaddr_in6     v6;
65          } src, dst;
66          u16     length;
67          u32     snd_nxt;
68          u32     snd_una;
69          u32     snd_wnd;
70          u32     rcv_wnd;
71          u32     snd_cwnd;
72          u32     ssthresh;
73          u32     srtt;
74  };
75
76  static struct {
77          spinlock_t      lock;
78          wait_queue_head_t wait;
79          ktime_t         start;
80          u32             lastcwnd;
81
82          unsigned long   head, tail;
83          struct tcp_log  *log;
84  } tcp_probe;
85
86
87  static inline int tcp_probe_used(void)
88  {
89          return (tcp_probe.head - tcp_probe.tail) & (bufsize - 1);
90  }
91
92  static inline int tcp_probe_avail(void)
93  {
94          return bufsize - tcp_probe_used() - 1;
95  }
96
97  #define tcp_probe_copy_fl_to_si4(inet, si4, mem)              \
98          do {                                                  \
99                  si4.sin_family = AF_INET;                     \
100                 si4.sin_port = inet->inet_##mem##port;        \
101                 si4.sin_addr.s_addr = inet->inet_##mem##addr; \
102         } while (0)
103
104
105 /*
106  * Hook inserted to be called before each receive packet.
107  * Note: arguments must match tcp_rcv_established()!
108  */
```

```
109 static void jtcp_rcv_established(struct sock *sk, struct sk_buff *skb,
110                                  const struct tcphdr *th, unsigned int len)
111 {
112         const struct tcp_sock *tp = tcp_sk(sk);
113         const struct inet_sock *inet = inet_sk(sk);
114
115         /* Only update if port or skb mark matches */
116         if (((port == 0 && fwmark == 0) ||
117               ntohs(inet->inet_dport) == port ||
118               ntohs(inet->inet_sport) == port ||
119              (fwmark > 0 && skb->mark == fwmark)) &&
120             (full || tp->snd_cwnd != tcp_probe.lastcwnd)) {
121
122                 spin_lock(&tcp_probe.lock);
123                 /* If log fills, just silently drop */
124                 if (tcp_probe_avail() > 1) {
125                         struct tcp_log *p = tcp_probe.log + tcp_probe.head;
126
127                         p->tstamp = ktime_get();
128                         switch (sk->sk_family) {
129                         case AF_INET:
130                                 tcp_probe_copy_fl_to_si4(inet, p->src.v4, s);
131                                 tcp_probe_copy_fl_to_si4(inet, p->dst.v4, d);
132                                 break;
133                         case AF_INET6:
134                                 memset(&p->src.v6, 0, sizeof(p->src.v6));
135                                 memset(&p->dst.v6, 0, sizeof(p->dst.v6));
136 #if IS_ENABLED(CONFIG_IPV6)
137                                 p->src.v6.sin6_family = AF_INET6;
138                                 p->src.v6.sin6_port = inet->inet_sport;
139                                 p->src.v6.sin6_addr = inet6_sk(sk)->saddr;
140
141                                 p->dst.v6.sin6_family = AF_INET6;
142                                 p->dst.v6.sin6_port = inet->inet_dport;
143                                 p->dst.v6.sin6_addr = sk->sk_v6_daddr;
144 #endif
145                                 break;
146                         default:
147                                 BUG();
148                         }
149
150                         p->length = skb->len;
151                         p->snd_nxt = tp->snd_nxt;
152                         p->snd_una = tp->snd_una;
153                         p->snd_cwnd = tp->snd_cwnd;
154                         p->snd_wnd = tp->snd_wnd;
155                         p->rcv_wnd = tp->rcv_wnd;
156                         p->ssthresh = tcp_current_ssthresh(sk);
157                         p->srtt = tp->srtt_us >> 3;
158
159                         tcp_probe.head = (tcp_probe.head + 1) & (bufsize - 1);
160                 }
161                 tcp_probe.lastcwnd = tp->snd_cwnd;
162                 spin_unlock(&tcp_probe.lock);
163
164                 wake_up(&tcp_probe.wait);
165         }
166
167         jprobe_return();
168 }
169
170 static struct jprobe tcp_jprobe = {
171         .kp = {
172                 .symbol_name    = "tcp_rcv_established",
173         },
```

```
174              .entry  = jtcp_rcv_established,
175 };
176
177 static int tcpprobe_open(struct inode *inode, struct file *file)
178 {
179              /* Reset (empty) log */
180              spin_lock_bh(&tcp_probe.lock);
181              tcp_probe.head = tcp_probe.tail = 0;
182              tcp_probe.start = ktime_get();
183              spin_unlock_bh(&tcp_probe.lock);
184
185              return 0;
186 }
187
188 static int tcpprobe_sprint(char *tbuf, int n)
189 {
190              const struct tcp_log *p
191                      = tcp_probe.log + tcp_probe.tail;
192              struct timespec tv
193                      = ktime_to_timespec(ktime_sub(p->tstamp, tcp_probe.start));
194
195              return scnprintf(tbuf, n,
196                      "%lu.%09lu %pISpc %pISpc %d %#x %#x %u %u %u %u %u\n",
197                      (unsigned long) tv.tv_sec,
198                      (unsigned long) tv.tv_nsec,
199                      &p->src, &p->dst, p->length, p->snd_nxt, p->snd_una,
200                      p->snd_cwnd, p->ssthresh, p->snd_wnd, p->srtt, p->rcv_wnd);
201 }
202
203 static ssize_t tcpprobe_read(struct file *file, char __user *buf,
204                              size_t len, loff_t *ppos)
205 {
206              int error = 0;
207              size_t cnt = 0;
208
209              if (!buf)
210                      return -EINVAL;
211
212              while (cnt < len) {
213                      char tbuf[256];
214                      int width;
215
216                      /* Wait for data in buffer */
217                      error = wait_event_interruptible(tcp_probe.wait,
218                                              tcp_probe_used() > 0);
219                      if (error)
220                              break;
221
222                      spin_lock_bh(&tcp_probe.lock);
223                      if (tcp_probe.head == tcp_probe.tail) {
224                              /* multiple readers race? */
225                              spin_unlock_bh(&tcp_probe.lock);
226                              continue;
227                      }
228
229                      width = tcpprobe_sprint(tbuf, sizeof(tbuf));
230
231                      if (cnt + width < len)
232                              tcp_probe.tail = (tcp_probe.tail + 1) & (bufsize - 1);
233
234                      spin_unlock_bh(&tcp_probe.lock);
235
236                      /* if record greater than space available
237                         return partial buffer (so far) */
238                      if (cnt + width >= len)
```

```
239                              break;
240
241                   if (copy_to_user(buf + cnt, tbuf, width))
242                           return -EFAULT;
243                   cnt += width;
244           }
245
246           return cnt == 0 ? error : cnt;
247 }
248
249 static const struct file_operations tcpprobe_fops = {
250           .owner   = THIS_MODULE,
251           .open    = tcpprobe_open,
252           .read    = tcpprobe_read,
253           .llseek  = noop_llseek,
254 };
255
256 static __init int tcpprobe_init(void)
257 {
258           int ret = -ENOMEM;
259
260           /* Warning: if the function signature of tcp_rcv_established,
261            * has been changed, you also have to change the signature of
262            * jtcp_rcv_established, otherwise you end up right here!
263            */
264           BUILD_BUG_ON(__same_type(tcp_rcv_established,
265                                    jtcp_rcv_established) == 0);
266
267           init_waitqueue_head(&tcp_probe.wait);
268           spin_lock_init(&tcp_probe.lock);
269
270           if (bufsize == 0)
271                   return -EINVAL;
272
273           bufsize = roundup_pow_of_two(bufsize);
274           tcp_probe.log = kcalloc(bufsize, sizeof(struct tcp_log), GFP_KERNEL);
275           if (!tcp_probe.log)
276                   goto err0;
277
278           if (!proc_create(procname, S_IRUSR, init_net.proc_net, &tcpprobe_fops))
279                   goto err0;
280
281           ret = register_jprobe(&tcp_jprobe);
282           if (ret)
283                   goto err1;
284
285           pr_info("probe registered (port=%d/fwmark=%u) bufsize=%u\n",
286                   port, fwmark, bufsize);
287           return 0;
288  err1:
289           remove_proc_entry(procname, init_net.proc_net);
290  err0:
291           kfree(tcp_probe.log);
292           return ret;
293 }
294 module_init(tcpprobe_init);
295
296 static __exit void tcpprobe_exit(void)
297 {
298           remove_proc_entry(procname, init_net.proc_net);
299           unregister_jprobe(&tcp_jprobe);
300           kfree(tcp_probe.log);
301 }
302 module_exit(tcpprobe_exit);
303
```

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](#)