# 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 3.6 3.7 3.8 3.9 3.10 3.11 3.12 3.13 3.14 3.15 3.16 3.17 3.18 3.19 4.0 4.1 *4.2*

# [Linux](#)/[net](#)/[ipv4](#)/[ip_input.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   *             The Internet Protocol (IP) module.
7   *
8   * Authors:    Ross Biro
9   *             Fred N. van Kempen, <waltje@uWalt.NL.Mugnet.ORG>
10  *             Donald Becker, <becker@super.org>
11  *             Alan Cox, <alan@lxorguk.ukuu.org.uk>
12  *             Richard Underwood
13  *             Stefan Becker, <stefanb@yello.ping.de>
14  *             Jorge Cwik, <jorge@laser.satlink.net>
15  *             Arnt Gulbrandsen, <agulbra@nvg.unit.no>
16  *
17  *
18  * Fixes:
19  *             Alan Cox        :       Commented a couple of minor bits of surplus code
20  *             Alan Cox        :       Undefining IP_FORWARD doesn't include the code
21  *                                     (just stops a compiler warning).
22  *             Alan Cox        :       Frames with >=MAX_ROUTE record routes, strict routes or loose routes
23  *                                     are junked rather than corrupting things.
24  *             Alan Cox        :       Frames to bad broadcast subnets are dumped
25  *                                     We used to process them non broadcast and
26  *                                     boy could that cause havoc.
27  *             Alan Cox        :       ip_forward sets the free flag on the
28  *                                     new frame it queues. Still crap because
29  *                                     it copies the frame but at least it
30  *                                     doesn't eat memory too.
31  *             Alan Cox        :       Generic queue code and memory fixes.
32  *             Fred Van Kempen :       IP fragment support (borrowed from NET2E)
33  *             Gerhard Koerting:       Forward fragmented frames correctly.
34  *             Gerhard Koerting:       Fixes to my fix of the above 8-).
35  *             Gerhard Koerting:       IP interface addressing fix.
36  *             Linus Torvalds  :       More robustness checks
37  *             Alan Cox        :       Even more checks: Still not as robust as it ought to be
38  *             Alan Cox        :       Save IP header pointer for later
39  *             Alan Cox        :       ip option setting
40  *             Alan Cox        :       Use ip_tos/ip_ttl settings
41  *             Alan Cox        :       Fragmentation bogosity removed
42  *                                     (Thanks to Mark.Bush@prg.ox.ac.uk)
43  *             Dmitry Gorodchanin :    Send of a raw packet crash fix.
44  *             Alan Cox        :       Silly ip bug when an overlength
45  *                                     fragment turns up. Now frees the
46  *                                     queue.
47  *             Linus Torvalds/ :       Memory leakage on fragmentation
48  *             Alan Cox        :       handling.
49  *             Gerhard Koerting:       Forwarding uses IP priority hints
50  *             Teemu Rantanen  :       Fragment problems.
51  *             Alan Cox        :       General cleanup, comments and reformat
52  *             Alan Cox        :       SNMP statistics
53  *             Alan Cox        :       BSD address rule semantics. Also see
54  *                                     UDP as there is a nasty checksum issue
55  *                                     if you do things the wrong way.
56  *             Alan Cox        :       Always defrag, moved IP_FORWARD to the config.in file
57  *             Alan Cox        :       IP options adjust sk->priority.
58  *             Pedro Roque     :       Fix mtu/length error in ip_forward.
```

```
 59  *              Alan Cox        :       Avoid ip_chk_addr when possible.
 60  *      Richard Underwood       :       IP multicasting.
 61  *              Alan Cox        :       Cleaned up multicast handlers.
 62  *              Alan Cox        :       RAW sockets demultiplex in the BSD style.
 63  *              Gunther Mayer   :       Fix the SNMP reporting typo
 64  *              Alan Cox        :       Always in group 224.0.0.1
 65  *      Pauline Middelink       :       Fast ip_checksum update when forwarding
 66  *                                      Masquerading support.
 67  *              Alan Cox        :       Multicast loopback error for 224.0.0.1
 68  *              Alan Cox        :       IP_MULTICAST_LOOP option.
 69  *              Alan Cox        :       Use notifiers.
 70  *              Bjorn Ekwall    :       Removed ip_csum (from slhc.c too)
 71  *              Bjorn Ekwall    :       Moved ip_fast_csum to ip.h (inline!)
 72  *              Stefan Becker   :       Send out ICMP HOST REDIRECT
 73  *      Arnt Gulbrandsen        :       ip_build_xmit
 74  *              Alan Cox        :       Per socket routing cache
 75  *              Alan Cox        :       Fixed routing cache, added header cache.
 76  *              Alan Cox        :       Loopback didn't work right in original ip_build_xmit - fixed it.
 77  *              Alan Cox        :       Only send ICMP_REDIRECT if src/dest are the same net.
 78  *              Alan Cox        :       Incoming IP option handling.
 79  *              Alan Cox        :       Set saddr on raw output frames as per BSD.
 80  *              Alan Cox        :       Stopped broadcast source route explosions.
 81  *              Alan Cox        :       Can disable source routing
 82  *              Takeshi Sone    :       Masquerading didn't work.
 83  *      Dave Bonn,Alan Cox      :       Faster IP forwarding whenever possible.
 84  *              Alan Cox        :       Memory leaks, tramples, misc debugging.
 85  *              Alan Cox        :       Fixed multicast (by popular demand 8))
 86  *              Alan Cox        :       Fixed forwarding (by even more popular demand 8))
 87  *              Alan Cox        :       Fixed SNMP statistics [I think]
 88  *      Gerhard Koerting        :       IP fragmentation forwarding fix
 89  *              Alan Cox        :       Device lock against page fault.
 90  *              Alan Cox        :       IP_HDRINCL facility.
 91  *      Werner Almesberger      :       Zero fragment bug
 92  *              Alan Cox        :       RAW IP frame length bug
 93  *              Alan Cox        :       Outgoing firewall on build_xmit
 94  *              A.N.Kuznetsov   :       IP_OPTIONS support throughout the kernel
 95  *              Alan Cox        :       Multicast routing hooks
 96  *              Jos Vos         :       Do accounting *before* call_in_firewall
 97  *      Willy Konynenberg       :       Transparent proxying support
 98  *
 99  *
100  *
101  * To Fix:
102  *              IP fragmentation wants rewriting cleanly. The RFC815 algorithm is much more efficient
103  *              and could be made very efficient with the addition of some virtual memory hacks to permit
104  *              the allocation of a buffer that can then be 'grown' by twiddling page tables.
105  *              Output fragmentation wants updating along with the buffer management to use a single
106  *              interleaved copy algorithm so that fragmenting has a one copy overhead. Actual packet
107  *              output should probably do its own fragmentation at the UDP/RAW layer. TCP shouldn't cause
108  *              fragmentation anyway.
109  *
110  *              This program is free software; you can redistribute it and/or
111  *              modify it under the terms of the GNU General Public License
112  *              as published by the Free Software Foundation; either version
113  *              2 of the License, or (at your option) any later version.
114  */
115
116  #define pr_fmt(fmt) "IPv4: " fmt
117
118  #include <linux/module.h>
119  #include <linux/types.h>
120  #include <linux/kernel.h>
121  #include <linux/string.h>
122  #include <linux/errno.h>
123  #include <linux/slab.h>
124
125  #include <linux/net.h>
126  #include <linux/socket.h>
127  #include <linux/sockios.h>
128  #include <linux/in.h>
129  #include <linux/inet.h>
130  #include <linux/inetdevice.h>
131  #include <linux/netdevice.h>
132  #include <linux/etherdevice.h>
133
134  #include <net/snmp.h>
135  #include <net/ip.h>
136  #include <net/protocol.h>
```

```
137 #include <net/route.h>
138 #include <linux/skbuff.h>
139 #include <net/sock.h>
140 #include <net/arp.h>
141 #include <net/icmp.h>
142 #include <net/raw.h>
143 #include <net/checksum.h>
144 #include <net/inet_ecn.h>
145 #include <linux/netfilter_ipv4.h>
146 #include <net/xfrm.h>
147 #include <linux/mroute.h>
148 #include <linux/netlink.h>
149
150 /*
151  *      Process Router Attention IP option (RFC 2113)
152  */
153 bool ip_call_ra_chain(struct sk_buff *skb)
154 {
155         struct ip_ra_chain *ra;
156         u8 protocol = ip_hdr(skb)->protocol;
157         struct sock *last = NULL;
158         struct net_device *dev = skb->dev;
159
160         for (ra = rcu_dereference(ip_ra_chain); ra; ra = rcu_dereference(ra->next)) {
161                 struct sock *sk = ra->sk;
162
163                 /* If socket is bound to an interface, only report
164                  * the packet if it came  from that interface.
165                  */
166                 if (sk && inet_sk(sk)->inet_num == protocol &&
167                     (!sk->sk_bound_dev_if ||
168                      sk->sk_bound_dev_if == dev->ifindex) &&
169                     net_eq(sock_net(sk), dev_net(dev))) {
170                         if (ip_is_fragment(ip_hdr(skb))) {
171                                 if (ip_defrag(skb, IP_DEFRAG_CALL_RA_CHAIN))
172                                         return true;
173                         }
174                         if (last) {
175                                 struct sk_buff *skb2 = skb_clone(skb, GFP_ATOMIC);
176                                 if (skb2)
177                                         raw_rcv(last, skb2);
178                         }
179                         last = sk;
180                 }
181         }
182
183         if (last) {
184                 raw_rcv(last, skb);
185                 return true;
186         }
187         return false;
188 }
189
190 static int ip_local_deliver_finish(struct sock *sk, struct sk_buff *skb)
191 {
192         struct net *net = dev_net(skb->dev);
193
194         __skb_pull(skb, skb_network_header_len(skb));
195
196         rcu_read_lock();
197         {
198                 int protocol = ip_hdr(skb)->protocol;
199                 const struct net_protocol *ipprot;
200                 int raw;
201
202 resubmit:
203                 raw = raw_local_deliver(skb, protocol);
204
205                 ipprot = rcu_dereference(inet_protos[protocol]);
206                 if (ipprot) {
207                         int ret;
208
209                         if (!ipprot->no_policy) {
210                                 if (!xfrm4_policy_check(NULL, XFRM_POLICY_IN, skb)) {
211                                         kfree_skb(skb);
212                                         goto out;
213                                 }
214                                 nf_reset(skb);
```

```
215                                }
216                                ret = ipprot->handler(skb);
217                                if (ret < 0) {
218                                        protocol = -ret;
219                                        goto resubmit;
220                                }
221                                IP_INC_STATS_BH(net, IPSTATS_MIB_INDELIVERS);
222                        } else {
223                                if (!raw) {
224                                        if (xfrm4_policy_check(NULL, XFRM_POLICY_IN, skb)) {
225                                                IP_INC_STATS_BH(net, IPSTATS_MIB_INUNKNOWNPROTOS);
226                                                icmp_send(skb, ICMP_DEST_UNREACH,
227                                                        ICMP_PROT_UNREACH, 0);
228                                        }
229                                        kfree_skb(skb);
230                                } else {
231                                        IP_INC_STATS_BH(net, IPSTATS_MIB_INDELIVERS);
232                                        consume_skb(skb);
233                                }
234                        }
235                }
236 out:
237            rcu_read_unlock();
238
239            return 0;
240 }
241
242 /*
243  *      Deliver IP Packets to the higher protocol layers.
244  */
245 int ip_local_deliver(struct sk_buff *skb)
246 {
247            /*
248             *      Reassemble IP fragments.
249             */
250
251            if (ip_is_fragment(ip_hdr(skb))) {
252                    if (ip_defrag(skb, IP_DEFRAG_LOCAL_DELIVER))
253                            return 0;
254            }
255
256            return NF_HOOK(NFPROTO_IPV4, NF_INET_LOCAL_IN, NULL, skb,
257                    skb->dev, NULL,
258                    ip_local_deliver_finish);
259 }
260
261 static inline bool ip_rcv_options(struct sk_buff *skb)
262 {
263            struct ip_options *opt;
264            const struct iphdr *iph;
265            struct net_device *dev = skb->dev;
266
267            /* It looks as overkill, because not all
268               IP options require packet mangling.
269               But it is the easiest for now, especially taking
270               into account that combination of IP options
271               and running sniffer is extremely rare condition.
272                                        --ANK (980813)
273            */
274            if (skb_cow(skb, skb_headroom(skb))) {
275                    IP_INC_STATS_BH(dev_net(dev), IPSTATS_MIB_INDISCARDS);
276                    goto drop;
277            }
278
279            iph = ip_hdr(skb);
280            opt = &(IPCB(skb)->opt);
281            opt->optlen = iph->ihl*4 - sizeof(struct iphdr);
282
283            if (ip_options_compile(dev_net(dev), opt, skb)) {
284                    IP_INC_STATS_BH(dev_net(dev), IPSTATS_MIB_INHDRERRORS);
285                    goto drop;
286            }
287
288            if (unlikely(opt->srr)) {
289                    struct in_device *in_dev = __in_dev_get_rcu(dev);
290
291                    if (in_dev) {
292                            if (!IN_DEV_SOURCE_ROUTE(in_dev)) {
```

```
293                                 if (IN_DEV_LOG_MARTIANS(in_dev))
294                                         net_info_ratelimited("source route option %pI4 -> %pI4\n",
295                                                              &iph->saddr,
296                                                              &iph->daddr);
297                                 goto drop;
298                         }
299                 }
300
301                 if (ip_options_rcv_srr(skb))
302                         goto drop;
303         }
304
305         return false;
306 drop:
307         return true;
308 }
309
310 int sysctl_ip_early_demux __read_mostly = 1;
311 EXPORT_SYMBOL(sysctl_ip_early_demux);
312
313 static int ip_rcv_finish(struct sock *sk, struct sk_buff *skb)
314 {
315         const struct iphdr *iph = ip_hdr(skb);
316         struct rtable *rt;
317
318         if (sysctl_ip_early_demux && !skb_dst(skb) && !skb->sk) {
319                 const struct net_protocol *ipprot;
320                 int protocol = iph->protocol;
321
322                 ipprot = rcu_dereference(inet_protos[protocol]);
323                 if (ipprot && ipprot->early_demux) {
324                         ipprot->early_demux(skb);
325                         /* must reload iph, skb->head might have changed */
326                         iph = ip_hdr(skb);
327                 }
328         }
329
330         /*
331          *      Initialise the virtual path cache for the packet. It describes
332          *      how the packet travels inside Linux networking.
333          */
334         if (!skb_dst(skb)) {
335                 int err = ip_route_input_noref(skb, iph->daddr, iph->saddr,
336                                                iph->tos, skb->dev);
337                 if (unlikely(err)) {
338                         if (err == -EXDEV)
339                                 NET_INC_STATS_BH(dev_net(skb->dev),
340                                                  LINUX_MIB_IPRPFILTER);
341                         goto drop;
342                 }
343         }
344
345 #ifdef CONFIG_IP_ROUTE_CLASSID
346         if (unlikely(skb_dst(skb)->tclassid)) {
347                 struct ip_rt_acct *st = this_cpu_ptr(ip_rt_acct);
348                 u32 idx = skb_dst(skb)->tclassid;
349                 st[idx&0xFF].o_packets++;
350                 st[idx&0xFF].o_bytes += skb->len;
351                 st[(idx>>16)&0xFF].i_packets++;
352                 st[(idx>>16)&0xFF].i_bytes += skb->len;
353         }
354 #endif
355
356         if (iph->ihl > 5 && ip_rcv_options(skb))
357                 goto drop;
358
359         rt = skb_rtable(skb);
360         if (rt->rt_type == RTN_MULTICAST) {
361                 IP_UPD_PO_STATS_BH(dev_net(rt->dst.dev), IPSTATS_MIB_INMCAST,
362                                 skb->len);
363         } else if (rt->rt_type == RTN_BROADCAST)
364                 IP_UPD_PO_STATS_BH(dev_net(rt->dst.dev), IPSTATS_MIB_INBCAST,
365                                 skb->len);
366
367         return dst_input(skb);
368
369 drop:
370         kfree_skb(skb);
```

```
371                return NET_RX_DROP;
372 }
373
374 /*
375  *        Main IP Receive routine.
376  */
377 int ip_rcv(struct sk_buff *skb, struct net_device *dev, struct packet_type *pt, struct net_device *orig_dev)
378 {
379        const struct iphdr *iph;
380        u32 len;
381
382        /* When the interface is in promisc. mode, drop all the crap
383         * that it receives, do not try to analyse it.
384         */
385        if (skb->pkt_type == PACKET_OTHERHOST)
386                goto drop;
387
388
389        IP_UPD_PO_STATS_BH(dev_net(dev), IPSTATS_MIB_IN, skb->len);
390
391        skb = skb_share_check(skb, GFP_ATOMIC);
392        if (!skb) {
393                IP_INC_STATS_BH(dev_net(dev), IPSTATS_MIB_INDISCARDS);
394                goto out;
395        }
396
397        if (!pskb_may_pull(skb, sizeof(struct iphdr)))
398                goto inhdr_error;
399
400        iph = ip_hdr(skb);
401
402        /*
403         *        RFC1122: 3.2.1.2 MUST silently discard any IP frame that fails the checksum.
404         *
405         *        Is the datagram acceptable?
406         *
407         *        1.        Length at least the size of an ip header
408         *        2.        Version of 4
409         *        3.        Checksums correctly. [Speed optimisation for later, skip loopback checksums]
410         *        4.        Doesn't have a bogus length
411         */
412
413        if (iph->ihl < 5 || iph->version != 4)
414                goto inhdr_error;
415
416        BUILD_BUG_ON(IPSTATS_MIB_ECT1PKTS != IPSTATS_MIB_NOECTPKTS + INET_ECN_ECT_1);
417        BUILD_BUG_ON(IPSTATS_MIB_ECT0PKTS != IPSTATS_MIB_NOECTPKTS + INET_ECN_ECT_0);
418        BUILD_BUG_ON(IPSTATS_MIB_CEPKTS != IPSTATS_MIB_NOECTPKTS + INET_ECN_CE);
419        IP_ADD_STATS_BH(dev_net(dev),
420                        IPSTATS_MIB_NOECTPKTS + (iph->tos & INET_ECN_MASK),
421                        max_t(unsigned short, 1, skb_shinfo(skb)->gso_segs));
422
423        if (!pskb_may_pull(skb, iph->ihl*4))
424                goto inhdr_error;
425
426        iph = ip_hdr(skb);
427
428        if (unlikely(ip_fast_csum((u8 *)iph, iph->ihl)))
429                goto csum_error;
430
431        len = ntohs(iph->tot_len);
432        if (skb->len < len) {
433                IP_INC_STATS_BH(dev_net(dev), IPSTATS_MIB_INTRUNCATEDPKTS);
434                goto drop;
435        } else if (len < (iph->ihl*4))
436                goto inhdr_error;
437
438        /* Our transport medium may have padded the buffer out. Now we know it
439         * is IP we can trim to the true length of the frame.
440         * Note this now means skb->len holds ntohs(iph->tot_len).
441         */
442        if (pskb_trim_rcsum(skb, len)) {
443                IP_INC_STATS_BH(dev_net(dev), IPSTATS_MIB_INDISCARDS);
444                goto drop;
445        }
446
447        skb->transport_header = skb->network_header + iph->ihl*4;
448
```

```
449            /* Remove any debris in the socket control block */
450            memset(IPCB(skb), 0, sizeof(struct inet_skb_parm));
451
452            /* Must drop socket now because of tproxy. */
453            skb_orphan(skb);
454
455            return NF_HOOK(NFPROTO_IPV4, NF_INET_PRE_ROUTING, NULL, skb,
456                           dev, NULL,
457                           ip_rcv_finish);
458
459  csum_error:
460            IP_INC_STATS_BH(dev_net(dev), IPSTATS_MIB_CSUMERRORS);
461  inhdr_error:
462            IP_INC_STATS_BH(dev_net(dev), IPSTATS_MIB_INHDRERRORS);
463  drop:
464            kfree_skb(skb);
465  out:
466            return NET_RX_DROP;
467  }
468
```

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