# Sheet 5

# Forwarding intro

1 **ip_rcv**

2

3 Right, so now we're going to look at routing. We start with net/ipv4/ip_input.c::ip_rcv, which is the main IP receive routine,

4 and which needs to make a decision as to whether to deal with this packet locally or forward it.

5

```
6  int ip_rcv(struct sk_buff *skb, struct net_device *dev, struct packet_type *pt)
7  {
8      struct iphdr *iph = skb->nh.iph;
9
```

10 There are several stages here. First, we're testing to see whether what we have is well formed.

```
11      /* When the interface is in promisc. mode, drop all the crap
12       * that it receives, do not try to analyse it.
13       */
14      if (skb->pkt_type == PACKET_OTHERHOST)
15          goto drop;
16
17      IP_INC_STATS_BH(IpInReceives);
18
19      if ((skb = skb_share_check(skb, GFP_ATOMIC)) == NULL)
20          goto out;
21
22      if (!pskb_may_pull(skb, sizeof(struct iphdr)))
23          goto inhdr_error;
24
```

```
25          iph = skb->nh.iph;
26
27
```

**28   In particular, we require that**

**29   1.      Length at least the size of an ip header**

**30   2.      Version of 4**

**31   3.      Checksums correctly. [Speed optimisation for later, skip loopback checksums]**

**32   4.      Doesn't have a bogus length**

**33   Also, RFC1122: 3.1.2.2 MUST silently discard any IP frame that fails the checksum.**

```
34          if (iph->ihl < 5 || iph->version != 4)
35              goto inhdr_error;
36
37          if (!pskb_may_pull(skb, iph->ihl*4))
38              goto inhdr_error;
39
40          if (ip_fast_csum((u8 *)iph, iph->ihl) != 0)
41              goto inhdr_error;
42
43          {
44              __u32 len = ntohs(iph->tot_len);
45              if (skb->len < len || len < (iph->ihl<<2))
46                  goto inhdr_error;
47
```

48  Our transport medium may have padded the buffer out. Now we know it is IP we can trim to the true length of the frame. Note
49  this now means skb->len holds ntohs(iph->tot_len).

```
50              if (skb->len > len) {
51                  __pskb_trim(skb, len);
52                  if (skb->ip_summed == CHECKSUM_HW)
53                      skb->ip_summed = CHECKSUM_NONE;
54              }
55          }
56
```

57  AOK, so the work carries on in [net/ipv4/ip_input.c::ip_rcv_finish](), described below after netfiltering

```
58      return NF_HOOK(PF_INET, NF_IP_PRE_ROUTING, skb, dev, NULL, ip_rcv_finish);
59
```

60  Otherwise, fail quietly.

```
61  inhdr_error:
62      IP_INC_STATS_BH(IpInHdrErrors);
63  drop:
64      kfree_skb(skb);
65  out:
66      return NET_RX_DROP;
67  }
68
69
```

69    **ip_rcv_finish**

70

71    At this point, we know that we have an IP packet, but we haven't decided what to do with it.

```
72    static inline int ip_rcv_finish(struct sk_buff *skb)
73    {
74          struct net_device *dev = skb->dev;
75          struct iphdr *iph = skb->nh.iph;
76
```

77    If we don't already know what to do with the packet, call net/ipv4/route.c::ip_route_input. This determines whether we
78    should forward this on another interface or whether it is for us. We'll talk about this later.

```
79          if (skb->dst == NULL) {
80                if (ip_route_input(skb, iph->daddr, iph->saddr, iph->tos, dev))
81                      goto drop;
82          }
83
84          <deal with IP options>
85
```

86    This is either call net/ipv4/ip_forward.c::ip_forward, or local delivery. We'll deal with forwarding below.

```
87          return skb->dst->input(skb);
88
89    inhdr_error:
90          IP_INC_STATS_BH(IpInHdrErrors);
91    drop:
```

```
92        kfree_skb(skb);
93        return NET_RX_DROP;
94   }
95
```

95  **ip_forward**

96

97  **We've decided that we need to forward the packet.**

```
98   int ip_forward(struct sk_buff *skb)
99   {
100      struct net_device      *dev2;   /* Output device */
101      struct iphdr           *iph;    /* Our header */
102      struct rtable          *rt;     /* Route we use */
103      struct ip_options      *opt = &(IPCB(skb)->opt);
104      unsigned short         mtu;
105
```

106 **Housekeeping**

```
107      if (IPCB(skb)->opt.router_alert && ip_call_ra_chain(skb))
108          return NET_RX_SUCCESS;
109
110      if (skb->pkt_type != PACKET_HOST)
111          goto drop;
112
113      skb->ip_summed = CHECKSUM_NONE;
114
115      iph = skb->nh.iph;
116      rt = (struct rtable*)skb->dst;
117
```

118  **According to the RFC, we must first decrease the TTL field. If that reaches zero, we must reply an ICMP control message**
119  **telling that the packet's lifetime expired.**

```
120         if (iph->ttl <= 1)
121             goto too_many_hops;
122
```

123  **Fail if we're supposed to be routing strictly and we can't meet the requirements. Strict routing permits no gatewaying**

```
124         if (opt->is_strictroute && rt->rt_dst != rt->rt_gateway)
125             goto sr_failed;
126
```

127  **Having picked a route we can now send the frame out after asking the firewall permission to do so.**

```
128         skb->priority = rt_tos2priority(iph->tos);
129         dev2 = rt->u.dst.dev;
130         mtu  = rt->u.dst.pmtu;
131
132         /*
133          *   We now generate an ICMP HOST REDIRECT giving the route
134          *   we calculated.
135          */
136         if (rt->rt_flags&RTCF_DOREDIRECT && !opt->srr)
137             ip_rt_send_redirect(skb);
138
```

139  **We are about to mangle packet. Copy it! (cow = copy on write)**

```
140         if (skb_cow(skb, dev2->hard_header_len))
141             goto drop;
142         iph = skb->nh.iph;
```

143

**Decrease ttl after skb cow done**

```
        ip_decrease_ttl(iph);
```

146

**Check if we need to fragment, because MTU is lower than the length of this packet. If we do, and the 'don't fragment' flag is set, then generate error.**

```
        if (skb->len > mtu && (ntohs(iph->frag_off) & IP_DF))
                goto frag_needed;
```

151

        **<NAT stuff>**

153

**Carry on with net/ipv4/ip_forward.c::ip_forward_finish after a netfilter.**

```
        return NF_HOOK(PF_INET, NF_IP_FORWARD, skb, skb->dev, dev2,
                                 ip_forward_finish);
```

157

**Generate diagnostic ICMP messages for unroutable packets**

```
frag_needed:
        IP_INC_STATS_BH(IpFragFails);
        icmp_send(skb, ICMP_DEST_UNREACH, ICMP_FRAG_NEEDED, htonl(mtu));
        goto drop;

sr_failed:
        icmp_send(skb, ICMP_DEST_UNREACH, ICMP_SR_FAILED, 0);
        goto drop;
```

```
167
168   too_many_hops:
169       icmp_send(skb, ICMP_TIME_EXCEEDED, ICMP_EXC_TTL, 0);
170   drop:
171       kfree_skb(skb);
172       return NET_RX_DROP;
173   }
174
```

174  **ip_forward_finish**

175

176  Get the packet on its way.

177

```
178  static inline int ip_forward_finish(struct sk_buff *skb)
179  {
180      struct ip_options * opt    = &(IPCB(skb)->opt);
181
182      IP_INC_STATS_BH(IpForwDatagrams);
183
```

184  If we haven't got any options, then come in here.

```
185      if (opt->optlen == 0) {
186          <deal with fast routing>
187
```

188  Send it

```
189          return (ip_send(skb));
190      }
191
```

192  Deal with options.

```
193      ip_forward_options(skb);
```

194  And send it.

```
195      return (ip_send(skb));
```

```
196    }
197
```

197    **ip_route_input**

198

199    OK, so how did we decide whether to forward or deliver locally? Look back and you'll see that we called **ip_route_input(skb,**

200    **iph->daddr, iph->saddr, iph->tos, dev) in net/ipv4/ip_input::ip_rcv_finish.**

201

202    One important thing to note here is that we have two structures we use for routing. The first is a transient route cache, held

203    in a hash table as below; the second is the permanent FIB (forwarding information base). The FIB is only consulted if we

204    don't have a matching hash table entry.

205

206    ```
int ip_route_input(struct sk_buff *skb, u32 daddr, u32 saddr, u8 tos, struct
```

207    ```
net_device *dev)
```

208    ```
{
```

209    ```
    struct rtable    *rth;
```

210    ```
    unsigned         hash;
```

211    ```
    int              iif = dev->ifindex;
```

212

213    Generate an index into the routing table, based on the destination address, sending address, input interface and tos fields.

214    ```
    tos &= IPTOS_RT_MASK;
```

215    ```
    hash = rt_hash_code(daddr, saddr ^ (iif << 5), tos);
```

216

217    Now chain down that entry looking for something that matches

218    ```
    read_lock(&rt_hash_table[hash].lock);
```

219    ```
    for (rth = rt_hash_table[hash].chain; rth; rth = rth->u.rt_next) {
```

```
220              if (rth->key.dst == daddr &&
221                  rth->key.src == saddr &&
222                  rth->key.iif == iif   &&
223                  rth->key.oif == 0     &&
224 #ifdef CONFIG_IP_ROUTE_FWMARK
225                  rth->key.fwmark == skb->nfmark &&
226 #endif
227                  rth->key.tos == tos) {
```

228   **We found something. Update time and usage stats**

```
229                  rth->u.dst.lastuse = jiffies;
230                  dst_hold(&rth->u.dst);
231                  rth->u.dst.__use++;
232                  rt_cache_stat[smp_processor_id()].in_hit++;
233                  read_unlock(&rt_hash_table[hash].lock);
```

234   **Set the destination for this packet**

```
235                  skb->dst = (struct dst_entry*)rth;
```

236   **And return success**

```
237                  return 0;
238              }
239          }
240      read_unlock(&rt_hash_table[hash].lock);
241
```

242      **<Deal with multicast traffic>**

243

244    **We need to look a bit harder.**

```
245        return ip_route_input_slow(skb, daddr, saddr, tos, dev);
246    }
247
```

247 **ip_route_input_slow**

248

249 The route we wanted was not in the route cache, so we need to consult the FIB. The FIB data is put in place by a routing
250 algorithm, and we'll look at that later. What we care about is looking up the FIB at present. This is complicated by the need
251 to do 'longest match' for CIDR addresses. See slides.

252 Take a look in the files include/net/ip_fib.h and net/ipv4/[fib_frontend.c, fib_hash.c, fib_rules.c, fib_semantics.c]

```
253 /*
254  *  NOTE. We drop all the packets that have local source addresses, because every
255  *  properly looped back packet must have correct destination already attached by
256  *  output routine.
257  *
258  *  Such an approach solves two big problems:
259  *  1. Not simplex devices are handled properly.
260  *  2. IP spoofing attempts are filtered with 100% of guarantee.
261  */
262
263 int ip_route_input_slow(struct sk_buff *skb, u32 daddr, u32 saddr, u8 tos, struct
264 net_device *dev)
265 {
266     struct rt_key          key;
267     struct fib_result      res;
268     struct in_device    *in_dev   = in_dev_get(dev);
269     struct in_device    *out_dev  = NULL;
270     unsigned              flags   = 0;
```

```
271        u32                      itag       = 0;
272        struct rtable          *rth;
273        unsigned                 hash;
274        u32                      spec_dst;
275        int                      err        = -EINVAL;
276        int                      free_res = 0;
277
```

278    Check to see if IP on this device is disabled.

```
279        if (!in_dev)
280            goto out;
281
```

282    Set up lookup key.

```
283        key.dst       = daddr;
284        key.src       = saddr;
285        key.tos       = tos;
286 #ifdef CONFIG_IP_ROUTE_FWMARK
287        key.fwmark    = skb->nfmark;
288 #endif
289        key.iif       = dev->ifindex;
290        key.oif       = 0;
291        key.scope     = RT_SCOPE_UNIVERSE;
292
293        hash = rt_hash_code(daddr, saddr ^ (key.iif << 5), tos);
294
```

295    Check for the most weird 'martians', which can be not detected by fib_lookup.

```
296        if (MULTICAST(saddr) || BADCLASS(saddr) || LOOPBACK(saddr))
297            goto martian_source;
298
299        if (daddr == 0xFFFFFFFF || (saddr == 0 && daddr == 0))
300            goto brd_input;
301
302        if (ZERONET(saddr))
303            goto martian_source;
304
305        if (BADCLASS(daddr) || ZERONET(daddr) || LOOPBACK(daddr))
306            goto martian_destination;
307
308
```

309    **This is the key statement -- lookup route in the FIB. We'll talk more about this later.**

```
310        if ((err = fib_lookup(&key, &res)) != 0) {
311            if (!IN_DEV_FORWARD(in_dev))
312                goto e_inval;
313            goto no_route;
314        }
```

315    **Mark res as needing to be deallocated. This gets done at label 'done:'**

```
316        free_res = 1;
317
318        rt_cache_stat[smp_processor_id()].in_slow_tot++;
319
320        <NAT stuff removed>
```

321

322 **Check for broadcast**

```
323     if (res.type == RTN_BROADCAST)
324         goto brd_input;
```

325

326 **Accept locally. But before we do, validate the source – check that it's not broadcast or our local address and that it arrived**
327 **on the right physical interface.**

```
328     if (res.type == RTN_LOCAL) {
329         int result;
330         result = fib_validate_source(saddr, daddr, tos, loopback_dev.ifindex, dev,
331                             &spec_dst, &itag);
332         if (result < 0)
333             goto martian_source;
334         if (result)
335             flags |= RTCF_DIRECTSRC;
336         spec_dst = daddr;
337         goto local_input;
338     }
```

339

340 **Check for more errors Only go past this point if we have a unicast direct route.**

```
341     if (!IN_DEV_FORWARD(in_dev))
342         goto e_inval;
343     if (res.type != RTN_UNICAST)
344         goto martian_destination;
```

345

346        **<multipath stuff removed>**

347

348  include/linux/inetdevice.h::in_dev_get increments a ref count and returns a include/linux/inetdevice.h::in_device

```
349        out_dev = in_dev_get(FIB_RES_DEV(res));
350        if (out_dev == NULL) {
351            if (net_ratelimit())
352                printk(KERN_CRIT "Bug in ip_route_input_slow(). Please, report\n");
353            goto e_inval;
354        }
```

355

356  **Validate the source – check that it's not broadcast or our local address and that it arrived on the right physical interface.**
357  **Also, calculate the logical interface this packet arrived on and calculate the 'specific destination' address.**

```
358        err = fib_validate_source(saddr, daddr, tos, FIB_RES_OIF(res), dev,
359                        &spec_dst, &itag);
360        if (err < 0)
361            goto martian_source;
362
363        if (err)
364            flags |= RTCF_DIRECTSRC;
365
366        if (out_dev == in_dev && err && !(flags & (RTCF_NAT | RTCF_MASQ)) &&
367                (IN_DEV_SHARED_MEDIA(out_dev) ||
368                 inet_addr_onlink(out_dev, saddr, FIB_RES_GW(res))))
369            flags |= RTCF_DOREDIRECT;
370
```

371 **Not IP (i.e. ARP). Do not create route, if it is invalid for proxy arp. Dynamic NAT (DNAT) routes are always valid.**

```
372         if (skb->protocol != __constant_htons(ETH_P_IP)) {
373             if (out_dev == in_dev && !(flags & RTCF_DNAT))
374                 goto e_inval;
375         }
376
```

377 **Set up cache entry.**

```
378         rth = dst_alloc(&ipv4_dst_ops);
379         if (!rth)
380             goto e_nobufs;
381
382         atomic_set(&rth->u.dst.__refcnt, 1);
383         rth->u.dst.flags  = DST_HOST;
384         rth->key.dst      = daddr;
385         rth->rt_dst       = daddr;
386         rth->key.tos      = tos;
387 #ifdef CONFIG_IP_ROUTE_FWMARK
388         rth->key.fwmark   = skb->nfmark;
389 #endif
390         rth->key.src      = saddr;
391         rth->rt_src       = saddr;
392         rth->rt_gateway   = daddr;
393
```

394         **<NAT stuff removed>**

395

```
396         rth->rt_iif       =
397         rth->key.iif      = dev->ifindex;
398         rth->u.dst.dev    = out_dev->dev;
399         dev_hold(rth->u.dst.dev);
400         rth->key.oif      = 0;
401         rth->rt_spec_dst  = spec_dst;
402
403         rth->u.dst.input  = ip_forward;
404         rth->u.dst.output = ip_output;
405
406         rt_set_nexthop(rth, &res, itag);
407
408         rth->rt_flags     = flags;
409
410     <FASTROUTE stuff removed>
411
412 intern:
413 Insert entry into the route cache
414         err = rt_intern_hash(hash, rth, (struct rtable**)&skb->dst);
415
416 done:
417 Reduce ref counts
418         in_dev_put(in_dev);
419         if (out_dev)
420             in_dev_put(out_dev);
```

```
421          if (free_res)
422              fib_res_put(&res);
423   out: return err;
424
425   brd_input:
426          <Deal with broadcast>
427
428   local_input:
429   Come here if we are willing to accept locally.
430          rth = dst_alloc(&ipv4_dst_ops);
431          if (!rth)
432              goto e_nobufs;
433
434          rth->u.dst.output= ip_rt_bug;
435
436          atomic_set(&rth->u.dst.__refcnt, 1);
437          rth->u.dst.flags  = DST_HOST;
438          rth->key.dst      = daddr;
439          rth->rt_dst       = daddr;
440          rth->key.tos      = tos;
441   #ifdef CONFIG_IP_ROUTE_FWMARK
442          rth->key.fwmark   = skb->nfmark;
443   #endif
444          rth->key.src      = saddr;
445          rth->rt_src       = saddr;
```

```
446  #ifdef CONFIG_IP_ROUTE_NAT
447      rth->rt_dst_map   = key.dst;
448      rth->rt_src_map   = key.src;
449  #endif
450  #ifdef CONFIG_NET_CLS_ROUTE
451      rth->u.dst.tclassid = itag;
452  #endif
453      rth->rt_iif       =
454      rth->key.iif      = dev->ifindex;
455      rth->u.dst.dev    = &loopback_dev;
456      dev_hold(rth->u.dst.dev);
457      rth->key.oif      = 0;
458      rth->rt_gateway   = daddr;
459      rth->rt_spec_dst  = spec_dst;
460      rth->u.dst.input  = ip_local_deliver;
461      rth->rt_flags     = flags|RTCF_LOCAL;
462      if (res.type == RTN_UNREACHABLE) {
463          rth->u.dst.input  = ip_error;
464          rth->u.dst.error  = -err;
465          rth->rt_flags     &= ~RTCF_LOCAL;
466      }
467      rth->rt_type = res.type;
468      goto intern;
469
470  no_route:
471      rt_cache_stat[smp_processor_id()].in_no_route++;
```

```
472          spec_dst = inet_select_addr(dev, 0, RT_SCOPE_UNIVERSE);
473          res.type = RTN_UNREACHABLE;
474          goto local_input;
475
476     martian_destination:
477          <Generate error>
478
479     e_inval:
480          err = -EINVAL;
481          goto done;
482
483     e_nobufs:
484          err = -ENOBUFS;
485          goto done;
486
487     martian_source:
488          <Generate error>
489          goto e_inval;
490     }
491
```