

# Sheet 6

## **FIB structures and lookup code**

1 Longest match – and how we do it in Linux. We achieve what we want using the Forwarding Information Base (FIB), which is  
 2 a complex structure in the kernel, containing the routing information we need indexed on its network mask. BTW all routes  
 3 with the same network mask are said to be in the same 'zone'.

4

5 First a look at some structures, then a look at the code that manipulates them.

## 6 **struct fib\_table**

7 This is [include/net/ip\\_fib.h::fib\\_table](#). It's the starting point for FIB traversal and is instantiated with data appropriate to the  
 8 type of network we are using at any given point in time.

```
9 struct fib_table
10 {
11     unsigned char    tb_id;
12     Timestamp
13     unsigned         tb_stamp;
14
```

15 Routines for lookup. These are set in [net/ipv4/fib\\_hash.c::fib\\_hash\\_init](#) to be [net/ipv4/fib\\_hash.c::\[fn\\_hash\\_lookup,](#)  
 16 [fn\\_hash\\_insert, fn\\_hash\\_delete,](#) etc.]

```
17     int              (*tb_lookup)(struct fib_table *tb, const struct rt_key *key,
18                                struct fib_result *res);
19     int              (*tb_insert)(struct fib_table *table, struct rtmsg *r,
20                                struct kern_rta *rta, struct nlmsg_hdr *n,
21                                struct netlink_skb_parms *req);
22     int              (*tb_delete)(struct fib_table *table, struct rtmsg *r,
23                                struct kern_rta *rta, struct nlmsg_hdr *n,
```

```
24         struct netlink_skb_parms *req);
25     int     (*tb_dump)(struct fib_table *table, struct sk_buff *skb,
26                 struct netlink_callback *cb);
27     int     (*tb_flush)(struct fib_table *table);
28     int     (*tb_get_info)(struct fib_table *table, char *buf, int first,
29                         int count);
30     void    (*tb_select_default)(struct fib_table *table,
31                                 const struct rt_key *key, struct fib_result *res);
32
```

33 **And this is set to be a [net/ipv4/fib\\_hash.c::fn\\_hash](#), described below.**

```
34     unsigned char  tb_data[0];
35 };
36
```

**36 struct fib\_hash**

37 This is [net/ipv4/fib\\_hash.c::fn\\_hash](#) and is the data we're talking about above. As you can see, we split the table into zones  
38 at a very high level. There are 33 possible netmasks (0x0000 to 0xFFFF) and a zone is defined by one of these. Also, the  
39 zones are all linked together, and the second field here points to the head of the list of zones.

```
40 struct fib_hash
41 {
42     struct fn_zone    *fn_zones[33];
43     struct fn_zone    *fn_zone_list;
44 };
45
46
```

46 **struct fib\_zone**47 This is [net/ipv4/fib\\_hash.c::fn\\_zone](#). It defines some housekeeping things about the hash table associated with each zone.

48 struct fn\_zone

49 {

50 Pointer to next non-empty zone in the hash structure where the netmask is *less* restrictive (= shorter) than this

51 struct fn\_zone \*fz\_next; /\* Next not empty zone \*/

52 This is a pointer to the hash table.

53 struct fib\_node \*\*fz\_hash; /\* Hash table pointer \*/

54 The number of entries in this zone

55 int fz\_nent; /\* Number of entries \*/

56 The number of buckets in the hash table associated with this zone (initially set to 16 for all zones but zone 0 in net/ipv4/fib\_hash.c::fn\_new\_zone, but reset to 256 or 1024 in net/ipv4/fib\_hash.c::fn\_rehash\_zone if no. entries grows)

57 int fz\_divisor; /\* Hash divisor \*/

58 Used so we can ensure that a hash value lies in the range [0, fz\_divisor-1] - see net/ipv4/fib\_hash.c::fn\_hash

59 u32 fz\_hashmask; /\* (1&lt;&lt;fz\_divisor) - 1 \*/

60 The index in the parent fn\_hash structure (i.e. 0 to 32)

61 int fz\_order; /\* Zone order \*/

62 This is the netmask for fz\_order = 0, fz\_mask = 0x0000, for for fz\_order = 1, fz\_mask = 0x8000, for fz\_order = 2, fz\_mask = 0xC000, ... , for fz\_order = 32, fz\_mask = 0xFFFF,

63 u32 fz\_mask;

64

65 #define FZ\_HASHMASK(fz) ((fz)-&gt;fz\_hashmask)

```
68 #define FZ_MASK(fz)    ((fz)->fz_mask)
69 };
70
```

**70 struct fib\_node**

71 This is [net/ipv4/fib\\_hash.c::fib\\_node](#). It's an entry in an open hash table that contains details about this particular route.

```
72 struct fib_node
```

```
73 {
```

74 It's an open hash table, so this is the link to the next item on the chain.

```
75     struct fib_node *fn_next;
```

```
76
```

77 Key is set to be the network part of an IP address against which addresses (masked with `fz_mask` from above) will be tested for equality.

```
79     fn_key_t          fn_key;
```

```
80
```

81 When we have something that matches the key, the details about this route are held in `fn_info`. Since many routes will have the the same next hop, this is a pointer to a shared structure

```
83     struct fib_info *fn_info;
```

```
84
```

```
85     u8              fn_tos;
```

```
86     u8              fn_type;
```

```
87     u8              fn_scope;
```

```
88     u8              fn_state;
```

```
89
```

```
90 #define FIB_INFO(f)    ((f)->fn_info)
```

```
91 };
```

```
92
```

**92 struct fib\_info**

93 This is [include/net/ip\\_fib.h::fib\\_info](#) This structure contains data specific to an interface and, therefore, common to many  
94 zones.

```
95 struct fib_info
```

```
96 {  
97     struct fib_info *fib_next;  
98     struct fib_info *fib_prev;
```

99 Index to network protocol (e.g. IP) used for this route.

```
100     int fib_protocol;
```

101 Pointer to next hop information

```
102     struct fib_nh fib_nh[0];
```

103

104 Housekeeping stuff. Since this is a shared structure, we care about reference counting carefully, amongst other things.

```
105     int fib_treeref;
```

```
106     atomic_t fib_clntref;
```

```
107     int fib_dead;
```

```
108     unsigned fib_flags;
```

```
109     u32 fib_prefsrc;
```

```
110     u32 fib_priority;
```

```
111     unsigned fib_metrics[RTAX_MAX];
```

```
112     int fib_nhs;
```

113 <Multipath stuff deleted>

114 <some #defines deleted>



```
115 };  
116
```

**116 struct fib\_nh**

117 This is [include/net/ip\\_fib.h::fib\\_nh](#) Next hop structure – defined in terms of the output device or the IP address of the next  
118 hop gateway.

```
119 struct fib_nh
120 {
121     struct net_device *nh_dev;
122     unsigned          nh_flags;
123     unsigned char     nh_scope;
124     <Multipath, class stuff deleted>
125     int               nh_oif;
126     u32               nh_gw;
127 };
128
```

**128 fib\_lookup**

129

130 This is [net/ipv4/fib\\_rules.c::fib\\_lookup](#) and was called from [ip\\_route\\_input\\_slow](#)

131

132 Different rules can be applied to forwarding to different destinations. There might be a rule prohibiting output, or one saying  
133 that we use NAT, but the simplest one and that with which we're really concerned is simple unicast (RTN\_UNICAST below)

134

```
135 int fib_lookup(const struct rt_key *key, struct fib_result *res)
```

```
136 {
```

```
137     int err;
```

```
138     struct fib_rule *r, *policy;
```

```
139     struct fib_table *tb;
```

140

```
141     u32 daddr = key->dst;
```

```
142     u32 saddr = key->src;
```

143

```
144     FRprintk("Lookup: %u.%u.%u.%u <- %u.%u.%u.%u ",
```

```
145             NIPQUAD(key->dst), NIPQUAD(key->src));
```

```
146     read_lock(&fib_rules_lock);
```

147

148 **Look for the relevant rule associated with this dest. By default we'll unicast.**

```
149     for (r = fib_rules; r; r=r->r_next) {
```

```
150         if (((saddr^r->r_src) & r->r_srcmask) ||
```

```
151         ((daddr^r->r_dst) & r->r_dstmask) ||
152 #ifdef CONFIG_IP_ROUTE_TOS
153         (r->r_tos && r->r_tos != key->tos) ||
154 #endif
155 #ifdef CONFIG_IP_ROUTE_FWMARK
156         (r->r_fwmark && r->r_fwmark != key->fwmark) ||
157 #endif
158         (r->r_ifindex && r->r_ifindex != key->iif))
159         continue;
160
161 FRprintk("tb %d r %d ", r->r_table, r->r_action);
162
163 This is where we decide whether we're going to do something or return an error.
164         switch (r->r_action) {
165             case RTN_UNICAST:
166             case RTN_NAT:
167                 policy = r;
168                 break;
169             case RTN_UNREACHABLE:
170                 read_unlock(&fib_rules_lock);
171                 return -ENETUNREACH;
172             default:
173             case RTN_BLACKHOLE:
174                 read_unlock(&fib_rules_lock);
175                 return -EINVAL;
176             case RTN_PROHIBIT:
```

```
177         read_unlock(&fib_rules_lock);
178         return -EACCES;
179     }
180
```

181 **Given that we've decided we're going to do something, get a handle on the correct FIB.**

```
182     if ((tb = fib_get_table(r->r_table)) == NULL)
183         continue;
184
```

185 **And perform the lookup. In our case, this is set to be [net/ipv4/fib\\_hash.c::fn\\_hash\\_lookup](#) – see below.**

```
186     err = tb->tb_lookup(tb, key, res);
187
188     if (err == 0) {
189         res->r = policy;
190         if (policy)
191             atomic_inc(&policy->r_clntref);
192         read_unlock(&fib_rules_lock);
193         return 0;
194     }
195     if (err < 0 && err != -EAGAIN) {
196         read_unlock(&fib_rules_lock);
197         return err;
198     }
199 }
200 FRprintk("FAILURE\n");
201     read_unlock(&fib_rules_lock);
```

```
202     return -ENETUNREACH;
203 }
204
```

204 **fn\_hash\_lookup**

205

206 This is [net/ipv4/fib\\_hash.c:fn\\_hash\\_lookup](#) and, actually, relatively straightforward. There are a number of static inline  
 207 functions used (all in [net/ipv4/fib\\_hash.c](#))

208	<code>fn_key_t</code>	<code>fz_key(u32 dst, struct fn_zone *fz)</code>	returns an address that has been masked by the netmask for a given zone
209			
210	<code>struct fib_node *</code>	<code>fz_chain(fn_key_t key, struct fn_zone *fz)</code>	hashes the key and returns the head of the chain of node structures that match for this zone
211			
212	<code>int</code>	<code>fn_key_eq(fn_key_t a, fn_key_t b)</code>	Compare keys and say if they're equal
213	<code>int</code>	<code>fn_key_leq(fn_key_t a, fn_key_t b)</code>	Same but for leq

214

215 As a matter of interest, the top two routines are defined thus:

```
216     fn_key_t fz_key(u32 dst, struct fn_zone *fz) {
217         fn_key_t k;  k.datum = dst & FZ_MASK(fz); return k; }
```

218

```
219     fz_chain(fn_key_t key, struct fn_zone *fz) {
220         return fz->fz_hash[fn_hash(key, fz).datum]; }
```

221

222

222 And the all important hash function is defined thus: N.B. in C, ^ is XOR

```
223     fn_hash(fn_key_t key, struct fn_zone *fz) {
224         u32 h = ntohl(key.datum)>>(32 - fz->fz_order);
225         h ^= (h>>20);
226         h ^= (h>>10);
227         h ^= (h>>5);
228         h &= FZ_HASHMASK(fz);
229         return *(fn_hash_idx_t*)&h;
230     }
```

231

232 The algorithm used for lookup is a simple linear search on a series of open hash tables, rather than anything massively  
233 sophisticated. Note that the elements of each chain in a hash table entry are held ordered by key value.

234

```
235 static int
236 fn_hash_lookup(struct fib_table *tb, const struct rt_key *key, struct fib_result
237 *res)
238 {
239     int err;
240     struct fn_zone *fz;
241     struct fn_hash *t = (struct fn_hash*)tb->tb_data;
242
243     read_lock(&fib_hash_lock);
244
```



245 **Start with the most restrictive zone and iterate over zones with smaller and smaller netmasks**

```
246     for (fz = t->fn_zone_list; fz; fz = fz->fz_next) {
247         struct fib_node *f;
```

248

249 **Mask the destination appropriately to produce the lookup key for this zone**

```
250         fn_key_t k = fz_key(key->dst, fz);
```

251

252 **Now, do a hash (implicit in fz\_chain) and walk down the open hash table chain returned looking for a match. As a matter of interest, the hash is defined as:**

```
254         for (f = fz_chain(k, fz); f; f = f->fn_next) {
```

255

256 **Did we find it, or did we go past it? If neither, then keep chaining down.**

```
257             if (!fn_key_eq(k, f->fn_key)) {
258                 if (fn_key_leq(k, f->fn_key))
259                     break;
260                 else
261                     continue;
262             }
```

263

264 **If we come here, we've found something where the keys match. However, we have to be careful, make sure that it's a proper match**

```
266 #ifdef CONFIG_IP_ROUTE_TOS
```

```
267     if (f->fn_tos && f->fn_tos != key->tos)
```

```
268         continue;
```

```
269 #endif
270     f->fn_state |= FN_S_ACCESSED;
271
272     if (f->fn_state & FN_S_ZOMBIE)
273         continue;
274     if (f->fn_scope < key->scope)
275         continue;
276
277 net/ipv4/fib\_semantics.c::fib\_semantic\_match is a routine to make really sure we're allowed to use this interface for this
278 packet; it also fills in some fields in res, notably res->fi which points to the fib_info structure passed as the second arg. If we
279 are allowed to proceed, then fill in a result structure with info about this node and return.
280     err = fib_semantic_match(f->fn_type, FIB_INFO(f), key, res);
281     if (err == 0) {
282         res->type      = f->fn_type;
283         res->scope     = f->fn_scope;
284         res->prefixlen = f->fn_order;
285
286         goto out;
287     }
288     if (err < 0)
289         goto out;
290 }
291 }
292 err = 1;
293 out:
```

```
294     read_unlock(&fib_hash_lock);
295     return err;
296 }
297
298
```