

IP Routing

The `ip_route_output_slow()` function, defined in `net/ipv4/route.c` is the major route resolver. Given a “routing key” as an input parameter, this routine builds a new route cache entry and stores a pointer to it in the parameter `**rp`. A Linux route is defined by (`dst`, `src`, `tos`).

```
1690 int ip_route_output_slow(struct rtable **rp, const
                             struct rt_key *oldkey)
1691 {
1692     struct rt_key key;
1693     struct fib_result res;
1694     unsigned flags = 0;
1695     struct rtable *rth;
1696     struct net_device *dev_out = NULL;
1697     unsigned hash;
1698     int free_res = 0;
1699     int err;
1700     u32 tos;
```

The function uses two important local variables: `key` is of `struct rt_key`, derived from the values pointed to by `oldkey` and is used to specify the characteristics of the desired route;

```
48 struct rt_key
49 {
50     __u32     dst;        /* Destination IP address */
51     __u32     src;        /* Source IP address      */
52     int       iif;        /* Input interface index  */
53     int       oif;        /* Output interface index */
54 #ifdef CONFIG_IP_ROUTE_FWMARK
55     __u32     fwmark;
56 #endif
57     __u8      tos;        /* Requested type of service */
58     __u8      scope;     /* Host, LAN, site, universe */
59 };
```

The variable `res` has type `struct fib_result` and is later used in building the new routing cache entry.

```
86 struct fib_result
87 {
88     unsigned char    prefixlen;
89     unsigned char    nh_sel;
90     unsigned char    type;
91     unsigned char    scope;
92     struct fib_info  *fi;
93 #ifdef CONFIG_IP_MULTIPLE_TABLES
94     struct fib_rule  *r;
95 #endif
96 };
```

The elements of the `fib_result` structure include:

<code>prefixlen</code>	prefix length or equivalently the number of leading 1 bits in the subnet mask
<code>nh_sel</code>	next hop (output dev index). This actually appears under <code>grep -r</code> to be one of the ever popular write only variables!
<code>scope</code>	An indication of the distance to the destination IP address (e.g. host, local network, site, universe). Higher scope values are more specific.
<code>type</code>	type of address (LOCAL, UNICAST, BROADCAST, MULTICAST)
<code>fi</code>	Pointer to the <code>fib_info</code> structure that contains protocol and hardware information specific to an interface.
<code>r</code>	Pointer to a <code>fib_rule</code> structure used for policy based routing.

The fib_rule structure is defined in net/ipv4/fib_rules.c. This structure is the key element defining the existence of a route with a given class of service between a specific source and destination address. **It is not used unless CONFIG_IP_MULTIPLE_TABLES has been defined.**

```
52 struct fib_rule
53 {
54     struct fib_rule *r_next;
55     atomic_t         r_clntref;
56     u32              r_preference;
57     unsigned char    r_table;
58     unsigned char    r_action;
59     unsigned char    r_dst_len;
60     unsigned char    r_src_len;
61     u32              r_src;
62     u32              r_srcmask;
63     u32              r_dst;
64     u32              r_dstmask;
65     u32              r_srcmap;
66     u8               r_flags;
67     u8               r_tos;
68 #ifdef CONFIG_IP_ROUTE_FWMARK
69     u32              r_fwmark;
70 #endif
71     int              r_index;
72 #ifdef CONFIG_NET_CLS_ROUTE
73     __u32            r_tclassid;
74 #endif
75     char             r_ifname[IFNAMSIZ];
76     int              r_dead;
77 };
```

The function `ip_route_output_slow()` begins by constructing the new routing key structure. Manipulation of the `tos` field is somewhat strange. TOS related constants are defined as follows:

```
23 #define IPTOS_TOS_MASK      0x1E
24 #define IPTOS_TOS(tos)     ((tos)&IPTOS_TOS_MASK)
25 #define IPTOS_LOWDELAY      0x10
26 #define IPTOS_THROUGHPUT    0x08
27 #define IPTOS_RELIABILITY   0x04
28 #define IPTOS_MINCOST       0x02

151 #define IPTOS_RT_MASK      (IPTOS_TOS_MASK & ~3)
40 #define RTO_ONLINK        0x01
```

`RTO_ONLINK` is a flag that indicates the destination is no more than one hop away and reachable via a link layer protocol. Thus, in line 1702 the input value of `tos` is being and'ed with `0x1d`. Then in line it is 1705 is and'ed with `0x1c`. Presumably there is a reason for distinguishing `tos` and `key.tos`.

```
1702     tos = oldkey->tos & (IPTOS_RT_MASK | RTO_ONLINK);
1703     key.dst = oldkey->dst;
1704     key.src = oldkey->src;
1705     key.tos = tos & IPTOS_RT_MASK;
```

The input interface identifier is forced to that of the loopback device. The variable `loopback_dev` is an instance of struct `net_device` and is globally defined in `drivers/net/Space.c`. The value of the `ifindex` field is a unique identifier assigned to the interface at initialization time. The output interface identifier is copied from the `oldkey` structure.

```
1706     key.ifindex = loopback_dev.ifindex;
1707     key.oif = oldkey->oif;
```

`CONFIG_IP_ROUTE_FWMARK` is an option to specify different route for packets with different (netfilter) mark values.

```
1708 #ifdef CONFIG_IP_ROUTE_FWMARK
1709     key.fwmark = oldkey->fwmark;
1710 #endif
```

The value of `key.scope` is an indication of the distance from the destination. Here there are only two possible choices, and they depend on the setting of `RTO_ONLINK`. If `RTO_ONLINK` is set then the scope must be `RT_SCOPE_LINK`. Otherwise it is `RT_SCOPE_UNIVERSE`. Thus the scope attribute of the new key does reflect the setting of the `RTO_ONLINK` bit in the `tos` field of the old key.

```
1711     key.scope = (tos & RTO_ONLINK) ? RT_SCOPE_LINK:
1712                   RT_SCOPE_UNIVERSE;
```

As described more fully in the kernel comments below and subsequent data definitions it is clear that a wider range of possible scopes is intended and that the higher the value of scope the more specific the target routing domain.

```
144 /* rtm_scope
145
146 Really not a scope, but sort of distance to the destination.
147 NOWHERE are reserved for non-existing dests, HOST is our
148 local addresses, LINK are dests on directly attached
149 link and UNIVERSE is everywhere in the Universe.
150
151 Intermediate values are also possible f.e. interior routes
152 could be assigned a value between UNIVERSE and LINK.
153 */
```

`RT_SCOPE_LINK`, `RT_SCOPE_UNIVERSE` stand for on-link routes and global routes respectively and are defined in `include/linux/rtnetlink.h`.

```
155 enum rt_scope_t
156 {
157     RT_SCOPE_UNIVERSE=0,
158 /* User defined values */
159     RT_SCOPE_SITE=200,
160     RT_SCOPE_LINK=253,
161     RT_SCOPE_HOST=254,
162     RT_SCOPE_NOWHERE=255
163 };
```

CONFIG_IP_MULTIPLE_TABLES is an option that allows the Linux router to be able to take the packet's source address into account. (Normally, a router decides what to do with a received packet based solely on the packet's final destination address.)

The routing tables are referred to as "classes". Currently, the number of classes is limited to 255, of which three classes are builtin¹:

```
RT_CLASS_LOCAL      = 255 – local interface addresses, broadcasts, nat addresses
RT_CLASS_MAIN       = 254 – all normal routes are put here by default.
RT_CLASS_DEFAULT    = 253 – If the ip_fib_model == 1, then normal default routes
                        are put there. If the ip_fib_model == 2, all gateway routes are
                        put there .
```

```
1713     res. fi          = NULL;

1714 #ifdef CONFIG_IP_MULTIPLE_TABLES
1715     res. r             = NULL;
1716 #endif
```

Check if the source address is defined, and if so determine its type. If the source address is a MULTICAST, BADCLASS and ZERONET address (these macros are defined in include/linux/in.h), return error.

```
182 #define LOOPBACK(x)
    (((x) & htonl(0xff000000)) == htonl(0x7f000000))
183 #define MULTICAST(x)
    (((x) & htonl(0xf0000000)) == htonl(0xe0000000))
184 #define BADCLASS(x)
    (((x) & htonl(0xf0000000)) == htonl(0xf0000000))
185 #define ZERONET(x)
    (((x) & htonl(0xff000000)) == htonl(0x00000000))
186 #define LOCAL_MCAST(x)
    (((x) & htonl(0xfffff00)) == htonl(0xe0000000))

1718     if (oldkey->src) {
1719         err = -EINVAL;
1720         if (MULTICAST(oldkey->src) ||
1721             BADCLASS(oldkey->src) ||
1722             ZERONET(oldkey->src))
1723             goto out;
```

¹ <http://lxr.linux.no/source/Documentation/networking/policy-routing.txt>

The `ip_dev_find()` function looks up the IP source address in the local table and returns a pointer to the `struct net_device` associated with the source address. This function is defined in `net/ipv4/fib_frontend.c`.

```
1725 /* It is equivalent to inet_addr_type(saddr) == RTN_LOCAL */
1726     dev_out = ip_dev_find(ol_dkey->src);
```

The input parameter here is the source IP address associated with the route being setup. Note that the address is subsequently put into the `dst` element of the new key structure that is built.

```
145 struct net_device * ip_dev_find(u32 addr)
146 {
147     struct rt_key key;
148     struct fib_result res;
149     struct net_device *dev = NULL;
150
151     memset(&key, 0, sizeof(key));
152     key.dst = addr;
153 #ifdef CONFIG_IP_MULTIPLE_TABLES
154     res.r = NULL;
155 #endif
156
```

The variable `local_table` is a reference to the statically defined local table.

```
ip_fib.h:
#define local_table (fib_tables[RT_TABLE_LOCAL])

157     if (!local_table ||
158         local_table->tb_lookup(local_table, &key, &res)) {
159 }
```

The call to `local_table->tb_lookup()` is a reference to the `fn_hash_lookup()` function. This function is used to determine if the destination entity identified by `key` exists in the specified table. All the `fib_tables` are searched by zone where a routing zone is the set of routing destinations that have the same length prefix (or equivalently netmask). The `fn_hash_lookup()` searches the specified table, starting with the most specific zone netmask looking for a match. The most specific existing zone is pointed by the `fn_zone_list` variable.

```

268 static int
269 fn_hash_lookup (struct fib_table *tb,
                const struct rt_key *key, struct fib_result *res)
270 {
271     int err;
272     struct fn_zone *fz;
273     struct fn_hash *t = (struct fn_hash*)tb->tb_data;
274

```

This outer loop processes every non-empty zone associated with the `fib_table` in longest prefix first order.

```

275     read_lock(&fib_hash_lock);
276     for (fz = t->fn_zone_list; fz; fz = fz->fz_next) {
277         struct fib_node *f;
278         fn_key_t k = fz_key(key->dst, fz);

```

The `fz_key()` function, defined in `fib_hash.c`, builds a test key by and-ing the address with the zone's netmask. The structures `fn_key_t` and `fn_hash_idx_t` are simply unsigned integers representing IP prefixes and hash table indices respectively.

```

60 typedef struct {
61     u32 datum;
62 } fn_key_t;

64 typedef struct {
65     u32 datum;
66 } fn_hash_idx_t;

123 static __inline__ fn_key_t fz_key(u32 dst, struct
                                fn_zone *fz)
124 {
125     fn_key_t k;
126     k.datum = dst & FZ_MASK(fz);
127     return k;
128 }

```

`FZ_MASK` is a macro defined in `fib_hash.c`

```

97 #define FZ_MASK(fz) ((fz)->fz_mask)

```


On returning to `fn_hash_lookup()`, this inner loop traverses the list of `fib_node` structures associated with the hash bucket of the routing key searching for the first key match. To initiate this process `fz_chain()` is called to retrieve the address of the first `fib_node` in the chain. It performs the hash function `fn_hash()` and ANDs this value with the zone's `fz_hashmask` to get an index into the zone's hash table of nodes. The syntax of this function is a bit dense. Note that `fn_hash` returns `fn_hash_idx_t` which was shown above to be a "structure" consisting of a single unsigned int member called `datum`. That value is used as an index into the hash table structure associated with the routing zone yielding the required pointer to the struct `fib_node`.

```

280         for (f = fz_chain(k, fz); f; f = f->fn_next) {
135 static __inline__ struct fib_node * fz_chain(fn_key_t
                                key, struct fn_zone *fz)
136 {
137     return fz->fz_hash[fn_hash(key, fz).datum];
138 }
110 static __inline__ fn_hash_idx_t fn_hash(fn_key_t key,
                                struct fn_zone *fz)
111 {
112     u32 h = ntohl (key.datum)>>(32 - fz->fz_order);
113     h ^= (h>>20);
114     h ^= (h>>10);
115     h ^= (h>>5);
116     h &= FZ_HASHMASK(fz);

```

`FZ_HASHMASK` is a macro defined in `fib_hash.c`

```

93 #define FZ_HASHMASK(fz) ((fz)->fz_hashmask)

```

`fn_hash_idx_t` is a structure containing the address as its element.

```

64 typedef struct {
65     u32 datum;
66 } fn_hash_idx_t;

117     return *(fn_hash_idx_t*)&h;
118 }

```

The first action of the inner loop is to compare search key with the key of the struct `fib_node`. Recall that the variable `k` is an instance of `fn_key_t`, a structure of the single element datum, whose value was previously set to the target IP address anded with the netmask associated with the zone. From this we can infer that the value of `f->fn_key` is the network address or CIDR network prefix associated with the routing table entity associated with this node. The nodes on any hash queue are sorted in decreasing order by prefix. Therefore, if they do not match and if the search key value is greater than that of the node key, the search continues on to the next node.

```

281         if (!fn_key_eq(k, f->fn_key))
282         {
283             if (fn_key_l eq(k, f->fn_key))
284                 break;
285             else
286                 continue;

```

Arriving here implies that there has been a match. `CONFIG_IP_ROUTE_TOS` makes use of TOS value as routing key and so if there is a TOS associated with the `fib_node` and it is not equal to the TOS of the key, the match is discarded and the search continues.

```

287 #ifdef CONFIG_IP_ROUTE_TOS
288     if (f->fn_tos && f->fn_tos != key->tos)
289         continue;
290 #endif

```

Update and test the state information of the `fib_node`. Zombie nodes are considered non-usable and likely relate to deleted routes or dead interfaces. Very little state information is present in `fib_node`. Only 2 bits are defined:

```

80 #define FN_S_ZOMBIE      1
81 #define FN_S_ACCESSED   2

291         f->fn_state |= FN_S_ACCESSED;
292
293         if (f->fn_state & FN_S_ZOMBIE)
294             continue;

```

Recall that higher values of scope means more specific or constrained routing. Thus the node scope is required to be at least as specific as the requested route scope. If the `fib_node` scope is less than that of the scope of the key, then this node is also not usable.

```

295         if (f->fn_scope < key->scope)
296             continue;
297

```

Finally the `fib_semantic_match()` function is called to ensure that this `fib_node` is usable within the semantic constraints imposed by the route key.

```
298             err = fib_semantic_match(f->fn_type,
                                     FIB_INFO(f), key, res);
```

The `fib_semantic_match()` function is defined in `net/ipv4/fib_semantics.c`. Its mission is to ensure that the candidate `fib_node` appears to represent an acceptable route. The tests include ensuring that the associated `fib_info`'s view of the next hop is that it is alive, the `fib_nh`'s view of the next hop is that its alive, and that if the output interface is specified in the routing key, it is the same interface as the one associated with the next hop structure.

```
569 int
570 fib_semantic_match (int type, struct fib_info *fi, const
                    struct rt_key *key, struct fib_result *res)
571 {
572     int err = fib_props[type].error;
573     if (err == 0) {
```

If the `fib_info` structure indicates that the next hop is dead, then failure is returned.

```
575         if (fi->fib_flags & RTNH_F_DEAD)
576             return 1;
577
```

The `fib_info` structure is connected to the results structure.

```
578         res->fi = fi;
579
580         switch (type) {
581 #ifdef CONFIG_IP_ROUTE_NAT
582         case RTN_NAT:
583             FIB_RES_RESET(*res);
584             atomic_inc(&fi->fib_clntref);
585             return 0;
586 #endif
```

Only the NAT type route is distinguished for the purposes of route semantics.

```
587         case RTN_UNICAST:
588         case RTN_LOCAL:
589         case RTN_BROADCAST:
590         case RTN_ANYCAST:
591         case RTN_MULTICAST:
```

Check if a next hop is feasible from this node. The macros used in this loop depend upon whether or not multipath routing is enabled. If not, there can be only one next hop associated with a fib_info structure.

```

57 #ifndef CONFIG_IP_ROUTE_MULTIPATH
58
59 #define for_nexthops(fi) { int nhsel; const struct fib_nh * nh; \
60 for (nhsel=0, nh = (fi)->fib_nhs; nhsel < (fi)->fib_nhs; nh++, nhsel++)
61
62 #else /* CONFIG_IP_ROUTE_MULTIPATH */
63 /* Hope, that gcc will optimize it to get rid of dummy loop */
64 #define for_nexthops(fi) {int nhsel=0;const struct fib_nh *nh = (fi)->fib_nh; \
65 for (nhsel=0; nhsel < 1; nhsel++)
66
67 #endif /* CONFIG_IP_R0
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000

```

If the route key requires a specific output interface and that is not the output interface associated with this fib_nh then the route is not usable. The break is taken if the route is usable.

```

595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000

```

The CONFIG_IP_ROUTE_MULTIPATH option allows the routing tables to specify alternative paths to travel for a given packet. The router considers all these paths to be of equal "cost" and chooses one of them in a non-deterministic fashion when selecting a route. **How is this done??**

```

598 #ifndef CONFIG_IP_ROUTE_MULTIPATH
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000

```

For non multi-path routing, this is the success return point. The loop will have been exited via the break and so nhsel will remain 0. The reference counter of the fib_info structure is incremented here.

```

605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000

```

This endfor is misleading. The actual loop ended at line 597. This closes the block in which the local variables preceding the for loop are declared.

```
610         endfor_nexthops(fi );
```

Falling out of the loop implies no fib_nh with acceptable semantics was found.

```
611         res->fi = NULL;
612         return 1;
613     default:
614         res->fi = NULL;
615         printk(KERN_DEBUG "impossible 102\n");
616         return -EINVAL;
617     }
618 }
619 return err;
620 }
```

This is the point of return from fib_semantic_match() to fn_hash_lookup(). If the the source address was found to be acceptable, the res structure was filled with the type and scope elements copied from the fib_node structure and the prefix length is copied from the fn_zone structure.

```
299         if (err == 0) {
300             res->type = f->fn_type;
301             res->scope = f->fn_scope;
302             res->prefixlen = fz->fz_order;
303             goto out;
304         }
305         if (err < 0)
306             goto out;
307     }
308 }
309 err = 1;
310 out:
311 read_unlock(&fib_hash_lock);
312 return err;
313 }
314 }
```

Here control returns to `ip_find_dev()`, since the source address is being processed, it is necessary that the returned route type be `RTN_LOCAL`. This seems like one convoluted way to find if a host owns a particular IP address. If the route is not `RTN_LOCAL`, a jump is made to the tag out bypassing the code which normally sets up the return value, `dev`. The value of `dev` was initialized to `NULL`, and a return value of `NULL` will cause `ip_route_output_slow` to return failure.

```
160     if (res.type != RTN_LOCAL)
161         goto out;
```

`FIB_RES_DEV`, a macro defined in `include/net/ip_fib.h`, extracts the struct netdevice pointer from the `fib_info` pointer contained in the results structure. Note that `dev->refcnt` is incremented here. Where the corresponding decrement occurs is not clear at present.

```
113 #define FIB_RES_DEV(res) (FIB_RES_NH(res).nh_dev)
106 #define FIB_RES_NH(res) ((res).fi->fib_nh[0])

162     dev = FIB_RES_DEV(res);
163     if (dev)
164         atomic_inc(&dev->refcnt);
165
```

The `fib_res_put()` function triggers a set of events that is not well understood at present.

```
166 out:
167     fib_res_put(&res);
168     return dev;
169 }
```

What is not well understood here is how routes dynamically become “dead” or come to have reference counts of 0. The best guess at the moment is that the `fib_info` structure is held by all but its creator for a very short interval of time. Nevertheless, it would be possible that whatever owned and normally keeps the reference count at 1 tried to delete the route while we owned it here. Thus when we release it, it really should go away, but qui sait.

```
268 static inline void fib_res_put(struct fib_result *res)
269 {
270     if (res->fi)
271         fib_info_put(res->fi);
272 #ifdef CONFIG_IP_MULTIPLE_TABLES
273     if (res->r)
274         fib_rule_put(res->r);
275 #endif
276 }
```

The fib_clntref is a reference counter and when its value reaches zero, the struct fib_info is deleted. In this context fib_clntref was incremented in the function fib_semantic_match(). The atomic_dec_and_test() function returns true if the value is zero.

```
262 static inline void fib_info_put(struct fib_info *fi)
263 {
264     if (atomic_dec_and_test(&fi->fib_clntref))
265         free_fib_info(fi);
266 }

106 void free_fib_info(struct fib_info *fi)
107 {
108     if (fi->fib_dead == 0) {
109         printk("Freeing alive fib_info %p\n", fi);
110         return;
111     }
```

Unless multipath routing is enabled, change_nexthops() will cause the enclosed block to be executed exactly one time and this fib_info structure's claim on the net_device will be dropped.

```
112     change_nexthops(fi) {
113         if (nh->nh_dev)
114             dev_put(nh->nh_dev);
115         nh->nh_dev = NULL;
116     } endfor_nexthops(fi);
117     fib_info_cnt--;
118     kfree(fi);
119 }
```

Release a fib_rule structure.

```
152 void fib_rule_put(struct fib_rule *r)
153 {
154     if (atomic_dec_and_test(&r->r_clntref)) {
155         if (r->r_dead)
156             kfree(r);
157     } else
158         printk("Freeing alive rule %p\n", r);
159 }
160 }
```

Finally return is made from `ip_find_dev()` to `ip_route_output_slow()`. Recall that we only embarked upon this path if the source IP address was not NULL. If the value of `dev_out` is NULL, then there is no usable network interface associated with the source IP address. The comment below discusses why it is not necessary that the device found here actually map to the output interface specified by the caller. He actually probably means `key.oif == dev_out->oif`.

```

1727         if (dev_out == NULL)
1728             goto out;
1729
1730 /* I removed check for oif == dev_out->oif here.
1731    It was wrong by three reasons:
1732    1. ip_dev_find(saddr) can return wrong i face, if saddr
1733       is assigned to multiple interfaces.
1734    2. Moreover, we are allowed to send packets with saddr
1735       of another i face. --ANK
1736 */

```

Since `oif == 0` means unspecified, what is happening here is a coerced conversion of a multicast and broadcast destination addresses to use the output interface associated with the device that was returned. In addition to the factors discussion below, it is also the case that proper multicast addresses must be associated with a specific interface.

```

1738         if (ol dkey->oif == 0
1739             && (MULTICAST(ol dkey->dst) ||
1740                ol dkey->dst == 0xFFFFFFFF)) {
1741     /* Special hack: user can direct multicasts
1742        and limited broadcast via necessary interface
1743        without fiddling with IP_MULTICAST_IF or IP_PKTINFO.
1744        This hack is not just for fun, it allows
1745        vic, vat and friends to work.
1746        They bind socket to loopback, set ttl to zero
1747        and expect that it will work.
1748        From the viewpoint of routing cache they are broken,
1749        because we are not allowed to build multicast path
1750        with loopback source addr (look, routing cache
1751        cannot know, that ttl is zero, so that packet
1752        will not leave this host and route is valid).
1753        Luckily, this hack is good workaround.
1754     */
1755         key.oif = dev_out->i f i n d e x;
1756         goto make_route;
1757     }

```

Release the device by invoking the `dev_put()` function defined in `include/linux/netdevice.h`

```

1758         if (dev_out)
1759             dev_put(dev_out);
1760         dev_out = NULL;
1761     } /* end if (ol dkey->src) */

```


If an output interface index is specified, attempt to retrieve a pointer to the associated struct net_device. A return value of NULL indicates the device is not found. If the device exists, its reference count is incremented, and the pointer is safe until dev_put is called to release it.

```

1762     if (oldkey->oif) {
1763         dev_out = dev_get_by_index(oldkey->oif);
1764         err = -ENODEV;
1765         if (dev_out == NULL)
1766             goto out;

```

The IPV4 specific data is retrieved by the in_dev_get() function which is defined in include/linux/inetdevice.h. This call returns the void *ip_ptr element of the net_device structure. This pointer points to an instance of struct in_device. Each net_device that supports IPV4 also has an associated struct in_device that carries the IPV4 dependencies of the device layer. An important element of the in_device is the ifa_list pointer. This pointer is the root of a list of struct ifa_list elements.

```

1767     if (__in_dev_get(dev_out) == NULL) {
1768         dev_put(dev_out);
1769         goto out;          /* Wrong error code */
1770     }

```

```

133 __in_dev_get(const struct net_device *dev)
134 {
135     return (struct in_device*)dev->ip_ptr;
136 }
137

```

```

26 struct in_device
27 {
28     struct net_device    *dev;
29     atomic_t             refcnt;
30     rwlock_t             lock;
31     int                  dead;
32     struct in_ifaddr     *ifa_list; /* IP ifaddr chain */
33     struct ip_mc_list    *mc_list; /* IP mcst filter chain */
34     unsigned long        mr_v1_seen;
35     struct neigh_parms   *arp_parms;
36     struct ipv4_devconf  cnf;
37 };

```

Each physical net_device may be assigned alias IP addresses and labels (eth0:1 eth0:2, .. etc). Each alias is represented by an instance of the struct in_ifaddr . The distinction between ifa_local and ifa_address is not well understood. Empirical analysis of “normal” network configurations fails to disclose any instances in which ifa_local and ifa_address differ.

```
60 struct in_ifaddr
61 {
62     struct in_ifaddr    *ifa_next;
63     struct in_device    *ifa_dev;
64     u32                  ifa_local;
65     u32                  ifa_address;
66     u32                  ifa_mask;
67     u32                  ifa_broadcast;
68     u32                  ifa_anycast;
69     unsigned char       ifa_scope;
70     unsigned char       ifa_flags;
71     unsigned char       ifa_preflen;
72     char                 ifa_label [IFNAMSIZ];
73 };
```

When a new interface is created by the inet_rtm_newaddr(struct sk_buff *skb, struct nlmsg_hdr *nlh, void *arg) function in net/ipv4/devinet, the two addresses are set to the values passed in via the netlinks protocol message (don't ask).

```
412     if (rta[IFA_ADDRESS-1] == NULL)
413         rta[IFA_ADDRESS-1] = rta[IFA_LOCAL-1];
414     memcpy(&ifa->ifa_local, RTA_DATA(rta[IFA_LOCAL-1]), 4);
415     memcpy(&ifa->ifa_address, RTA_DATA(rta[IFA_ADDRESS-1]), 4);
416     ifa->ifa_preflen = ifm->ifa_preflen;
```

Continuing along in the code block in which the oif index was explicitly specified, if the destination address is a **LOCAL** multicast address or broadcast address, retrieve the IP address of the output device. Recall that dev_out is a pointer to the struct net_device associated with the explicitly specified output interface. The call to inet_select_address() will return the ifa_local associated with the first interface that is found associated with the net_device that has scope no more restrictive (numerically less than or equal to) than LINK. The use of RT_SCOPE_LINK seems a bit unusual here. It will turn out that the scope is used only for LOCAL MCAST and BCAST. For UCAST destinations the scope will be set to RT_SCOPE_HOST when inet_select_address() is called.

```

1772         if (LOCAL_MCAST(okey->dst) || okey->dst ==
1773             0xFFFFFFFF) {
1774             if (!key.src)
1775                 key.src = inet_select_addr(dev_out, 0,
1776                                         RT_SCOPE_LINK);
1777             goto make_route;
1778         }

181 /* Some random defines to make it easier in the kernel.. */
182 #define LOOPBACK(x)      (((x) & htonl(0xff000000)) == htonl(0x7f000000))
183 #define MULTICAST(x)    (((x) & htonl(0xf0000000)) == htonl(0xe0000000))
184 #define BADCLASS(x)    (((x) & htonl(0xf0000000)) == htonl(0xf0000000))
185 #define ZERONET(x)     (((x) & htonl(0xff000000)) == htonl(0x00000000))
186 #define LOCAL_MCAST(x) (((x) & htonl(0xfffff000)) == htonl(0xe0000000))

```

When the destination address is LOCAL multicast or broadcast, the inet_select_addr() function, defined in net/ipv4/devinet.c, returns the local address associated with the specified output device. In this case dev points to the output device, the dst address is NULL and the scope is RT_SCOPE_LINK. The return value is the selected IP address or is NULL upon failure.

```

718 u32 inet_select_addr(const struct net_device *dev, u32
719                      dst, int scope)
720 {
721     u32 addr = 0;
722     struct net_device *in_dev;
723     read_lock(&inetdev_lock);
724     in_dev = __in_dev_get(dev);
725     if (in_dev == NULL) {
726         read_unlock(&inetdev_lock);
727         return 0;
728     }
729 }

```

At this point `in_dev` points to a valid `in_device` structure. The `for_primary_ifa` macro runs the interface address chain associated with the `in_device`. Recall that routing scope values are ordered with the **most specific scope (i.e. this host) having the highest value**. The scope passed in was `RT_SCOPE_LINK`. Thus interfaces having a more specific address scope (`HOST` or `NOWHERE`) are rejected (for reasons, yet unknown). The address matching logic is with respect to the network mask associated with the `in_ifaddr` structure. In practice it would appear that only 2 distinct values of scope are assigned to interfaces. Scope 0 (`UNIVERSE`) is assigned to physical interfaces and scope 254 (`HOST`) to the loopback interface, `lo`. Thus in the scope matching logic below, physical interfaces are always acceptable and the loopback interface is acceptable only if the input scope is also `HOST`.

```

730     read_lock(&in_dev->lock);
731     for_primary_ifa(in_dev) {
732         if (ifa->ifa_scope > scope)
733             continue;

```

The value of `dst` that was passed in was 0. Therefore `!dst` is true and the value of `addr` is set to the `ifa_local` field of the interface. **Note that the address matching test is against `ifa_address`, but if a match occurs `addr` is set to `ifa_local`.**

```

88 static __inline__ int inet_ifa_match(u32 addr, struct
                                     in_ifaddr *ifa)
89 {
90     return !((addr^ifa->ifa_address)&ifa->ifa_mask);
91 }

734         if (!dst || inet_ifa_match(dst, ifa)) {
735             addr = ifa->ifa_local;
736             break;
737         }
738         if (!addr)
739             addr = ifa->ifa_local;
740     } endfor_ifa(in_dev);
741     read_unlock(&in_dev->lock);
742     read_unlock(&netdev_lock);
743

```

For the control path we are investigating it appears that `addr` should always be non-zero here and thus a return should take place.

```

744     if (addr)
745         return addr;

```

If control should reach here, it indicates that `dst` was non-zero and didn't match the `ifa_address` field of any interface address structure associated with the device. `dev_base` is a global variable pointing to the list of all instances of struct `net_device`. Here the selection criterion appears to be finding an interface whose scope is not LINK and whose scope is numerically less than or equal to the scope that was passed in.

```

746
747 /* Not loopback addresses on loopback should be preferred
748    in this case. It is important that lo is the 1st intf
749    in dev_base list.
750 */
751 read_lock(&dev_base_lock);
752 read_lock(&netdev_lock);
753 for (dev = dev_base; dev; dev = dev->next) {
754     if ((in_dev=__in_dev_get(dev)) == NULL)
755         continue;
756
757     read_lock(&in_dev->lock);
758     for_primary_ifa(in_dev) {
759         if (ifa->ifa_scope != RT_SCOPE_LINK &&
760             ifa->ifa_scope <= scope) {
761             read_unlock(&in_dev->lock);
762             read_unlock(&netdev_lock);
763             read_unlock(&dev_base_lock);
764             return ifa->ifa_local;
765         }
766     } endfor_ifa(in_dev);
767     read_unlock(&in_dev->lock);
768 }
769 read_unlock(&netdev_lock);
770 read_unlock(&dev_base_lock);
771

```

Return failure if an acceptable address cannot be found.

```

772     return 0;
773 }
774

```

Well wasn't that an interesting excursion! Recall that this code block was executed only if the routing key specified an output interface and that the objective was to find an IP source address that is in some sense compatible with previously selected output device. We just dispensed with local multicast and broadcast source addresses. If the destination is general MULTICAST, then the address is selected from the output device using the key's scope. If the destination is unspecified, the scope RT_SCOPE_HOST is passed in.

```

1778         if (!key.src) {
1779             if (MULTICAST(ol dkey->dst))
1780                 key.src = inet_select_addr(dev_out, 0,
1781                                         key.scope);
1782             else if (!ol dkey->dst)
1783                 key.src = inet_select_addr(dev_out, 0,
1784                                         RT_SCOPE_HOST);
1785         }
1786     } /* if (ol dkey->oi f) */
1787

```

If the destination address is unspecified, **the destination is set to the source address** (which is presumably on this machine). If the source is also NULL then they are both set to the loopback address.

```

1788         if (!key.dst) {
1789             key.dst = key.src;
1790
1791             if (!key.dst)
1792                 key.dst = key.src = htonl(INADDR_LOOPBACK);
1793             if (dev_out)
1794                 dev_put(dev_out);
1795

```

Use loopback device for sending packet to this machine.

```

1794         dev_out = &loopback_dev;
1795         dev_hold(dev_out);
1796         key.oi f = loopback_dev.i f i n d e x;
1797         res.type = RTN_LOCAL;
1798         flags |= RTCF_LOCAL;
1799         goto make_route;
1800     }

```

Finally, the function `fib_lookup()` defined in `include/net/ip_fib.h` is invoked to try to resolve the destination address.

```
1802     if (fib_lookup(&key, &res)) {
1803         res.fib = NULL;
```

Since the destination may be on this host as well as elsewhere in the Internet, the `fib_lookup()` function calls `tb_lookup()` on **both** the local table and the main table. Both `tb_lookup` functions resolve to `fn_hash_lookup` which was encountered earlier. Since `fn_hash_lookup()` returns 0 on success and non-zero on failure. The operation fails only if both lookups fail. Theoretically, at least, the lookup should not succeed in both tables but if it does, it would appear that the main table has precedence.

```
155 static inline int fib_lookup(const struct rt_key *key,
156                               struct fib_result *res)
157 {
158     if (local_table->tb_lookup(local_table, key, res) &&
159         main_table->tb_lookup(main_table, key, res))
160         return -ENETUNREACH;
161     return 0;
```

Falling into this implies that the `fib_lookup` failed. Check to see if an output interface was specified (haven't we been here before??) and, if so, get the source address from the device.

```
1804     if (oldkey->oif) {
1805 /* Apparently, routing tables are wrong. Assume,
1806    that the destination is on link.
1807
1808    WHY? DW.
1809    Because we are allowed to send to iface
1810    even if it has NO routes and NO assigned
1811    addresses. When oif is specified, routing
1812    tables are looked up with only one purpose:
1813    to catch if destination is gatewayed, rather than
1814    direct. Moreover, if MSG_DONTROUTE is set,
1815    we send packet, ignoring both routing tables
1816    and ifaddr state. --ANK
1819
1820    We could make it even if oif is unknown,
1821    likely IPv6, but we do not.
1822 */
1823     if (key.src == 0)
1824         key.src = inet_select_addr(dev_out, 0,
1825                                   RT_SCOPE_LINK);
1826     res.type = RTN_UNICAST;
1827     goto make_route;
1828 }
```

```

1829         if (dev_out)
1830             dev_put(dev_out);
1831         err = -ENETUNREACH;
1832         goto out;
1833     }
1834     free_res = 1;
1835

```

probably some error occurred during lookup ??

```

1836         if (res.type == RTN_NAT)
1837             goto e_inval;
1838

```

This packet is routed locally (RTN_LOCAL), so destination and source are the same.

```

1839         if (res.type == RTN_LOCAL) {
1840             if (!key.src)
1841                 key.src = key.dst;
1842             if (dev_out)
1843                 dev_put(dev_out);
1844             dev_out = &loopback_dev;
1845             dev_hold(dev_out);
1846             key.oif = dev_out->oifindex;

```

Release reference into FIB table by calling fib_info_put().

```

1847         if (res.fi)
1848             fib_info_put(res.fi);
1849         res.fi = NULL;
1850         flags |= RTCF_LOCAL;
1851         goto make_route;
1852     }
1853
1854 #ifdef CONFIG_IP_ROUTE_MULTIPATH
1855     if (res.fi->fib_nhs > 1 && key.oif == 0)
1856         fib_select_multipath(&key, &res);
1857     else
1858 #endif

```


If the prefix length is 0 (implying default route), and the type is UNICAST, and no output interface index was specified then its necessary to select among (possibly multiple) default routes.

```
1859     if(!res.prefixlen && res.type == RTN_UNICAST && !key.oif)
1860         fib_select_default(&key, &res);
```

The fib_select_default() function is defined in include/net/ip_fib.h. It appears that it may be a no operation if the if conditions are false. The FIB_RES_GW() macro will return the nh_gw of the next hop structure. As noted earlier three different entities, the node, the next hop, and the interface all have scope values. The value of nh_scope appears to be the most specific, having the value 254 for all local interface entries and local net entries in the main table. It does appear to have a 253 value for those table entries that do specify routing through a gateway either to a remote net or the default route.

```
163 static inline void fib_select_default(const struct rt_key
      *key, struct fib_result *res)
164 {
165     if (FIB_RES_GW(*res) &&
        FIB_RES_NH(*res).nh_scope == RT_SCOPE_LINK)
166         main_table->tb_select_default(main_table, key, res);
167 }
```

This function calls the main table's tb_select_default() which is a reference to the function fn_hash_select_default() defined in net/ipv4/fib_hash.c.

```
340 static void
341 fn_hash_select_default(struct fib_table *tb,
      const struct rt_key *key, struct fib_result *res)
342 {
343     int order, last_idx;
344     struct fib_node *f;
345     struct fib_info *fi = NULL;
346     struct fib_info *last_resort;
347     struct fn_hash *t = (struct fn_hash*)tb->tb_data;
348     struct fn_zone *fz = t->fn_zones[0];
```

fz points to the default netmask (fn_zones[0]). If that zone list is empty, there are no default routes and there is no more that can be done.

```
349
350     if (fz == NULL)
351         return;
352
353     last_idx = -1;
354     last_resort = NULL;
355     order = -1;
356
357     read_lock(&fib_hash_lock);
```

Iterate through all the nodes for the order zero zone. Needless to say this implies the existence of more than one default route. To successfully find something here would require finding a `nh_scope` of `RT_SCOPE_LINK` which we have not seen in our examples.

```

358     for (f = fz->fz_hash[0]; f; f = f->fn_next) {
359         struct fib_info *next_fi = FIB_INFO(f);
360
361         if ((f->fn_state & FN_S_ZOMBIE) ||
362             f->fn_scope != res->scope ||
363             f->fn_type != RTN_UNICAST)
364             continue;
365
366         if (next_fi->fib_priority > res->fib_priority)
367             break;
368         if (!next_fi->fib_nh[0].nh_gw ||
369             next_fi->fib_nh[0].nh_scope != RT_SCOPE_LINK)
370             continue;
371         f->fn_state |= FN_S_ACCESSED;
372
373         if (fi == NULL) {
374             if (next_fi != res->fi)
375                 break;
376         } else if (!fib_detect_death
377                  (fi, order, &last_resort, &last_idx)) {

```

`fib_detect_death()` checks whether the route (checking all nexthops) contains alive paths and whether the route can be used as last resort if there are no valid alternative routes in the group.

```

317 static int fib_detect_death(struct fib_info *fi, int order,
318                             struct fib_info **last_resort, int *last_idx)
319 {
320     struct neighbour *n;
321     int state = NUD_NONE;
322
323     n = neigh_lookup(&arp_tbl, &fi->fib_nh[0].nh_gw,
324                     fi->fib_dev);

```

`neigh_lookup()` function is defined in `net/core/neighbour.c`. This function gets the physical address of the neighbour by matching the IP address of the search key with the one in the table.

```

267 struct neighbour *neigh_lookup (struct neigh_table
268                                 *tbl, const void *pkey,
269                                 struct net_device *dev)
270 {
271     struct neighbour *n;
272     u32 hash_val;
273     int key_len = tbl->key_len;
274
275     hash_val = tbl->hash(pkey, dev);

```

tbl->hash() is a reference to arp_hash() function defined in net/ipv4/arp.c. This function takes the IP address of the neighbor and the associated device and returns a hash value that is used as an index for the hash buckets of the arp table.

```

214 static u32 arp_hash(const void *pkey, const struct net_device
                        *dev)
215 {
216     u32 hash_val;
217
218     hash_val = *(u32*)pkey;
219     hash_val ^= (hash_val >>16);
220     hash_val ^= hash_val >>8;
221     hash_val ^= hash_val >>3;
222     hash_val = (hash_val ^dev->i fi ndex)&NEI GH_HASHMASK;
223
224     return hash_val;
225 }

```

On return to neigh_lookup the search proceeds.

```

275
276     read_l ock_bh(&tbl ->l ock);
277     for (n = tbl ->hash_buckets[hash_val]; n; n = n->next) {
278         if (dev == n->dev &&
279             memcmp(n->primary_key, pkey, key_l en) == 0) {
280             nei gh_hol d(n);
281             break;
282         }
283     }
284     read_unl ock_bh(&tbl ->l ock);
285     return n;
286 }

```

After returning from neigh_lookup() that returns a pointer to an entry in the neighbour table, the state of the neighbor is checked.

```

324     if (n) {
325         state = n->nud_state;
326         nei gh_rel ease(n);
327     }
328     if (state==NUD_REACHABLE)
329         return 0;
330     if ((state & NUD_VALID) && order != fn_hash_l ast_dfl t)
331         return 0;
332     if ((state & NUD_VALID) || (*l ast_i dx<0 && order >
333         fn_hash_l ast_dfl t)) {
334         *l ast_resort = fi;
335         *l ast_i dx = order;
336     }
337     return 1;
338 }

```

We return from fib_detect_death() into fib_select_default()

```
376         if (res->fi)
377             fib_info_put(res->fi);
378         res->fi = fi;
379         atomic_inc(&fib->fib_clntref);
380         fn_hash_last_dfl_t = order;
381         goto out;
382     }
383     fi = next_fi;
384     order++;
385 }
386
387 if (order<=0 || fi==NULL) {
388     fn_hash_last_dfl_t = -1;
389     goto out;
390 }
391
392 if (!fib_detect_death(fi, order, &last_resort,
393     &last_idx)) {
394     if (res->fi)
395         fib_info_put(res->fi);
396     res->fi = fi;
397     atomic_inc(&fib->fib_clntref);
398     fn_hash_last_dfl_t = order;
399     goto out;
400 }
401 if (last_idx >= 0) {
402     if (res->fi)
403         fib_info_put(res->fi);
404     res->fi = last_resort;
405     if (last_resort)
406         atomic_inc(&last_resort->fib_clntref);
407 }
408 fn_hash_last_dfl_t = last_idx;
409 out:
410 read_unlock(&fib_hash_lock);
411 }
```

Finally, back in the main line of `ip_route_output_slow()` a check is made to see if the source IP address remains NULL.

```
1861
1862     if (!key.src)
```

If so an attempt is made to derive the source address from the `fib_prefsrc` field of the `fib_info` structure. If that field is also NULL then our old friend `inet_select_addr()` is asked to recover it from the `net_device` and `nh_gw` parameters. **This makes no sense to me because the value of `nh_gw` should be an IP address that is owned by a different host!**

```
1863         key.src = FIB_RES_PREFSRC(res);
```

`FIB_RES_PREFSRC` is a macro defined in `include/net/ip_fib.h`

```
111 #define FIB_RES_PREFSRC(res) ((res).fib_prefsrc ? :
    __fib_res_prefsrc(&res))

624 u32 __fib_res_prefsrc(struct fib_result *res)
625 {
626     return inet_select_addr(FIB_RES_DEV(*res),
    FIB_RES_GW(*res), res->scope);
627 }
```

If a net device is held in `dev_out`, release it here.

```
1864
1865     if (dev_out)
1866         dev_put(dev_out);
```

Set the value of `key.oif` from the `net_device` pointed to by the `fib_info` structure than lives in the `res` structure.

```
1867     dev_out = FIB_RES_DEV(res);
1868     dev_hold(dev_out);
1869     key.oif = dev_out->iifindex;
```

Now we are ready to create the route and add it to the route cache. After filling in the appropriate data, we determine the hash id and install the new route in the cache.

First ensure that if the source address is a loopback address then the selected output device carries the IFF_LOOPBACK flag. **Couldn't this have been done earlier???**

```
1871 make_route:
1872     if (LOOPBACK(key.src) &&
           !(dev_out->flags & IFF_LOOPBACK))
1873         goto e_inval;
1874
1875     if (key.dst == 0xFFFFFFFF)
1876         res.type = RTN_BROADCAST;
1877     else if (MULTICAST(key.dst))
1878         res.type = RTN_MULTICAST;
1879     else if (BADCLASS(key.dst) || ZERONET(key.dst))
1880         goto e_inval;
1881
1882     if (dev_out->flags & IFF_LOOPBACK)
1883         flags |= RTCF_LOCAL;
1884
```

If the result type is BROADCAST, then any fib_info structure that is held is released.

```
1885     if (res.type == RTN_BROADCAST) {
1886         flags |= RTCF_BROADCAST | RTCF_LOCAL;
1887         if (res.fi) {
1888             fib_info_put(res.fi);
1889             res.fi = NULL;
1890         }
1891     } else if (res.type == RTN_MULTICAST) {
1892         flags |= RTCF_MULTICAST | RTCF_LOCAL;
1893         read_lock(&netdev_lock);
1894         if (!__in_dev_get(dev_out) ||
1895             !ip_check_mc(__in_dev_get(dev_out),
                           oldkey->dst))
1896             flags &= ~RTCF_LOCAL;
1897         read_unlock(&netdev_lock);
1898         /* If multicast route do not exist use
1899            default one, but do not gateway in
1900            this case. Yes, it is hack.
1901            */
1902         if (res.fi && res.prefixlen < 4) {
1903             fib_info_put(res.fi);
1904             res.fi = NULL;
1905         }
1906     }
```

Allocate memory from the slab allocator for a route cache entry.

```
1908     rth = dst_alloc(&pv4_dst_ops);
1909     if (!rth)
1910         goto e_nobufs;
1911
1912     atomic_set(&rth->u.dst.__refcnt, 1);
```

Copy (most of) the elements of the key structure that was used to create the route to the key structure embedded in the rth.

```
1913     rth->u.dst.flags = DST_HOST;
1914     rth->key.dst      = oldkey->dst;
1915     rth->key.tos      = tos;
1916     rth->key.src      = oldkey->src;
1917     rth->key.iif      = 0;
1918     rth->key.oif      = oldkey->oif;
1919     #ifdef CONFIG_IP_ROUTE_FWMARK
1920     rth->key.fwmark   = oldkey->fwmark;
1921     #endif
```

Copy the elements used to route the packet to the rt_ fields of the route cache element.

```
1922     rth->rt_dst       = key.dst;
1923     rth->rt_src       = key.src;
1924     #ifdef CONFIG_IP_ROUTE_NAT
1925     rth->rt_dst_map   = key.dst;
1926     rth->rt_src_map   = key.src;
1927     #endif
1928     rth->rt_iif = oldkey->oif ? : dev_out->iifindex;
1929     rth->u.dst.dev     = dev_out;
1930     dev_hold(dev_out);
1931     rth->rt_gateway   = key.dst;
1932     rth->rt_spec_dst   = key.src;
1933
```

Setup the function that will be used to transmit the packet.

```
1934     rth->u.dst.output = ip_output;
1935
1936     rt_cache_stat[smp_processor_id()].out_slow_tot++;
1937
```

If the flags indicate that this route terminates on this machine, then the input handler is set to `ip_local_deliver`.

```
1938     if (flags & RTCF_LOCAL) {
1939         rth->u.dst.input = ip_local_deliver;
1940         rth->rt_spec_dst = key.dst;
1941     }
1942     if (flags & (RTCF_BROADCAST | RTCF_MULTICAST)) {
1943         rth->rt_spec_dst = key.src;
1944         if (flags & RTCF_LOCAL &&
1945             !(dev_out->flags & IFF_LOOPBACK)) {
1946             rth->u.dst.output = ip_mc_output;
1947             rt_cache_stat
1948                 [smp_processor_id()].out_slow_mc++;
1949         }
1950     }
```

`CONFIG_IP_MROUTE` option is used if you want your machine to act as a router for IP packets that have multicast destination addresses.

```
1948 #ifdef CONFIG_IP_MROUTE
1949     if (res.type == RTN_MULTICAST) {
1950         struct in_device *in_dev =
1951             in_dev_get(dev_out);
1952         if (in_dev) {
1953             if (IN_DEV_MFORWARD(in_dev) &&
1954                 !LOCAL_MCAST(ol_dkey->dst)) {
1955                 rth->u.dst.input = ip_mr_input;
1956                 rth->u.dst.output = ip_mc_output;
1957             }
1958             in_dev_put(in_dev);
1959         }
1960 #endif
1961     }
1962 }
```


Call `rt_set_nexthop()` defined in `net/ipv4/route.c` to set next neighbor parameters like `pmtu` and `mss`.

```
1963     rt_set_nexthop(rth, &res, 0);

1180 static void rt_set_nexthop(struct rtable *rt, struct
                                fib_result *res, u32 itag)
1181 {
1182     struct fib_info *fi = res->fi;
1183
```

The bulk of this code seems to be attempting to address potential problems associated with missing or invalid elements in the `fib_info` structure.

```
1184     if (fi) {
1185         if (FIB_RES_GW(*res) &&
1186             FIB_RES_NH(*res).nh_scope == RT_SCOPE_LINK)
1187             rt->rt_gateway = FIB_RES_GW(*res);
1188         memcpy(&rt->u.dst.mxlck, fi->fib_metrics,
1189              sizeof(fi->fib_metrics));
```

`fib_mtu` is actually a macro referencing the `RTAX_MTU` element of the `fib_metrics` array. If the value is zero it is copied from the net device. **Oddly, it appears that `rt->u.dst.pmtu` has not been previously set in this module... so it is also set in the else clause!**

```
1190         if (fi->fib_mtu == 0) {
1191             rt->u.dst.pmtu = rt->u.dst.dev->mtu;
1192             if (rt->u.dst.mxlck & (1 << RTAX_MTU) &&
1193                 rt->rt_gateway != rt->rt_dst &&
1194                 rt->u.dst.pmtu > 576)
1195                 rt->u.dst.pmtu = 576;
1196         }
1197 #ifdef CONFIG_NET_CLS_ROUTE
1198     rt->u.dst.tclassid = FIB_RES_NH(*res).nh_tclassid;
1199 #endif
1200     } else
1201         rt->u.dst.pmtu = rt->u.dst.dev->mtu;
1202
1203     if (rt->u.dst.pmtu > IP_MAX_MTU)
1204         rt->u.dst.pmtu = IP_MAX_MTU;
1205     if (rt->u.dst.advmss == 0)
1206         rt->u.dst.advmss = max_t(unsigned int,
1207                                 rt->u.dst.dev->mtu - 40,
1208                                 ip_rt_min_advmss);
```

```

1208         if (rt->u.dst.advmss > 65535 - 40)
1209             rt->u.dst.advmss = 65535 - 40;
1210
1211 #ifdef CONFIG_NET_CLS_ROUTE
1212 #ifdef CONFIG_IP_MULTIPLE_TABLES
1213     set_class_tag(rt, fib_rules_tclass(res));
1214 #endif
1215     set_class_tag(rt, itag);
1216 #endif
1217     rt->rt_type = res->type;
1218 }

```

On return to `ip_route_output_slow()`, use the source address, destination address, and tos to determine and return a hash value by invoking the `rt_hash_code()` function defined in `net/ipv4/route.c`. We had visited this function earlier in UDP connect and was called by the `ip_route_output_key()` function.

```

1965     rth->rt_flags = flags;
1967     hash = rt_hash_code(ol dkey->dst, ol dkey->src ^
                        (ol dkey->oi f << 5), tos);

```

The hash code returned is used by `rt_intern_hash()` function to search in the respective hash queue of routing cache (`rt_hash_table`) to find an entry that matches the entry that was just created. The `rp` parameter was passed in to `ip_route_output_slow()` as the location at which a pointer to the new route cache entry should be returned.

```

1968     err = rt_intern_hash(hash, rth, rp);

601 static int rt_intern_hash(unsigned hash, struct rtable
                        *rt, struct rtable **rp)
602 {
603     struct rtable *rth, **rthp;
604     unsigned long now = jiffies;
605     int attempts = !in_softirq();
606
607 restart:

```

Recall that the route cache is based upon a table of structures. Each structure contains a pointer to the first struct rtable element in the hash queue and a lock for the queue. Here the queue is locked and the value of rthp is set to point to the chain header (as opposed to set to the chain header!) As the while loop continues rthp will be advanced.

```

608     rthp = &rt_hash_table[hash].chain;
609
610     write_lock_bh(&rt_hash_table[hash].lock);

```

This loop appears to be looking for the possible case that the route already exists! This could conceivably occur due to race conditions involving multiple callers of ip_route_output(). If an existing entry with the same key is found, the existing entry is used and the newly created one is dropped.

```

611     while ((rth = *rthp) != NULL) {
612         if (memcmp(&rth->key, &rt->key,
613                 sizeof(rt->key)) == 0) {
614             /* Put it first */
615             *rthp = rth->u.rt_next;
616             rth->u.rt_next = rt_hash_table[hash].chain;
617             rt_hash_table[hash].chain = rth;

```

Update the reference count and the last use of the existing entry.

```

618             rth->u.dst.__use++;
619             dst_hold(&rth->u.dst);
620             rth->u.dst.lastuse = now;
621
622             write_unlock_bh(&rt_hash_table
623                             [hash].lock);
624
625             rt_drop(rt);
626             *rp = rth;
627             return 0;
628         }
629     }
630     rthp = &rth->u.rt_next;

```

```

631     /* Try to bind route to arp only if it is output
632        route or unicast forwarding path.
633     */
634     if (rt->rt_type == RTN_UNICAST || rt->key.iif == 0) {
635         int err = arp_bind_neighbour(&rt->u.dst);

```

The `arp_bind_neighbour()` function defined in `net/ipv4/arp.c` is invoked. This function tries to locate an entry in the ARP table for the destination address, if one exists.

```

429 int arp_bind_neighbour(struct dst_entry *dst)
430 {
431     struct net_device *dev = dst->dev;
432     struct neighbour *n = dst->neighbour;
433
434     if (dev == NULL)
435         return -EINVAL;
436     if (n == NULL) {
437         u32 nexthop = ((struct rtable*)dst)->rt_gateway;
438         if (dev->flags & (IFF_LOOPBACK | IFF_POINTOPOINT))
439             nexthop = 0;
440         n = __neigh_lookup_errno(
441 #ifdef CONFIG_ATM_CLIP
442             dev->type == ARPHRD_ATM ? &clip_tbl :
443 #endif
444             &arp_tbl, &nexthop, dev);

```

`neigh_lookup_errno()` is defined in `include/net/neighbour.h`

```

266 static inline struct neighbour *
267 _neigh_lookup_errno(struct neighbour_table *tbl, const void *pkey,
268                    struct net_device *dev)
269 {
270     struct neighbour *n = neigh_lookup(tbl, pkey, dev);
271
272     if (n)
273         return n;
274
275     return neigh_create(tbl, pkey, dev);
276 }

```

```

445     if (IS_ERR(n))
446         return PTR_ERR(n);
447     dst->neighbour = n;
448 }
449 return 0;
450 }

```

```

636         if (err) {
637             write_unlock_bh(&rt_hash_table[hash].lock);
638
639             if (err != -ENOBUFS) {
640                 rt_drop(rt);
641                 return err;
642             }
644 /* Neighbour tables are full and nothing
645    can be released. Try to shrink route cache,
646    it is most likely it holds some neighbour records.
647 */
648
649         if (attempts-- > 0) {
650             int saved_elasticity =
651                 ip_rt_gc_elasticity;
652             int saved_int = ip_rt_gc_min_interval;
653             ip_rt_gc_elasticity = 1;
654             ip_rt_gc_min_interval = 0;
655             rt_garbage_collect();
656             ip_rt_gc_min_interval = saved_int;
657             ip_rt_gc_elasticity = saved_elasticity;
658             goto restart;
659         }
660
661         if (net_ratelimit())
662             printk("Neighbour table overflow.\n");
663         rt_drop(rt);
664         return -ENOBUFS;
665     }

```

Here the new route is inserted at the head of the hash queue.

```
666         rt->u.rt_next = rt_hash_table[hash].chain;
667 #if RT_CACHE_DEBUG >= 2
668     if (rt->u.rt_next) {
669         struct rtable *trt;
670         printk("rt_cache @%02x: %u.%u.%u.%u", hash,
671             NIPQUAD(rt->rt_dst));
672         for (trt = rt->u.rt_next; trt;
673             trt = trt->u.rt_next)
674             printk(" . %u.%u.%u.%u", NIPQUAD(trt->rt_dst));
675     }
676 #endif
677     rt_hash_table[hash].chain = rt;
678     write_unlock_bh(&rt_hash_table[hash].lock);
679     *rp = rt;
680     return 0;
681 }
```

Release reference to FIB table and the device, if holding.

```
1969 done:
1970     if (free_res)
1971         fib_res_put(&res);
1972     if (dev_out)
1973         dev_put(dev_out);
1974 out: return err;
1975
1976 e_inval:
1977     err = -EINVAL;
1978     goto done;
1979 e_nobufs:
1980     err = -ENOBUFS;
1981     goto done;
1982 }
```

To summarize, the `ip_route_output_slow()` function does the following:

- Creates a routing table cache key

- If the source address is specified, calls `ip_dev_find()` to determine the output device.

- If the oif is specified, use `dev_get_by_index` to retrieve output device and select source address (if the dest address was not NULL (p.22).

- If the destination address is not known, set up loopback

- Calls `fib_lookup()` to find route to destination.

- Allocates memory for new routing cache entry and initializes it.

- Calls `rt_setnexthop()` to find next destination.

- Returns `rt_intern_hash()`, which creates a new route in the routing cache.

