

Ethernet Device Initialization

The `init_etherdev()` function is called by most Ethernet drivers at initialization time to initialize and possibly allocate a `net_device` structure. It is a convenience front end that forwards the call to the more generic `init_netdev()` providing the naming string `eth%d` that is eventually used in the construction of the interface name `ethn`. This module resides in `drivers/net/net_init.c`. If the `dev` parameter is `NULL`, then the `net_device` structure is allocated and initialized. Otherwise it is only initialized.

```
210 struct net_device *init_etherdev(struct net_device *dev,
211                                 int sizeof_priv)
212 {
213     return init_netdev(dev, sizeof_priv, "eth%d",
214                       ether_setup);
215 }

129 static struct net_device *init_netdev(struct net_device *dev,
130                                       int sizeof_priv, char *mask,
131                                       void (*setup)(struct net_device *))
132 {
133     int new_device = 0;
134     /*
135     *Allocate a device if one is not provided.
136     */
137
138     if (dev == NULL) {
139         dev=init_alloc_dev(sizeof_priv);
140         if(dev==NULL)
141             return NULL;
142         new_device = 1;
143     }
144
```

On return to `init_netdev()` an attempt is made to try to allocate a name. For ethernet devices, `mask` points to the string `eth%d`.

```
145 /*
146 *Allocate a name
147 */
148
149     if (dev->name[0] == '\0' || dev->name[0] == ' ') {
150         strcpy(dev->name, mask);
151         if (dev_alloc_name(dev, mask)<0) {
152             if (new_device)
153                 kfree(dev);
154             return NULL;
155         }
156     }
157
```

On return to `init_netdev()`, the call at line 158 is used to copy any boot time options. This is a vestige of the bad old days of ISA NICs in which it might be necessary to supply I/O addresses and irqs manually.

```
158     netdev_boot_setup_check(dev);
159
```

The `setup()` function referenced here is an input parameter to this routine. For ethernet devices the call is actually to `ether_setup()`.

```
160 /*
161  *Configure via the caller provided setup function then
162  *register if needed.
163  */
164
165     setup(dev);
166
```

Finally, if this was not a previously allocated device, `init_netdev()` attempts to register the device with the protocols.

```
167     if (new_device) {
168         int err;
169
170         rtnl_lock();
171         err = register_netdevice(dev);
172         rtnl_unlock();
173
174         if (err < 0) {
175             kfree(dev);
176             dev = NULL;
177         }
178     }
179     return dev;
180 }
181
```

Allocation of the *net_device* structure

The `init_alloc_dev()` function allocates a struct `net_device` plus whatever memory is requested for the device driver private area. The private area is a 32 byte aligned structure that follows the struct `net_device` and is accessed via the `dev->priv` pointer.

```
101 static struct net_device *init_alloc_dev(int sizeof_priv)
102 {
103     struct net_device *dev;
104     int alloc_size;
105
106     /* ensure 32-byte alignment of the private area */
107     alloc_size = sizeof(*dev) + sizeof_priv + 31;
108
109     dev = (struct net_device *) kmalloc(alloc_size,
110                                       GFP_KERNEL);
111     if (dev == NULL)
112     {
113         printk(KERN_ERR "init_alloc_dev: Unable to allocate
114                device memory.\n");
115         return NULL;
116     }
117 }
```

Initialize the `net_device` structure and link in the private area if it was provided.

```
116     memset(dev, 0, alloc_size);
117
118     if (sizeof_priv)
119         dev->priv = (void *) (((long)(dev + 1) + 31) & ~31);
120
121     return dev;
122 }
123
```

The `dev_alloc_name()` function attempts to allocate an available device name in the space `name0` through `name99` where the name of an Ethernet device is `eth`. It uses a serial search through the namespace which terminates when an unallocated name is encountered.

```

560 int dev_alloc_name(struct net_device *dev, const char *name)
561 {
562     int i;
563     char buf[32];
564     char *p;
565
566     /*
567      * Verify the string as this thing may have come from
568      * the user. There must be either one "%d" and no other "%"
569      * characters, or no "%" characters at all.
570     */
571     p = strchr(name, '%');
572     if (p && (p[1] != 'd' || strchr(p+2, '%')))
573         return -EINVAL;
574

```

Here a series of names (e.g., `eth0`, `eth1`, `eth2`, ...) are generated. The name parameter that is passed to the `sprintf()` function would be `eth%d` in this case, and the value of `i` the index. The `__dev_get_by_name()` function returns `NULL` if the name is not already in use.

```

575     /*
576      * If you need over 100 please also fix the algorithm...
577     */
578     for (i = 0; i < 100; i++) {
579         sprintf(buf, sizeof(buf), name, i);
580         if (__dev_get_by_name(buf) == NULL) {
581             strcpy(dev->name, buf);
582             return i;
583         }
584     }
585     return -ENFILE; /* Over 100 of the things .. bail out! */
586 }
587
411 struct net_device *__dev_get_by_name(const char *name)
412 {
413     struct net_device *dev;
414
415     for (dev = dev_base; dev != NULL; dev = dev->next) {
416         if (strncmp(dev->name, name, IFNAMSIZ) == 0)
417             return dev;
418     }
419     return NULL;
420 }

```

The ether_setup() function just fills in spots in the dev structure that are common to all ethernet drivers. The labels of the form eth_ are references to functions the live in net/ethernet/eth.c

```
405 void ether_setup(struct net_device *dev)
406 {
407 /* Fill in the fields of the device structure with ethernet-
generic values.
408 This should be in a common file instead of per-driver. */
409
410     dev->change_mtu = eth_change_mtu;
411     dev->hard_header = eth_header;
412     dev->rebuild_header = eth_rebuild_header;
413     dev->set_mac_address = eth_mac_addr;
414     dev->hard_header_cache = eth_header_cache;
415     dev->header_cache_update = eth_header_cache_update;
416     dev->hard_header_parse = eth_header_parse;
417
418     dev->type= ARPHRD_ETHER;
419     dev->hard_header_len = ETH_HLEN;
420     dev->mtu= 1500; /* eth_mtu */
421     dev->addr_len= ETH_ALEN;
422     dev->tx_queue_len= 100; /* Ethernet wants good queues */
423
424     memset(dev->broadcast, 0xFF, ETH_ALEN);
425
```

Indicate that the interface supports both hardware broadcast and multicast.

```
426 /* New-style flags. */
427     dev->flags= IFF_BROADCAST|IFF_MULTICAST;
428 }
```

Registering a net_device

We shall see that the call to `register_netdevice()` will trigger indirect calls to a relatively large collection of change-of-state handlers associated with various protocol routines. The recipients of these indirect calls will in turn use the `netlinks` interface to send routing updates to the FIB manager.

```
2488 int register_netdevice(struct net_device *dev)
2489 {
2490     struct net_device *d, **dp;
2491     #ifdef CONFIG_NET_DIVERT
2492     int ret;
2493     #endif
2494
2495     spin_lock_init(&dev->queue_lock);
2496     spin_lock_init(&dev->xmit_lock);
2497     dev->xmit_lock_owner = -1;
2498     #ifdef CONFIG_NET_FASTROUTE
2499     dev->fastpath_lock = RW_LOCK_UNLOCKED;
2500     #endif
2501
```

The value of `dev_boot_phase` is statically initialized to 1 at compile time. It is reset to 0 during the call to `net_dev_init()` ensuring that the code is executed exactly once at boot time when the first network driver initializes.

```
2502     if (dev_boot_phase)
2503         net_dev_init();
2504
```

It would be nice to understand diverters.

```
2505 #ifdef CONFIG_NET_DIVERT
2506     ret = alloc_dvert_blk(dev);
2507     if (ret)
2508         return ret;
2509 #endif /* CONFIG_NET_DIVERT */
2510
```

The iflink element is an alternative identifying index that can be set by the device driver. It is initialized to -1 before calling the drivers initialization routine, dev->init(). If control reaches this point via any path (including the init_etherdev() path) which includes the creation of the net_device structure, dev->init() will necessarily be NULL. Device drivers which allocate the net_device structure and later register can specific a callback that will be activated here.

```

2511     dev->i f l i n k = -1;
2512
2513 /* Init, if this function is available */
2514     if (dev->i n i t && dev->i n i t(dev) != 0) {
2515 #i fdef CONFIG_NET_DIVERT
2516         free_d i v e r t_b l k(dev);
2517 #endi f
2518         return -EIO;
2520

```

Each interface is given a unique identifier number. This number is also inherited by dev->iflink if the device driver didn't provide the info in its init routine.

```

2521     dev->i f i n d e x = dev_new_i n d e x();
2522     if (dev->i f l i n k == -1)
2523         dev->i f l i n k = dev->i f i n d e x;
2524
2525 /* Check for existence, and append to tail of chain */
2526     for (dp=&dev_base; (d=*dp) != NULL; dp=&d->next) {
2527         if (d == dev || strcmp(d->name, dev->name) == 0) {
2528 #i fdef CONFIG_NET_DIVERT
2529             free_d i v e r t_b l k(dev);
2530 #endi f
2531             return -EEXIST;
2532         }
2533     }
2534 /*
2535  *nil rebuild_header routine,
2536  *that should be never called and used as just bug trap.
2537  */
2538
2539     if (dev->r e b u i l d_h e a d e r == NULL)
2540         dev->r e b u i l d_h e a d e r = default_r e b u i l d_h e a d e r;
2541

```

```

2542 /*
2543  *Default initial state at registry is that the
2544  *device is present.
2545  */
2546
2547     set_bit(__LINK_STATE_PRESENT, &dev->state);
2548
2549     dev->next = NULL;
2550     dev_init_scheduler(dev);
2551     write_lock_bh(&dev_base_lock);
2552     *dp = dev;
2553     dev_hold(dev);
2554     dev->deadbeaf = 0;
2555     write_unlock_bh(&dev_base_lock);
2556

```

The call to `notifier_call_chain()` results in a call to the callback function associated with every notifier block in the `netdev_chain` passing them the event code `NETDEV_REGISTER` and a pointer to the struct `netdevice`. As will be shown there may be twenty or more such functions, but in this case they really don't do very much.

```

2557 /* Notify protocols, that a new device appeared. */
2558     notifier_call_chain(&netdev_chain, NETDEV_REGISTER, dev);
2559
2560     net_run_sbinfo_hotplug(dev, "register");
2561
2562     return 0;
2563 }

```


Device scheduler initialization

The scheduler is initially configure to support no queuing at all. This appears to get rectified during the call to `dev_activate()`.

```
486 void dev_init_scheduler(struct net_device *dev)
487 {
488     write_lock(&qdisc_tree_lock);
489     spin_lock_bh(&dev->queue_lock);
490     dev->qdisc = &noop_qdisc;
491     spin_unlock_bh(&dev->queue_lock);
492     dev->qdisc_sleeping = &noop_qdisc;
493     dev->qdisc_list = NULL;
494     write_unlock(&qdisc_tree_lock);
495
496     dev_watchdog_init(dev);
497 }
498
```

Notifier Chains

The notifier chain facility is a general mechanism provided by the kernel. It is designed to provide a way for kernel elements to express interest in being informed about the occurrence of general asynchronous events. The basic building block of the mechanism is the struct `notifier_block` which is defined in `include/linux/notifier.h`. The block contains a pointer to the function to be called when the event occurs. The parameters passed to the function include:

- a pointer to the notifier block itself,
- an event code such as `NETDEV_REGISTER` or `NETDEV_UNREGISTER`,
- and a pointer to an unspecified private data type which in the case of the network chain points to the associated struct `netdevice`.

```
14 struct notifier_block
15 {
16     int (*notifier_call)(struct notifier_block *self,
17                          unsigned long, void *);
17     struct notifier_block *next;
18     int priority;
19 };

181
182 static struct notifier_block *netdev_chain=NULL;
183
```

The kernel function `notifier_chain_register()` assembles related notifier blocks into notifier chains. Modules within the networking subsystem use the `register_netdevice_notifier()` function defined in `net/core/dev.c` to add their own notifier blocks to the `netdev_chain` which is statically initialized as `NULL` in `dev.c`.

```
850 int register_netdevice_notifier(struct notifier_block *nb)
851 {
852     return notifier_chain_register(&netdev_chain, nb);
853 }

181
182 static struct notifier_block *netdev_chain=NULL;
183
```

Adding the notifier_block to the chain.

The kernel routine notifier_chain_register() links the notifier block into the specified chain in priority order.

```
63
64 int notifier_chain_register(struct notifier_block **list,
        struct notifier_block *n)
65 {
66     write_lock(&notifier_lock);
67     while(*list)
68     {
69         if(n->priority > (*list)->priority)
70             break;
71         list= &((*list)->next);
72     }
73     n->next = *list;
74     *list=n;
75     write_unlock(&notifier_lock);
76     return 0;
77 }
```

Here are the notifiers associated with net_device events.

```
41 /* netdevice notifier chain */
42 #define NETDEV_UP 0x0001
        /* For now you can't veto a device up/down */
43 #define NETDEV_DOWN 0x0002
44 #define NETDEV_REBOOT 0x0003
        /* Tell a protocol stack a network interface
45     detected a hardware crash and restarted
46     - we can use this eg to kick tcp sessions
47     once done */
48 #define NETDEV_CHANGE 0x0004
        /* Notify devstate change */
49 #define NETDEV_REGISTER 0x0005
50 #define NETDEV_UNREGISTER 0x0006
51 #define NETDEV_CHANGEMTU 0x0007
52 #define NETDEV_CHANGEADDR 0x0008
53 #define NETDEV_GOING_DOWN 0x0009
54 #define NETDEV_CHANGE_NAME 0x000A
55
```

Invoking notifier_call_chain()

When a function such as `netdev_init()` makes the call to `notifier_call_chain()`, it results in a callback being made for every notifier block that is in the chain. These notifier callback functions typically contain a `switch()` block which they use to select and process only those event types in which they are interested.

```
2557 /* Notify protocols, that a new device appeared. */
2558     notifier_call_chain(&netdev_chain, NETDEV_REGISTER, dev);
```

This structure is illustrated below in the `rtnetlink_event()` callback.

```
487 static int rtnetlink_event(struct notifier_block *this,
                             unsigned long event, void *ptr)
488 {
489     struct net_device *dev = ptr;
490     switch (event) {
491     case NETDEV_UNREGISTER:
492         rtmsg_info(RTM_DELLINK, dev, ~0U);
493         break;
494     case NETDEV_REGISTER:
495         rtmsg_info(RTM_NEWLINK, dev, ~0U);
496         break;
497     case NETDEV_UP:
498     case NETDEV_DOWN:
499         rtmsg_info(RTM_NEWLINK, dev,
                    IFF_UP|IFF_RUNNING);
500         break;
501     case NETDEV_CHANGE:
502     case NETDEV_GOING_DOWN:
503         break;
504     default:
505         rtmsg_info(RTM_NEWLINK, dev, 0);
506         break;
507     }
508     return NOTIFY_DONE;
509 }
510
```

The entire collection of callers of `register_netdevice_notifier()` is quite large. Each of the modules shown below has a callback function in the netdev chain. However, only the notifiers shown in red have any impact on IP_V4.

Referenced (in 35 files total) in:

- include/linux/netdevice.h, line 454
- net/netsyms.c, line 465
- net/appletalk/aarp.c, line 859
- net/appletalk/ddp.c, line 1974
- net/ax25/af_ax25.c, line 1851
- net/core/dev.c, line 850
- net/core/dst.c, line 214
- net/core/rtnetlink.c, line 526
- net/ipv4/devinet.c, line 1140
- net/ipv4/ipmr.c, line 1756
- net/ipv4/fib_frontend.c, line 652
- net/ipv4/fib_rules.c, line 466
- net/ipv4/netfilter/ip_queue.c, line 647
- net/ipv4/netfilter/ipfwadm_core.c, line 1385
- net/ipv4/netfilter/ipt_MASQUERADE.c, line 190
- net/ipx/af_ipx.c, line 2562
- net/netrom/af_netrom.c, line 1311
- net/dechnet/af_dechnet.c, line 2260
- net/dechnet/dn_rules.c, line 363
- net/ipv6/ipv6_sockglue.c, line 563
- net/ipv6/netfilter/ip6_queue.c, line 703
- net/bridge/br.c, line 51
- net/econet/af_econet.c, line 1125
- net/x25/af_x25.c, line 1324
- net/rose/af_rose.c, line 1463
- net/wanrouter/af_wanpipe.c, line 2762
- net/packet/af_packet.c, line 1896
- net/irda/af_irda.c, line 2590
- net/atm/clip.c:
 - line 739
 - line 740
- net/atm/mpc.c, line 768
- net/8021q/vlan.c, line 99
- drivers/net/wan/lapbether.c, line 478
- drivers/net/hamradio/bpqether.c, line 614
- drivers/net/pppoe.c, line 1065
- drivers/net/bonding.c, line 2010

- `register_netdevice_notifier`
- `dst_dev_event()`
- `rtnetlink_dev_notifier()`
- `ip_netdev_notifier()`
- `ip_mr_notifier()`
- `fib_netdev_notifier()`
- `fib_rules_notifier()`
- `ipq_dev_notifier()`

Actions associated with NETDEV_REGISTER

net/core/dst.c, line 214 dst_dev_event()

No action is taken on REGISTER. On UNREGISTER/DOWN dst->output is set to BLACKHOLE.

net/core/rtnetlink.c, line 526 rtnetlink_event()

```
494     case NETDEV_REGISTER:
495         rtmmsg_info(RTM_NEWLINK, dev, ~0U);
496         break;
```

net/ipv4/devinet.c, line 1140 inetdev_event()

```
802     case NETDEV_REGISTER:
803         printk(KERN_DEBUG "inetdev_event: bug\n");
804         dev->ip_ptr = NULL;
805         break;
```

net/ipv4/ipmr.c, line 1756 ipmr_device_event()

Multicast routing support via mrouted.

net/ipv4/fib_frontend.c, line 652 fib_netdev_event()

No action.

net/ipv4/fib_rules.c, line 466 fib_rules_event()

```
388     else if (event == NETDEV_REGISTER)
389         fib_rules_attach(dev);
```

Recall that fib_rules aren't in play unless IP_MULTIPLE_TABLES is configured.

net/ipv4/netfilter/ip_queue.c, line 647 ipq_rcv_dev_event()

No action is taken on REGISTER. The packet queue is dumped on DOWN.

The functions `rtnetlink_event()` and `inetdev_event()` (among many others) are indirectly called. This whole gruesome collection is shown on the previous pages. Here we focus upon `rtnetlink_event()` when called with the parameter `NETDEV_REGISTER`.

```
487 static int rtnetlink_event(struct notifier_block *this,
                             unsigned long event, void *ptr)
488 {
489     struct net_device *dev = ptr;
490     switch (event) {
491     case NETDEV_UNREGISTER:
492         rtmsg_info(RTM_DELLINK, dev, ~0U);
493         break;
494     case NETDEV_REGISTER:
495         rtmsg_info(RTM_NEWLINK, dev, ~0U);
496         break;
497     case NETDEV_UP:
498     case NETDEV_DOWN:
499         rtmsg_info(RTM_NEWLINK, dev,
                    IFF_UP|IFF_RUNNING);
500         break;
501     case NETDEV_CHANGE:
502     case NETDEV_GOING_DOWN:
503         break;
504     default:
505         rtmsg_info(RTM_NEWLINK, dev, 0);
506         break;
507     }
508     return NOTIFY_DONE;
509 }
510
```

Construction of netlink messages

The `rtmsg_ifinfo()` function is called in response to register, unregister, interface up, and interface down events. For register and interface up the message type is `RTM_NEWLINK`. The change parameter is `0xffffffff`.

```
247 void rtmsg_ifinfo(int type, struct net_device *dev,
                    unsigned change)
248 {
249     struct sk_buff *skb;
250     int size = NLMSG_GOODSIZE;
251
252     skb = alloc_skb(size, GFP_KERNEL);
253     if (!skb)
254         return;
255
256     if (rtnetlink_fill_ifinfo(skb, dev,
                              type, 0, 0, change) < 0) {
257         kfree_skb(skb);
258         return;
259     }
```

The updating of the control buffer appears to be establishing the target recipients of this message. The call to `netlink_broadcast()` actually effects the delivery.

```
260     NETLINK_CB(skb).dst_groups = RTMGRP_LINK;
261     netlink_broadcast(rtnl, skb, 0, RTMGRP_LINK, GFP_KERNEL);
262 }
```


Netlink message headers

Each of these messages begins with a header of the following layout. In the present context the type is always RTM_NEWLINK.

```
26 struct nlmsg_hdr
27 {
28     __u32  nlmsg_len;          /* Len of msg including hdr */
29     __u16  nlmsg_type;        /* Message content */
30     __u16  nlmsg_flags;       /* Additional flags */
31     __u32  nlmsg_seq;         /* Sequence number */
32     __u32  nlmsg_pid;         /* Sending process PID */
33 };
```

For messages of the interface information class, a fixed structure follows the netlink header.

```
419 struct ifinfo_msg
420 {
421     unsigned char  ifi_family;
422     unsigned char  __ifi_pad;
423     unsigned short ifi_type;      /* ARPHRD_* */
424     int            ifi_index;     /* Link index */
425     unsigned       ifi_flags;     /* IFF_* flags */
426     unsigned       ifi_change;    /* IFF_* change mask */
427 };
```

```
152 static int rtnetlink_fill_ifinfo(struct sk_buff *skb,
153     struct net_device *dev,
154     int type, u32 pid, u32 seq, u32 change)
155 {
156     struct ifinfo_msg *r;
157     struct nlmsg_hdr *nlh;
158     unsigned char *b = skb->tail;
```

The NLMSG_PUT macro builds the header. In this context pid which plays an important role in routing of these messages is 0.

```
159     nlh = NLMSG_PUT(skb, pid, seq, type, sizeof(*r));
160     if (pid) nlh->nlmsg_flags |= NLM_F_MULTI;
```

Fill in the interface information header.

```
161     r = NLMSG_DATA(nlh);
162     r->ifi_family = AF_UNSPEC;
163     r->ifi_type = dev->type;
164     r->ifi_index = dev->ifindex;
165     r->ifi_flags = dev->flags;
166     r->ifi_change = change;
167
168     if (!netif_running(dev) || !netif_carrier_ok(dev))
169         r->ifi_flags &= ~IFF_RUNNING;
170     else
171         r->ifi_flags |= IFF_RUNNING;
172
173     RTA_PUT(skb, IFLA_IFNAME, strlen(dev->name)+1,
174             dev->name);
175     if (dev->addr_len) {
176         RTA_PUT(skb, IFLA_ADDRESS, dev->addr_len,
177               dev->dev_addr);
178         RTA_PUT(skb, IFLA_BROADCAST, dev->addr_len,
179               dev->broadcast);
180     }
181     if (1) {
182         unsigned mtu = dev->mtu;
183         RTA_PUT(skb, IFLA_MTU, sizeof(mtu), &mtu);
184     }
185     if (dev->ifindex != dev->iflink)
186         RTA_PUT(skb, IFLA_LINK, sizeof(int), &dev->iflink);
187     if (dev->qdisc_slipping)
188         RTA_PUT(skb, IFLA_QDISC,
189               strlen(dev->qdisc_slipping->ops->id) + 1,
190               dev->qdisc_slipping->ops->id);
191     if (dev->master)
192         RTA_PUT(skb, IFLA_MASTER, sizeof(int),
193               &dev->master->ifindex);
194     if (dev->get_stats) {
195         struct net_device_stats *stats =
196             dev->get_stats(dev);
197         if (stats)
198             RTA_PUT(skb, IFLA_STATS, sizeof(*stats), stats);
199     }
200     nlh->nmsg_len = skb->tail - b;
201     return skb->len;
202 }
203
204 nlmsg_failure:
205 rtattr_failure:
206     skb_trim(skb, b - skb->data);
207     return -1;
208 }
```

RTA_PUT is a macro used to add information elements to the message.

```
564 #define RTA_PUT(skb, attrtype, attrlen, data) \
565 ({ if (skb_tailroom(skb) < (int)RTA_SPACE(attrlen)) goto \
    rtattr_failure; \
566     __rta_fill(skb, attrtype, attrlen, data); })
```

It relies upon a collection of related macros...

```
64
65 #define RTA_ALIGNTO 4
66 #define RTA_ALIGN(len) ( ((len)+RTA_ALIGNTO-1) &
    ~(RTA_ALIGNTO-1) )
67 #define RTA_OK(rta, len) ((len) > 0 && (rta)->rta_len >=
    sizeof(struct rtattr) && \
68     (rta)->rta_len <= (len))
69 #define RTA_NEXT(rta, attrlen) ((attrlen) -=
    RTA_ALIGN((rta)->rta_len), \
70     (struct
    rtattr*)((char*)(rta) + RTA_ALIGN((rta)->rta_len)))
71 #define RTA_LENGTH(len) (RTA_ALIGN(sizeof(struct rtattr)) +
    (len))
72 #define RTA_SPACE(len) RTA_ALIGN(RTA_LENGTH(len))
73 #define RTA_DATA(rta) ((void*)((char*)(rta) +
    RTA_LENGTH(0)))
74 #define RTA_PAYLOAD(rta) ((int)((rta)->rta_len) -
    RTA_LENGTH(0))
75
```

The RTA_PUT macro invokes the __rta_fill function to add the data to the packet.

```
107 void __rta_fill(struct sk_buff *skb, int attrtype, int
    attrlen, const void *data)
108 {
109     struct rtattr *rta;
110     int size = RTA_LENGTH(attrlen);
111
```

Here is where the actual (T, L, V) data is stored in the message.

```
112     rta = (struct rtattr*)skb_put(skb, RTA_ALIGN(size));
113     rta->rta_type = attrtype;
114     rta->rta_len = size;
115     memcpy(RTA_DATA(rta), data, attrlen);
116 }
```

```

478 void netlink_broadcast(struct sock *ssk, struct sk_buff
479     *skb, u32 pid, u32 group, int allocation)
480 {
481     struct sock *sk;
482     struct sk_buff *skb2 = NULL;
483     int protocol = ssk->protocol;
484     int failure = 0;
485
486     /* While we sleep in clone, do not allow to change socket
487        list */
488     netlink_lock_table();
489
490     for (sk = nl_table[protocol]; sk; sk = sk->next) {
491         if (ssk == sk)
492             continue;
493
494         if (sk->protoinfo.af_netlink->pid == pid ||
495             !(sk->protoinfo.af_netlink->groups&group))
496             continue;
497
498         if (failure) {
499             netlink_overrun(sk);
500             continue;
501         }
502
503         sock_hold(sk);
504         if (skb2 == NULL) {
505             if (atomic_read(&skb->users) != 1) {
506                 skb2 = skb_clone(skb, allocation);
507             } else {
508                 skb2 = skb;
509                 atomic_inc(&skb->users);
510             }
511         }
512         if (skb2 == NULL) {
513             netlink_overrun(sk)
514     /* Clone failed. Notify ALL listeners. */
515             failure = 1;
516         } else if (netlink_broadcast_deliver(sk, skb2)) {
517             netlink_overrun(sk);
518         } else
519             skb2 = NULL;
520         sock_put(sk);
521     }
522
523     netlink_unlock_table();
524
525     if (skb2)
526         kfree_skb(skb2);
527     kfree_skb(skb);
528 }

```

```

458 static __inline__ int netlink_broadcast_deliver(struct sock
      *sk, struct sk_buff
459 {

```

This is not normally configured in kernels I build. Here is the description from make xconfig:

This option will be removed soon. Any programs that want to use character special nodes like /dev/tap0 or /dev/route (all with major number 36) need this option, and need to be rewritten soon to use the real netlink socket.

```

460 #ifdef NL_EMULATE_DEV
461     if (sk->protoinfo.af_netlink->handler) {
462         skb_orphan(skb);
463         sk->protoinfo.af_netlink->handler(sk->protocol, skb);
464         return 0;
465     } else
466 #endif
467     if (atomic_read(&sk->rmem_alloc) <= sk->rcvbuf &&
468         !test_bit(0, &sk->protoinfo.af_netlink->state)) {
469         skb_orphan(skb);
470         skb_set_owner_r(skb, sk);
471         skb_queue_tail(&sk->receive_queue, skb);
472         sk->data_ready(sk, skb->len);
473         return 0;
474     }
475     return -1;
476 }

```

This code in net/ipv4/fib_frontend.c contains the glue that eventually causes fib_magic to be called. The resulting call causes the local and main tables to be updated.

```
631 struct notifier_block fib_netaddr_notifier = {
632     notifier_call: fib_netaddr_event,
633 };
634
635 struct notifier_block fib_netdev_notifier = {
636     notifier_call: fib_netdev_event,
637 };
638
639 void __init ip_fib_init(void)
640 {
641     #ifdef CONFIG_PROC_FS
642         proc_net_create("route", 0, fib_get_procinfo);
643     #endif /* CONFIG_PROC_FS */
644
645     #ifdef CONFIG_IP_MULTIPLE_TABLES
646         local_table = fib_hash_init(RT_TABLE_LOCAL);
647         main_table = fib_hash_init(RT_TABLE_MAIN);
648     #else
649         fib_rules_init();
650     #endif
651
652     register_netdevice_notifier(&fib_netdev_notifier);
653     register_netaddr_notifier(&fib_netaddr_notifier);
654 }
```


For the purposes of FIB (routing table) management , two important notifiers reside in net/ipv4/fib_frontend.c. These notifiers were registered during IP initialization at boot time.

```

631 struct notifier_block fib_netaddr_notifier = {
632     notifier_call : fib_netaddr_event,
633 };
634
635 struct notifier_block fib_netdev_notifier = {
636     notifier_call : fib_netdev_event,
637 };
638
639 void __init ip_fib_init(void)
640 {
641     #ifdef CONFIG_PROC_FS
642         proc_net_create("route", 0, fib_get_procinfo);
643     #endif /* CONFIG_PROC_FS */
644
645     #ifndef CONFIG_IP_MULTIPLE_TABLES
646         local_table = fib_hash_init(RT_TABLE_LOCAL);
647         main_table = fib_hash_init(RT_TABLE_MAIN);
648     #else
649         fib_rules_init();
650     #endif
651
652     register_netdev_notifier(&fib_netdev_notifier);
653     register_netaddr_notifier(&fib_netaddr_notifier);
654 }
655
14 struct notifier_block
15 {
16     int (*notifier_call)(struct notifier_block *self, unsigned
17                          long, void *);
17     struct notifier_block *next;
18     int priority;
19 };
866 int register_netdev_notifier(struct notifier_block *nb)
867 {
868     return notifier_chain_register(&netdev_chain, nb);
869 }

```



```

575 static int fib_inetaddr_event(struct notifier_block *this,
576                               unsigned long event, v
577 {
578     struct in_ifaddr *ifa = (struct in_ifaddr*)ptr;
579     switch (event) {
580     case NETDEV_UP:
581         fib_add_ifaddr(ifa);
582         rt_cache_flush(-1);
583         break;
584     case NETDEV_DOWN:
585         fib_del_ifaddr(ifa);
586         if (ifa->ifa_dev && ifa->ifa_dev->ifa_list == NULL){
587             /* Last address was deleted from this interface.
588              * Disable IP.
589              */
590             fib_disable_ip(ifa->ifa_dev->dev, 1);
591         } else {
592             rt_cache_flush(-1);
593         }
594         break;
595     }
596     return NOTIFY_DONE;
597 }

459 static void fib_add_ifaddr(struct in_ifaddr *ifa)
460 {
461     struct in_device *in_dev = ifa->ifa_dev;
462     struct net_device *dev = in_dev->dev;
463     struct in_ifaddr *prim = ifa;
464     u32 mask = ifa->ifa_mask;
465     u32 addr = ifa->ifa_local;
466     u32 prefix = ifa->ifa_address&mask;
467
468     if (ifa->ifa_flags&IFA_F_SECONDARY) {
469         prim = inet_ifa_byprefix(in_dev, prefix, mask);
470         if (prim == NULL) {
471             printk(KERN_DEBUG "fib_add_ifaddr: bug: prim ==
472                             NULL\n");
473             return;
474         }
475     }
476     fib_magic(RTM_NEWROUTE, RTN_LOCAL, addr, 32, prim);
477
478     if (!(dev->flags&IFF_UP))
479         return;
480

```

```

481 /* Add broadcast address, if it is explicitly assigned. */
482   if (ifa->ifa_broadcast && ifa->ifa_broadcast !=
483       0xFFFFFFFF)
484       fib_magic(RTM_NEWROUTE, RTN_BROADCAST,
485                ifa->ifa_broadcast, 32, prim);
486   if (!ZERONET(prefix) &&
487       !(ifa->ifa_flags&IFA_F_SECONDARY) &&
488       (prefix != addr || ifa->ifa_prefixlen < 32)) {
489       fib_magic(RTM_NEWROUTE, dev->flags&IFF_LOOPBACK ?
490                RTN_LOCAL : RTN_UNICAST, prefix,
491                ifa->ifa_prefixlen, prim);
492   }
493   /* Add network specific broadcasts, when it takes a sense */
494   if (ifa->ifa_prefixlen < 31) {
495       fib_magic(RTM_NEWROUTE, RTN_BROADCAST, prefix,
496                32, prim);
497       fib_magic(RTM_NEWROUTE, RTN_BROADCAST,
498                prefix|~mask, 32, prim);
499   }
500 }

```

```

417 static void fib_magic(int cmd, int type, u32 dst, int
                               dst_len, struct in_ifaddr *ifa)
418 {
419     struct fib_table *tb;
420     struct {
421         struct nlmsg_hdr nlh;
422         struct rtm_rtm;
423     } req;
424     struct kern_rta rta;
425
426     memset(&req.rtm, 0, sizeof(req.rtm));
427     memset(&rta, 0, sizeof(rta));
428
429     if (type == RTN_UNICAST)
430         tb = fib_new_table(RT_TABLE_MAIN);
431     else
432         tb = fib_new_table(RT_TABLE_LOCAL);
433
434     if (tb == NULL)
435         return;
436
437     req.nlh.nlmsg_len = sizeof(req);
438     req.nlh.nlmsg_type = cmd;
439     req.nlh.nlmsg_flags =
440         NLM_F_REQUEST | NLM_F_CREATE | NLM_F_APPEND;
441     req.nlh.nlmsg_pid = 0;
442     req.nlh.nlmsg_seq = 0;
443
444     req.rtm.rtm_dst_len = dst_len;
445     req.rtm.rtm_table = tb->tb_id;
446     req.rtm.rtm_protocol = RTPROT_KERNEL;
447     req.rtm.rtm_scope = (type != RTN_LOCAL ? RT_SCOPE_LINK :
448         RT_SCOPE_HOST);
449     req.rtm.rtm_type = type;
450
451     rta.rta_dst = &dst;
452     rta.rta_prefsrc = &ifa->ifa_local;
453     rta.rta_oif = &ifa->ifa_dev->dev->index;
454
455     if (cmd == RTM_NEWROUTE)
456         tb->tb_insert(tb, &req.rtm, &rta, &req.nlh, NULL);
457     else
458         tb->tb_delete(tb, &req.rtm, &rta, &req.nlh, NULL);
459 }

```