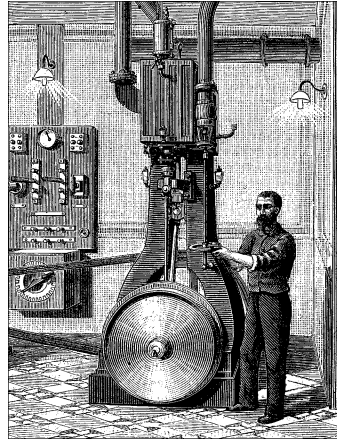# NETWORK DRIVERS

We are now through discussing char and block drivers and are ready to move on to the fascinating world of networking. Network interfaces are the third standard class of Linux devices, and this chapter describes how they interact with the rest of the kernel.

The role of a network interface within the system is similar to that of a mounted block device. A block device registers its features in the `blk_dev` array and other kernel structures, and it then "transmits" and "receives" blocks on request, by means of its *request* function. Similarly, a network interface must register itself in specific data structures in order to be invoked when packets are exchanged with the outside world.

There are a few important differences between mounted disks and packet-delivery interfaces. To begin with, a disk exists as a special file in the */dev* directory, whereas a network interface has no such entry point. The normal file operations (read, write, and so on) do not make sense when applied to network interfaces, so it is not possible to apply the Unix "everything is a file" approach to them. Thus, network interfaces exist in their own namespace and export a different set of operations.

Although you may object that applications use the *read* and *write* system calls when using sockets, those calls act on a software object that is distinct from the interface. Several hundred sockets can be multiplexed on the same physical interface.

But the most important difference between the two is that block drivers operate only in response to requests from the kernel, whereas network drivers receive packets asynchronously from the outside. Thus, while a block driver *is asked* to send a buffer toward the kernel, the network device *asks* to push incoming packets toward the kernel. The kernel interface for network drivers is designed for this different mode of operation.

Network drivers also have to be prepared to support a number of administrative tasks, such as setting addresses, modifying transmission parameters, and maintaining traffic and error statistics. The API for network drivers reflects this need, and thus looks somewhat different from the interfaces we have seen so far.

The network subsystem of the Linux kernel is designed to be completely protocol independent. This applies to both networking protocols (IP versus IPX or other protocols) and hardware protocols (Ethernet versus token ring, etc.). Interaction between a network driver and the kernel proper deals with one network packet at a time; this allows protocol issues to be hidden neatly from the driver and the physical transmission to be hidden from the protocol.

This chapter describes how the network interfaces fit in with the rest of the Linux kernel and shows a memory-based modularized network interface, which is called (you guessed it) *snull*. To simplify the discussion, the interface uses the Ethernet hardware protocol and transmits IP packets. The knowledge you acquire from examining *snull* can be readily applied to protocols other than IP, and writing a non-Ethernet driver is only different in tiny details related to the actual network protocol.

This chapter doesn't talk about IP numbering schemes, network protocols, or other general networking concepts. Such topics are not (usually) of concern to the driver writer, and it's impossible to offer a satisfactory overview of networking technology in less than a few hundred pages. The interested reader is urged to refer to other books describing networking issues.

The networking subsystem has seen many changes over the years as the kernel developers have striven to provide the best performance possible. The bulk of this chapter describes network drivers as they are implemented in the 2.4 kernel. Once again, the sample code works on the 2.0 and 2.2 kernels as well, and we cover the differences between those kernels and 2.4 at the end of the chapter.

One note on terminology is called for before getting into network devices. The networking world uses the term *octet* to refer to a group of eight bits, which is generally the smallest unit understood by networking devices and protocols. The term byte is almost never encountered in this context. In keeping with standard usage, we will use octet when talking about networking devices.

## *How snull Is Designed*

This section discusses the design concepts that led to the *snull* network interface. Although this information might appear to be of marginal use, failing to understand this driver might lead to problems while playing with the sample code.

The first, and most important, design decision was that the sample interfaces should remain independent of real hardware, just like most of the sample code

used in this book. This constraint led to something that resembles the loopback interface. *snull* is not a loopback interface, however; it simulates conversations with real remote hosts in order to better demonstrate the task of writing a network driver. The Linux loopback driver is actually quite simple; it can be found in *drivers/net/loopback.c.*

Another feature of *snull* is that it supports only IP traffic. This is a consequence of the internal workings of the interface—*snull* has to look inside and interpret the packets to properly emulate a pair of hardware interfaces. Real interfaces don't depend on the protocol being transmitted, and this limitation of *snull* doesn't affect the fragments of code that are shown in this chapter.

## Assigning IP Numbers

The *snull* module creates two interfaces. These interfaces are different from a simple loopback in that whatever you transmit through one of the interfaces loops back to the other one, not to itself. It looks like you have two external links, but actually your computer is replying to itself.

Unfortunately, this effect can't be accomplished through IP-number assignment alone, because the kernel wouldn't send out a packet through interface A that was directed to its own interface B. Instead, it would use the loopback channel without passing through *snull.* To be able to establish a communication through the *snull* interfaces, the source and destination addresses need to be modified during data transmission. In other words, packets sent through one of the interfaces should be received by the other, but the receiver of the outgoing packet shouldn't be recognized as the local host. The same applies to the source address of received packets.

To achieve this kind of "hidden loopback," the *snull* interface toggles the least significant bit of the third octet of both the source and destination addresses; that is, it changes both the network number and the host number of class C IP numbers. The net effect is that packets sent to network A (connected to `sn0`, the first interface) appear on the `sn1` interface as packets belonging to network B.

To avoid dealing with too many numbers, let's assign symbolic names to the IP numbers involved:

- `snullnet0` is the class C network that is connected to the `sn0` interface. Similarly, `snullnet1` is the network connected to `sn1`. The addresses of these networks should differ only in the least significant bit of the third octet.

- `local0` is the IP address assigned to the `sn0` interface; it belongs to `snullnet0`. The address associated with `sn1` is `local1`. `local0` and `local1` must differ in the least significant bit of their third octet and in the fourth octet.

- `remote0` is a host in `snullnet0`, and its fourth octet is the same as that of `local1`. Any packet sent to `remote0` will reach `local1` after its class C address has been modified by the interface code. The host `remote1` belongs to `snullnet1`, and its fourth octet is the same as that of `local0`.

The operation of the *snull* interfaces is depicted in Figure 14-1, in which the host-name associated with each interface is printed near the interface name.
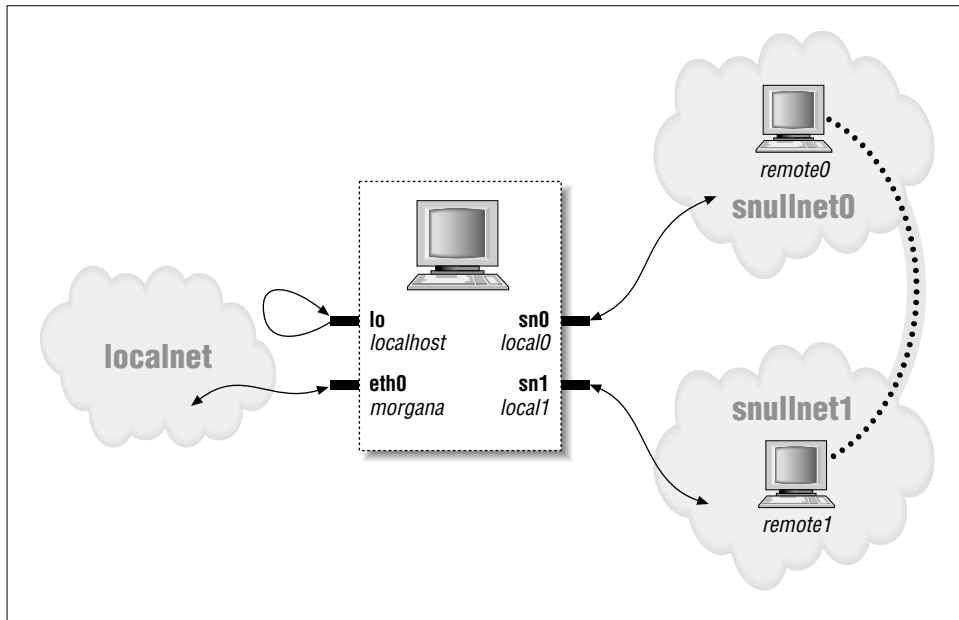


*Figure 14-1. How a host sees its interfaces*

Here are possible values for the network numbers. Once you put these lines in */etc/networks*, you can call your networks by name. The values shown were chosen from the range of numbers reserved for private use.

```
snullnet0       192.168.0.0
snullnet1       192.168.1.0
```

The following are possible host numbers to put into */etc/hosts*:

```
192.168.0.1    local0
192.168.0.2    remote0
192.168.1.2    local1
192.168.1.1    remote1
```

The important feature of these numbers is that the host portion of `local0` is the same as that of `remote1`, and the host portion of `local1` is the same as that of `remote0`. You can use completely different numbers as long as this relationship applies.

Be careful, however, if your computer is already connected to a network. The numbers you choose might be real Internet or intranet numbers, and assigning them to your interfaces will prevent communication with the real hosts. For example, although the numbers just shown are not routable Internet numbers, they could already be used by your private network if it lives behind a firewall.

Whatever numbers you choose, you can correctly set up the interfaces for operation by issuing the following commands:

```
ifconfig sn0 local0
ifconfig sn1 local1
case "`uname -r`" in 2.0.*)
    route add -net snullnet0 dev sn0
    route add -net snullnet1 dev sn1
esac
```

There is no need to invoke *route* with 2.2 and later kernels because the route is automatically added. Also, you may need to add the `netmask 255.255.255.0` parameter if the address range chosen is not a class C range.

At this point, the "remote" end of the interface can be reached. The following screendump shows how a host reaches `remote0` and `remote1` through the *snull* interface.

```
morgana% ping -c 2 remote0
64 bytes from 192.168.0.99: icmp_seq=0 ttl=64 time=1.6 ms
64 bytes from 192.168.0.99: icmp_seq=1 ttl=64 time=0.9 ms
2 packets transmitted, 2 packets received, 0% packet loss

morgana% ping -c 2 remote1
64 bytes from 192.168.1.88: icmp_seq=0 ttl=64 time=1.8 ms
64 bytes from 192.168.1.88: icmp_seq=1 ttl=64 time=0.9 ms
2 packets transmitted, 2 packets received, 0% packet loss
```

Note that you won't be able to reach any other "host" belonging to the two networks because the packets are discarded by your computer after the address has been modified and the packet has been received. For example, a packet aimed at 192.168.0.32 will leave through `sn0` and reappear at `sn1` with a destination address of 192.168.1.32, which is not a local address for the host computer.

## The Physical Transport of Packets

As far as data transport is concerned, the *snull* interfaces belong to the Ethernet class.

*snull* emulates Ethernet because the vast majority of existing networks—at least the segments that a workstation connects to—are based on Ethernet technology, be it 10baseT, 100baseT, or gigabit. Additionally, the kernel offers some

generalized support for Ethernet devices, and there's no reason not to use it. The advantage of being an Ethernet device is so strong that even the *plip* interface (the interface that uses the printer ports) declares itself as an Ethernet device.

The last advantage of using the Ethernet setup for *snull* is that you can run *tcp-dump* on the interface to see the packets go by. Watching the interfaces with *tcp-dump* can be a useful way to see how the two interfaces work. (Note that on 2.0 kernels, *tcpdump* will not work properly unless *snull*'s interfaces show up as `ethx`. Load the driver with the `eth=1` option to use the regular Ethernet names, rather than the default `snx` names.)

As was mentioned previously, *snull* only works with IP packets. This limitation is a result of the fact that *snull* snoops in the packets and even modifies them, in order for the code to work. The code modifies the source, destination, and check-sum in the IP header of each packet without checking whether it actually conveys IP information. This quick-and-dirty data modification destroys non-IP packets. If you want to deliver other protocols through *snull*, you must modify the module's source code.

# *Connecting to the Kernel*

We'll start looking at the structure of network drivers by dissecting the *snull* source. Keeping the source code for several drivers handy might help you follow the discussion and to see how real-world Linux network drivers operate. As a place to start, we suggest *loopback.c*, *plip.c*, and *3c509.c*, in order of increasing complexity. Keeping *skeleton.c* handy might help as well, although this sample driver doesn't actually run. All these files live in *drivers/net*, within the kernel source tree.

## *Module Loading*

When a driver module is loaded into a running kernel, it requests resources and offers facilities; there's nothing new in that. And there's also nothing new in the way resources are requested. The driver should probe for its device and its hard-ware location (I/O ports and IRQ line)—but without registering them—as described in "Installing an Interrupt Handler" in Chapter 9. The way a network driver is registered by its module initialization function is different from char and block drivers. Since there is no equivalent of major and minor numbers for net-work interfaces, a network driver does not request such a number. Instead, the driver inserts a data structure for each newly detected interface into a global list of network devices.

Each interface is described by a `struct net_device` item. The structures for `sn0` and `sn1`, the two *snull* interfaces, are declared like this:

```
struct net_device snull_devs[2] = {
    { init: snull_init, },  /* init, nothing more */
    { init: snull_init, }
};
```

The initialization shown seems quite simple—it sets only one field. In fact, the `net_device` structure is huge, and we will be filling in other pieces of it later on. But it is not helpful to cover the entire structure at this point; instead, we will explain each field as it is used. For the interested reader, the definition of the structure may be found in `<linux/netdevice.h>`.

The first `struct net_device` field we will look at is `name`, which holds the interface name (the string identifying the interface). The driver can hardwire a name for the interface or it can allow dynamic assignment, which works like this: if the name contains a `%d` format string, the first available name found by replacing that string with a small integer is used. Thus, `eth%d` is turned into the first available `eth`*n* name; the first Ethernet interface is called `eth0`, and the others follow in numeric order. The *snull* interfaces are called `sn0` and `sn1` by default. However, if `eth=1` is specified at load time (causing the integer variable `snull_eth` to be set to 1), *snull_init* uses dynamic assignment, as follows:

```
if (!snull_eth) { /* call them "sn0" and "sn1" */
    strcpy(snull_devs[0].name, "sn0");
    strcpy(snull_devs[1].name, "sn1");
} else { /* use automatic assignment */
    strcpy(snull_devs[0].name, "eth%d");
    strcpy(snull_devs[1].name, "eth%d");
}
```

The other field we initialized is `init`, a function pointer. Whenever you register a device, the kernel asks the driver to initialize itself. Initialization means probing for the physical interface and filling the `net_device` structure with the proper values, as described in the following section. If initialization fails, the structure is not linked to the global list of network devices. This peculiar way of setting things up is most useful during system boot; every driver tries to register its own devices, but only devices that exist are linked to the list.

Because the real initialization is performed elsewhere, the initialization function has little to do, and a single statement does it:

```
for (i=0; i<2;  i++)
    if ( (result = register_netdev(snull_devs + i)) )
        printk("snull: error %i registering device \"%s\"\n",
                result, snull_devs[i].name);
    else device_present++;
```

## *Initializing Each Device*

Probing for the device should be performed in the *init* function for the interface (which is often called the "probe" function). The single argument received by *init* is a pointer to the device being initialized; its return value is either 0 or a negative error code, usually -ENODEV.

No real probing is performed for the *snull* interface, because it is not bound to any hardware. When you write a real driver for a real interface, the usual rules for probing devices apply, depending on the peripheral bus you are using. Also, you should avoid registering I/O ports and interrupt lines at this point. Hardware registration should be delayed until device open time; this is particularly important if interrupt lines are shared with other devices. You don't want your interface to be called every time another device triggers an IRQ line just to reply "no, it's not mine."

The main role of the initialization routine is to fill in the `dev` structure for this device. Note that for network devices, this structure is always put together at runtime. Because of the way the network interface probing works, the `dev` structure cannot be set up at compile time in the same manner as a `file_operations` or `block_device_operations` structure. So, on exit from `dev->init`, the `dev` structure should be filled with correct values. Fortunately, the kernel takes care of some Ethernet-wide defaults through the function *ether_setup*, which fills several fields in `struct net_device`.

The core of *snull_init* is as follows:

```
ether_setup(dev); /* assign some of the fields */

dev->open            = snull_open;
dev->stop            = snull_release;
dev->set_config      = snull_config;
dev->hard_start_xmit = snull_tx;
dev->do_ioctl        = snull_ioctl;
dev->get_stats       = snull_stats;
dev->rebuild_header  = snull_rebuild_header;
dev->hard_header     = snull_header;
#ifdef HAVE_TX_TIMEOUT
dev->tx_timeout      = snull_tx_timeout;
dev->watchdog_timeo = timeout;
#endif
/* keep the default flags, just add NOARP */
dev->flags           |= IFF_NOARP;
dev->hard_header_cache = NULL;      /* Disable caching */
SET_MODULE_OWNER(dev);
```

The single unusual feature of the code is setting `IFF_NOARP` in the flags. This specifies that the interface cannot use ARP, the Address Resolution Protocol. ARP is

a low-level Ethernet protocol; its job is to turn IP addresses into Ethernet Medium Access Control (MAC) addresses. Since the "remote" systems simulated by *snull* do not really exist, there is nobody available to answer ARP requests for them. Rather than complicate *snull* with the addition of an ARP implementation, we chose to mark the interface as being unable to handle that protocol. The assignment to `hard_header_cache` is there for a similar reason: it disables the caching of the (nonexistent) ARP replies on this interface. This topic is discussed in detail later in this chapter in "MAC Address Resolution."

The initialization code also sets a couple of fields (`tx_timeout` and `watchdog_timeo`) that relate to the handling of transmission timeouts. We will cover this topic thoroughly later in this chapter in "Transmission Timeouts."

Finally, this code calls `SET_MODULE_OWNER`, which initializes the `owner` field of the `net_device` structure with a pointer to the module itself. The kernel uses this information in exactly the same way it uses the `owner` field of the `file_operations` structure—to maintain the module's usage count.

We'll look now at one more `struct net_device` field, `priv`. Its role is similar to that of the `private_data` pointer that we used for char drivers. Unlike `fops->private_data`, this `priv` pointer is allocated at initialization time instead of open time, because the data item pointed to by `priv` usually includes the statistical information about interface activity. It's important that statistical information always be available, even when the interface is down, because users may want to display the statistics at any time by calling *ifconfig*. The memory wasted by allocating `priv` during initialization instead of on open is irrelevant because most probed interfaces are constantly up and running in the system. The *snull* module declares a `snull_priv` data structure to be used for `priv`:

```
struct snull_priv {
    struct net_device_stats stats;
    int status;
    int rx_packetlen;
    u8 *rx_packetdata;
    int tx_packetlen;
    u8 *tx_packetdata;
    struct sk_buff *skb;
    spinlock_t lock;
};
```

The structure includes an instance of `struct net_device_stats`, which is the standard place to hold interface statistics. The following lines in *snull_init* allocate and initialize `dev->priv`:

```
dev->priv = kmalloc(sizeof(struct snull_priv), GFP_KERNEL);
if (dev->priv == NULL)
    return -ENOMEM;
memset(dev->priv, 0, sizeof(struct snull_priv));
spin_lock_init(& ((struct snull_priv *) dev->priv)->lock);
```

## *Module Unloading*

Nothing special happens when the module is unloaded. The module cleanup function simply unregisters the interfaces from the list after releasing memory associated with the private structure:

```
void snull_cleanup(void)
{
    int i;

    for (i=0; i<2;  i++) {
        kfree(snull_devs[i].priv);
        unregister_netdev(snull_devs + i);
    }
    return;
}
```

## *Modularized and Nonmodularized Drivers*

Although char and block drivers are the same regardless of whether they're modular or linked into the kernel, that's not the case for network drivers.

When a driver is linked directly into the Linux kernel, it doesn't declare its own `net_device` structures; the structures declared in *drivers/net/Space.c* are used instead. *Space.c* declares a linked list of all the network devices, both driver-specific structures like `plip1` and general-purpose `eth` devices. Ethernet drivers don't care about their `net_device` structures at all, because they use the general-purpose structures. Such general `eth` device structures declare *ethif_probe* as their *init* function. A programmer inserting a new Ethernet interface in the mainstream kernel needs only to add a call to the driver's initialization function to *ethif_probe*. Authors of non-`eth` drivers, on the other hand, insert their `net_device` structures in *Space.c*. In both cases only the source file *Space.c* has to be modified if the driver must be linked to the kernel proper.

At system boot, the network initialization code loops through all the `net_device` structures and calls their probing (`dev->init`) functions by passing them a pointer to the device itself. If the probe function succeeds, the kernel initializes the next available `net_device` structure to use that interface. This way of setting up drivers permits incremental assignment of devices to the names `eth0`, `eth1`, and so on, without changing the `name` field of each device.

When a modularized driver is loaded, on the other hand, it declares its own `net_device` structures (as we have seen in this chapter), even if the interface it controls is an Ethernet interface.

The curious reader can learn more about interface initialization by looking at *Space.c* and *net_init.c*.

# *The net_device Structure in Detail*

The `net_device` structure is at the very core of the network driver layer and deserves a complete description. At a first reading, however, you can skip this section, because you don't need a thorough understanding of the structure to get started. This list describes all the fields, but more to provide a reference than to be memorized. The rest of this chapter briefly describes each field as soon as it is used in the sample code, so you don't need to keep referring back to this section.

`struct net_device` can be conceptually divided into two parts: visible and invisible. The visible part of the structure is made up of the fields that can be explicitly assigned in static `net_device` structures. All structures in *drivers/net/Space.c* are initialized in this way, without using the tagged syntax for structure initialization. The remaining fields are used internally by the network code and usually are not initialized at compilation time, not even by tagged initialization. Some of the fields are accessed by drivers (for example, the ones that are assigned at initialization time), while some shouldn't be touched.

## *The Visible Head*

The first part of `struct net_device` is composed of the following fields, in this order:

`char name[IFNAMSIZ];`
  The name of the device. If the name contains a `%d` format string, the first available device name with the given base is used; assigned numbers start at zero.

`unsigned long rmem_end;`
`unsigned long rmem_start;`
`unsigned long mem_end;`
`unsigned long mem_start;`
  Device memory information. These fields hold the beginning and ending addresses of the shared memory used by the device. If the device has different receive and transmit memories, the `mem` fields are used for transmit memory and the `rmem` fields for receive memory. `mem_start` and `mem_end` can be specified on the kernel command line at system boot, and their values are retrieved by *ifconfig.* The `rmem` fields are never referenced outside of the driver itself. By convention, the `end` fields are set so that `end - start` is the amount of available on-board memory.

`unsigned long base_addr;`
  The I/O base address of the network interface. This field, like the previous ones, is assigned during device probe. The *ifconfig* command can be used to display or modify the current value. The `base_addr` can be explicitly assigned on the kernel command line at system boot or at load time. The field is not used by the kernel, like the memory fields shown previously.

```
unsigned char irq;
```
The assigned interrupt number. The value of `dev->irq` is printed by *ifconfig* when interfaces are listed. This value can usually be set at boot or load time and modified later using *ifconfig*.

```
unsigned char if_port;
```
Which port is in use on multiport devices. This field is used, for example, with devices that support both coaxial (`IF_PORT_10BASE2`) and twisted-pair (`IF_PORT_10BASET`) Ethernet connections. The full set of known port types is defined in `<linux/netdevice.h>`.

```
unsigned char dma;
```
The DMA channel allocated by the device. The field makes sense only with some peripheral buses, like ISA. It is not used outside of the device driver itself, but for informational purposes (in *ifconfig*).

```
unsigned long state;
```
Device state. The field includes several flags. Drivers do not normally manipulate these flags directly; instead, a set of utility functions has been provided. These functions will be discussed shortly when we get into driver operations.

```
struct net_device *next;
```
Pointer to the next device in the global linked list. This field shouldn't be touched by the driver.

```
int (*init)(struct net_device *dev);
```
The initialization function, described earlier.

## *The Hidden Fields*

The `net_device` structure includes many additional fields, which are usually assigned at device initialization. Some of these fields convey information about the interface, while some exist only for the benefit of the driver (i.e., they are not used by the kernel); other fields, most notably the device methods, are part of the kernel-driver interface.

We will list the three groups separately, independent of the actual order of the fields, which is not significant.

### *Interface information*

Most of the information about the interface is correctly set up by the function *ether_setup*. Ethernet cards can rely on this general-purpose function for most of these fields, but the `flags` and `dev_addr` fields are device specific and must be explicitly assigned at initialization time.

Some non-Ethernet interfaces can use helper functions similar to *ether_setup*. *drivers/net/net_init.c* exports a number of such functions, including the following:

```
void ltalk_setup(struct net_device *dev);
```
Sets up the fields for a LocalTalk device.

```
void fc_setup(struct net_device *dev);
```
Initializes for fiber channel devices.

```
void fddi_setup(struct net_device *dev);
```
Configures an interface for a Fiber Distributed Data Interface (FDDI) network.

```
void hippi_setup(struct net_device *dev);
```
Prepares fields for a High-Performance Parallel Interface (HIPPI) high-speed interconnect driver.

```
void tr_configure(struct net_device *dev);
```
Handles setup for token ring network interfaces. Note that the 2.4 kernel also exports a function *tr_setup*, which, interestingly, does nothing at all.

Most devices will be covered by one of these classes. If yours is something radically new and different, however, you will need to assign the following fields by hand.

```
unsigned short hard_header_len;
```
The hardware header length, that is, the number of octets that lead the transmitted packet before the IP header, or other protocol information. The value of `hard_header_len` is 14 (`ETH_HLEN`) for Ethernet interfaces.

```
unsigned mtu;
```
The maximum transfer unit (MTU). This field is used by the network layer to drive packet transmission. Ethernet has an MTU of 1500 octets (`ETH_DATA_LEN`).

```
unsigned long tx_queue_len;
```
The maximum number of frames that can be queued on the device's transmission queue. This value is set to 100 by *ether_setup*, but you can change it. For example, *plip* uses 10 to avoid wasting system memory (*plip* has a lower throughput than a real Ethernet interface).

```
unsigned short type;
```
The hardware type of the interface. The `type` field is used by ARP to determine what kind of hardware address the interface supports. The proper value for Ethernet interfaces is `ARPHRD_ETHER`, and that is the value set by *ether_setup*. The recognized types are defined in `<linux/if_arp.h>`.

```
unsigned char addr_len;
unsigned char broadcast[MAX_ADDR_LEN];
unsigned char dev_addr[MAX_ADDR_LEN];
```
Hardware (MAC) address length and device hardware addresses. The Ethernet address length is six octets (we are referring to the hardware ID of the

*437*

interface board), and the broadcast address is made up of six `0xff` octets; *ether_setup* arranges for these values to be correct. The device address, on the other hand, must be read from the interface board in a device-specific way, and the driver should copy it to `dev_addr`. The hardware address is used to generate correct Ethernet headers before the packet is handed over to the driver for transmission. The *snull* device doesn't use a physical interface, and it invents its own hardware address.

`unsigned short flags;`
> Interface flags, detailed next.

The `flags` field is a bit mask including the following bit values. The `IFF_` prefix stands for "interface flags." Some flags are managed by the kernel, and some are set by the interface at initialization time to assert various capabilities and other features of the interface. The valid flags, which are defined in `<linux/if.h>`, are as follows:

`IFF_UP`
> This flag is read-only for the driver. The kernel turns it on when the interface is active and ready to transfer packets.

`IFF_BROADCAST`
> This flag states that the interface allows broadcasting. Ethernet boards do.

`IFF_DEBUG`
> This marks debug mode. The flag can be used to control the verbosity of your *printk* calls or for other debugging purposes. Although no official driver currently uses this flag, it can be set and reset by user programs via *ioctl*, and your driver can use it. The *misc-progs/netifdebug* program can be used to turn the flag on and off.

`IFF_LOOPBACK`
> This flag should be set only in the loopback interface. The kernel checks for `IFF_LOOPBACK` instead of hardwiring the `lo` name as a special interface.

`IFF_POINTOPOINT`
> This flag signals that the interface is connected to a point-to-point link. It is set by *ifconfig*. For example, *plip* and the PPP driver have it set.

`IFF_NOARP`
> This means that the interface can't perform ARP. For example, point-to-point interfaces don't need to run ARP, which would only impose additional traffic without retrieving useful information. *snull* runs without ARP capabilities, so it sets the flag.

IFF_PROMISC
> This flag is set to activate promiscuous operation. By default, Ethernet inter-faces use a hardware filter to ensure that they receive broadcast packets and packets directed to that interface's hardware address only. Packet sniffers such as *tcpdump* set promiscuous mode on the interface in order to retrieve all packets that travel on the interface's transmission medium.

IFF_MULTICAST
> This flag is set by interfaces that are capable of multicast transmission. *ether_setup* sets IFF_MULTICAST by default, so if your driver does not sup-port multicast, it must clear the flag at initialization time.

IFF_ALLMULTI
> This flag tells the interface to receive all multicast packets. The kernel sets it when the host performs multicast routing, only if IFF_MULTICAST is set. IFF_ALLMULTI is read-only for the interface. We'll see the multicast flags used in "Multicasting" later in this chapter.

IFF_MASTER
IFF_SLAVE
> These flags are used by the load equalization code. The interface driver doesn't need to know about them.

IFF_PORTSEL
IFF_AUTOMEDIA
> These flags signal that the device is capable of switching between multiple media types, for example, unshielded twisted pair (UTP) versus coaxial Ether-net cables. If IFF_AUTOMEDIA is set, the device selects the proper medium automatically.

IFF_DYNAMIC
> This flag indicates that the address of this interface can change; used with dialup devices.

IFF_RUNNING
> This flag indicates that the interface is up and running. It is mostly present for BSD compatibility; the kernel makes little use of it. Most network drivers need not worry about IFF_RUNNING.

IFF_NOTRAILERS
> This flag is unused in Linux, but it exists for BSD compatibility.

When a program changes IFF_UP, the *open* or *stop* device method is called. When IFF_UP or any other flag is modified, the *set_multicast_list* method is invoked. If the driver needs to perform some action because of a modification in the flags, it must take that action in *set_multicast_list*. For example, when IFF_PROMISC is set or reset, *set_multicast_list* must notify the onboard hardware filter. The responsibilities of this device method are outlined in "Multicasting."

### The device methods

As happens with the char and block drivers, each network device declares the functions that act on it. Operations that can be performed on network interfaces are listed in this section. Some of the operations can be left `NULL`, and some are usually untouched because *ether_setup* assigns suitable methods to them.

Device methods for a network interface can be divided into two groups: fundamental and optional. Fundamental methods include those that are needed to be able to use the interface; optional methods implement more advanced functionalities that are not strictly required. The following are the fundamental methods:

`int (*open)(struct net_device *dev);`
> Opens the interface. The interface is opened whenever *ifconfig* activates it. The *open* method should register any system resource it needs (I/O ports, IRQ, DMA, etc.), turn on the hardware, and increment the module usage count.

`int (*stop)(struct net_device *dev);`
> Stops the interface. The interface is stopped when it is brought down; operations performed at open time should be reversed.

`int (*hard_start_xmit) (struct sk_buff *skb, struct`
`    net_device *dev);`
> This method initiates the transmission of a packet. The full packet (protocol headers and all) is contained in a socket buffer (`sk_buff`) structure. Socket buffers are introduced later in this chapter.

`int (*hard_header) (struct sk_buff *skb, struct net_device`
`    *dev, unsigned short type, void *daddr, void *saddr,`
`    unsigned len);`
> This function builds the hardware header from the source and destination hardware addresses that were previously retrieved; its job is to organize the information passed to it as arguments into an appropriate, device-specific hardware header. *eth_header* is the default function for Ethernet-like interfaces, and *ether_setup* assigns this field accordingly.

`int (*rebuild_header)(struct sk_buff *skb);`
> This function is used to rebuild the hardware header before a packet is transmitted. The default function used by Ethernet devices uses ARP to fill the packet with missing information. The *rebuild_header* method is used rarely in the 2.4 kernel; *hard_header* is used instead.

`void (*tx_timeout)(struct net_device *dev);`
> This method is called when a packet transmission fails to complete within a reasonable period, on the assumption that an interrupt has been missed or the interface has locked up. It should handle the problem and resume packet transmission.

```
struct net_device_stats *(*get_stats)(struct net_device
    *dev);
```
Whenever an application needs to get statistics for the interface, this method is called. This happens, for example, when *ifconfig* or *netstat -i* is run. A sample implementation for *snull* is introduced in "Statistical Information" later in this chapter.

```
int (*set_config)(struct net_device *dev, struct ifmap
    *map);
```
Changes the interface configuration. This method is the entry point for configuring the driver. The I/O address for the device and its interrupt number can be changed at runtime using *set_config*. This capability can be used by the system administrator if the interface cannot be probed for. Drivers for modern hardware normally do not need to implement this method.

The remaining device operations may be considered optional.

```
int (*do_ioctl)(struct net_device *dev, struct ifreq *ifr,
    int cmd);
```
Perform interface-specific *ioctl* commands. Implementation of those commands is described later in "Custom ioctl Commands." The corresponding field in `struct net_device` can be left as `NULL` if the interface doesn't need any interface-specific commands.

```
void (*set_multicast_list)(struct net_device *dev);
```
This method is called when the multicast list for the device changes and when the flags change. See "Multicasting" for further details and a sample implementation.

```
int (*set_mac_address)(struct net_device *dev, void *addr);
```
This function can be implemented if the interface supports the ability to change its hardware address. Many interfaces don't support this ability at all. Others use the default *eth_mac_addr* implementation (from *drivers/net/net_init.c*). *eth_mac_addr* only copies the new address into `dev->dev_addr`, and it will only do so if the interface is not running. Drivers that use *eth_mac_addr* should set the hardware MAC address from `dev->dev_addr` when they are configured.

```
int (*change_mtu)(struct net_device *dev, int new_mtu);
```
This function is in charge of taking action if there is a change in the MTU (maximum transfer unit) for the interface. If the driver needs to do anything particular when the MTU is changed, it should declare its own function; otherwise, the default will do the right thing. *snull* has a template for the function if you are interested.

```
int (*header_cache) (struct neighbour *neigh, struct
     hh_cache *hh);
```
*header_cache* is called to fill in the `hh_cache` structure with the results of an ARP query. Almost all drivers can use the default *eth_header_cache* implementation.

```
int (*header_cache_update) (struct hh_cache *hh, struct
     net_device *dev, unsigned char *haddr);
```
This method updates the destination address in the `hh_cache` structure in response to a change. Ethernet devices use *eth_header_cache_update*.

```
int (*hard_header_parse) (struct sk_buff *skb, unsigned char
     *haddr);
```
The *hard_header_parse* method extracts the source address from the packet contained in `skb`, copying it into the buffer at `haddr`. The return value from the function is the length of that address. Ethernet devices normally use *eth_header_parse*.

### Utility fields

The remaining `struct net_device` data fields are used by the interface to hold useful status information. Some of the fields are used by *ifconfig* and *netstat* to provide the user with information about the current configuration. An interface should thus assign values to these fields.

```
unsigned long trans_start;
unsigned long last_rx;
```
Both of these fields are meant to hold a jiffies value. The driver is responsible for updating these values when transmission begins and when a packet is received, respectively. The `trans_start` value is used by the networking subsystem to detect transmitter lockups. `last_rx` is currently unused, but the driver should maintain this field anyway to be prepared for future use.

```
int watchdog_timeo;
```
The minimum time (in jiffies) that should pass before the networking layer decides that a transmission timeout has occurred and calls the driver's *tx_timeout* function.

```
void *priv;
```
The equivalent of `filp->private_data`. The driver owns this pointer and can use it at will. Usually the private data structure includes a `struct net_device_stats` item. The field is used in "Initializing Each Device," later in this chapter.

```
struct dev_mc_list *mc_list;
int mc_count;
```
> These two fields are used in handling multicast transmission. `mc_count` is the count of items in `mc_list`. See "Multicasting" for further details.

```
spinlock_t xmit_lock;
int xmit_lock_owner;
```
> The `xmit_lock` is used to avoid multiple simultaneous calls to the driver's *hard_start_xmit* function. `xmit_lock_owner` is the number of the CPU that has obtained `xmit_lock`. The driver should make no changes to these fields.

```
struct module *owner;
```
> The module that "owns" this device structure; it is used to maintain the use count for the module.

There are other fields in `struct net_device`, but they are not used by network drivers.

## *Opening and Closing*

Our driver can probe for the interface at module load time or at kernel boot. Before the interface can carry packets, however, the kernel must open it and assign an address to it. The kernel will open or close an interface in response to the *ifconfig* command.

When *ifconfig* is used to assign an address to the interface, it performs two tasks. First, it assigns the address by means of `ioctl(SIOCSIFADDR)` (Socket I/O Control Set Interface Address). Then it sets the `IFF_UP` bit in `dev->flag` by means of `ioctl(SIOCSIFFLAGS)` (Socket I/O Control Set Interface Flags) to turn the interface on.

As far as the device is concerned, `ioctl(SIOCSIFADDR)` does nothing. No driver function is invoked—the task is device independent, and the kernel performs it. The latter command (`ioctl(SIOCSIFFLAGS)`), though, calls the *open* method for the device.

Similarly, when the interface is shut down, *ifconfig* uses `ioctl(SIOCSIFFLAGS)` to clear `IFF_UP`, and the *stop* method is called.

Both device methods return 0 in case of success and the usual negative value in case of error.

As far as the actual code is concerned, the driver has to perform many of the same tasks as the char and block drivers do. *open* requests any system resources it needs and tells the interface to come up; *stop* shuts down the interface and releases system resources. There are a couple of additional steps to be performed, however.

First, the hardware address needs to be copied from the hardware device to `dev->dev_addr` before the interface can communicate with the outside world. The hardware address can be assigned at probe time or at open time, at the driver's will. The *snull* software interface assigns it from within *open*; it just fakes a hardware number using an ASCII string of length `ETH_ALEN`, the length of Ethernet hardware addresses.

The *open* method should also start the interface's transmit queue (allow it to accept packets for transmission) once it is ready to start sending data. The kernel provides a function to start the queue:

```
void netif_start_queue(struct net_device *dev);
```

The *open* code for *snull* looks like the following:

```
int snull_open(struct net_device *dev)
{
    MOD_INC_USE_COUNT;

    /* request_region(), request_irq(), .... (like fops->open) */

    /*
     * Assign the hardware address of the board: use "\0SNULx", where
     * x is 0 or 1. The first byte is '\0' to avoid being a multicast
     * address (the first byte of multicast addrs is odd).
     */
    memcpy(dev->dev_addr, "\0SNUL0", ETH_ALEN);
    dev->dev_addr[ETH_ALEN-1] += (dev - snull_devs); /* the number */

    netif_start_queue(dev);
    return 0;
}
```

As you can see, in the absence of real hardware, there is little to do in the *open* method. The same is true of the *stop* method; it just reverses the operations of *open*. For this reason the function implementing *stop* is often called *close* or *release*.

```
int snull_release(struct net_device *dev)
{
    /* release ports, irq and such -- like fops->close */

    netif_stop_queue(dev); /* can't transmit any more */
    MOD_DEC_USE_COUNT;
    return 0;
}
```

The function:

```
void netif_stop_queue(struct net_device *dev);
```

is the opposite of *netif_start_queue*; it marks the device as being unable to transmit any more packets. The function must be called when the interface is closed (in the *stop* method) but can also be used to temporarily stop transmission, as explained in the next section.

# Packet Transmission

The most important tasks performed by network interfaces are data transmission and reception. We'll start with transmission because it is slightly easier to understand.

Whenever the kernel needs to transmit a data packet, it calls the *hard_start_transmit* method to put the data on an outgoing queue. Each packet handled by the kernel is contained in a socket buffer structure (`struct sk_buff`), whose definition is found in `<linux/skbuff.h>`. The structure gets its name from the Unix abstraction used to represent a network connection, the *socket*. Even if the interface has nothing to do with sockets, each network packet belongs to a socket in the higher network layers, and the input/output buffers of any socket are lists of `struct sk_buff` structures. The same `sk_buff` structure is used to host network data throughout all the Linux network subsystems, but a socket buffer is just a packet as far as the interface is concerned.

A pointer to `sk_buff` is usually called `skb`, and we follow this practice both in the sample code and in the text.

The socket buffer is a complex structure, and the kernel offers a number of functions to act on it. The functions are described later in "The Socket Buffers;" for now a few basic facts about `sk_buff` are enough for us to write a working driver.

The socket buffer passed to *hard_start_xmit* contains the physical packet as it should appear on the media, complete with the transmission-level headers. The interface doesn't need to modify the data being transmitted. `skb->data` points to the packet being transmitted, and `skb->len` is its length, in octets.

The *snull* packet transmission code is follows; the physical transmission machinery has been isolated in another function because every interface driver must implement it according to the specific hardware being driven.

```
int snull_tx(struct sk_buff *skb, struct net_device *dev)
{
    int len;
    char *data;
    struct snull_priv *priv = (struct snull_priv *) dev->priv;
    len = skb->len < ETH_ZLEN ? ETH_ZLEN : skb->len;
    data = skb->data;
    dev->trans_start = jiffies; /* save the timestamp */

    /* Remember the skb, so we can free it at interrupt time */
    priv->skb = skb;
```

```
      /* actual delivery of data is device specific, and not shown here */
      snull_hw_tx(data, len, dev);

      return 0; /* Our simple device cannot fail */
  }
```

The transmission function thus performs only some sanity checks on the packet and transmits the data through the hardware-related function. That function (*snull_hw_tx*) is omitted here since it is entirely occupied with implementing the trickery of the *snull* device (including manipulating the source and destination addresses) and has little of interest to authors of real network drivers. It is present, of course, in the sample source for those who want to go in and see how it works.

## Controlling Transmission Concurrency

The *hard_start_xmit* function is protected from concurrent calls by a spinlock (xmit_lock) in the `net_device` structure. As soon as the function returns, however, it may be called again. The function returns when the software is done instructing the hardware about packet transmission, but hardware transmission will likely not have been completed. This is not an issue with *snull*, which does all of its work using the CPU, so packet transmission is complete before the transmission function returns.

Real hardware interfaces, on the other hand, transmit packets asynchronously and have a limited amount of memory available to store outgoing packets. When that memory is exhausted (which, for some hardware, will happen with a single outstanding packet to transmit), the driver will need to tell the networking system not to start any more transmissions until the hardware is ready to accept new data.

This notification is accomplished by calling *netif_stop_queue*, the function introduced earlier to stop the queue. Once your driver has stopped its queue, it *must* arrange to restart the queue at some point in the future, when it is again able to accept packets for transmission. To do so, it should call:

```
  void netif_wake_queue(struct net_device *dev);
```

This function is just like *netif_start_queue*, except that it also pokes the networking system to make it start transmitting packets again.

Most modern network interfaces maintain an internal queue with multiple packets to transmit; in this way they can get the best performance from the network. Network drivers for these devices support having multiple transmisions outstanding at any given time, but device memory can fill up whether or not the hardware supports multiple outstanding transmission. Whenever device memory fills to the point that there is no room for the largest possible packet, the driver should stop the queue until space becomes available again.

## *Transmission Timeouts*

Most drivers that deal with real hardware have to be prepared for that hardware to fail to respond occasionally. Interfaces can forget what they are doing, or the system can lose an interrupt. This sort of problem is common with some devices designed to run on personal computers.

Many drivers handle this problem by setting timers; if the operation has not completed by the time the timer expires, something is wrong. The network system, as it happens, is essentially a complicated assembly of state machines controlled by a mass of timers. As such, the networking code is in a good position to detect transmission timeouts automatically.

Thus, network drivers need not worry about detecting such problems themselves. Instead, they need only set a timeout period, which goes in the `watch-dog_timeo` field of the `net_device` structure. This period, which is in jiffies, should be long enough to account for normal transmission delays (such as collisions caused by congestion on the network media).

If the current system time exceeds the device's `trans_start` time by at least the timeout period, the networking layer will eventually call the driver's *tx_timeout* method. That method's job is to do whatever is needed to clear up the problem and to ensure the proper completion of any transmissions that were already in progress. It is important, in particular, that the driver not lose track of any socket buffers that have been entrusted to it by the networking code.

*snull* has the ability to simulate transmitter lockups, which is controlled by two load-time parameters:

```
static int lockup = 0;
MODULE_PARM(lockup, "i");

#ifdef HAVE_TX_TIMEOUT
static int timeout = SNULL_TIMEOUT;
MODULE_PARM(timeout, "i");
#endif
```

If the driver is loaded with the parameter `lockup=n`, a lockup will be simulated once every `n` packets transmitted, and the `watchdog_timeo` field will be set to the given `timeout` value. When simulating lockups, *snull* also calls *netif_stop_queue* to prevent other transmission attempts from occurring.

The *snull* transmission timeout handler looks like this:

```
void snull_tx_timeout (struct net_device *dev)
{
    struct snull_priv *priv = (struct snull_priv *) dev->priv;

    PDEBUG("Transmit timeout at %ld, latency %ld\n", jiffies,
                    jiffies - dev->trans_start);
```

```
        priv->status = SNULL_TX_INTR;
        snull_interrupt(0, dev, NULL);
        priv->stats.tx_errors++;
        netif_wake_queue(dev);
        return;
    }
```

When a transmission timeout happens, the driver must mark the error in the inter-
face statistics and arrange for the device to be reset to a sane state so that new
packets can be transmitted. When a timeout happens in *snull*, the driver calls
*snull_interrupt* to fill in the "missing" interrupt and restarts the transmit queue with
*netif_wake_queue*.

# Packet Reception

Receiving data from the network is trickier than transmitting it because an
`sk_buff` must be allocated and handed off to the upper layers from within an
interrupt handler. The usual way to receive a packet is through an interrupt, unless
the interface is a purely software one like *snull* or the loopback interface.
Although it is possible to write polling drivers, and a few exist in the official ker-
nel, interrupt-driven operation is much better, both in terms of data throughput
and computational demands. Because most network interfaces are interrupt
driven, we won't talk about the polling implementation, which just exploits kernel
timers.

The implementation of *snull* separates the "hardware" details from the device-
independent housekeeping. The function *snull_rx* is thus called after the hardware
has received the packet and it is already in the computer's memory. *snull_rx*
receives a pointer to the data and the length of the packet; its sole responsibility is
to send the packet and some additional information to the upper layers of net-
working code. This code is independent of the way the data pointer and length
are obtained.

```
    void snull_rx(struct net_device *dev, int len, unsigned char *buf)
    {
        struct sk_buff *skb;
        struct snull_priv *priv = (struct snull_priv *) dev->priv;

        /*
         * The packet has been retrieved from the transmission
         * medium. Build an skb around it, so upper layers can handle it
         */
        skb = dev_alloc_skb(len+2);
        if (!skb) {
            printk("snull rx: low on mem - packet dropped\n");
            priv->stats.rx_dropped++;
            return;
        }
        memcpy(skb_put(skb, len), buf, len);
```

```
        /* Write metadata, and then pass to the receive level */
        skb->dev = dev;
        skb->protocol = eth_type_trans(skb, dev);
        skb->ip_summed = CHECKSUM_UNNECESSARY; /* don't check it */
        priv->stats.rx_packets++;
        priv->stats.rx_bytes += len;
        netif_rx(skb);
        return;
    }
```

The function is sufficiently general to act as a template for any network driver, but some explanation is necessary before you can reuse this code fragment with confidence.

The first step is to allocate a buffer to hold the packet. Note that the buffer allocation function (*dev_alloc_skb*) needs to know the data length. The information is used by the function to allocate space for the buffer. *dev_alloc_skb* calls *kmalloc* with atomic priority; it can thus be used safely at interrupt time. The kernel offers other interfaces to socket-buffer allocation, but they are not worth introducing here; socket buffers are explained in detail in "The Socket Buffers," later in this chapter.

Once there is a valid `skb` pointer, the packet data is copied into the buffer by calling *memcpy*; the *skb_put* function updates the end-of-data pointer in the buffer and returns a pointer to the newly created space.

If you are writing a high-performance driver for an interface that can do full bus-mastering I/O, there is a possible optimization that is worth considering here. Some drivers allocate socket buffers for incoming packets prior to their reception, then instruct the interface to place the packet data directly into the socket buffer's space. The networking layer cooperates with this strategy by allocating all socket buffers in DMA-capable space. Doing things this way avoids the need for a separate copy operation to fill the socket buffer, but requires being careful with buffer sizes because you won't know in advance how big the incoming packet is. The implementation of a *change_mtu* method is also important in this situation, since it allows the driver to respond to a change in the maximum packet size.

The network layer needs to have some information spelled out before it will be able to make sense of the packet. To this end, the `dev` and `protocol` fields must be assigned before the buffer is passed upstairs. Then we need to specify how checksumming is to be performed or has been performed on the packet (*snull* does not need to perform any checksums). The possible policies for `skb->ip_summed` are as follows:

CHECKSUM_HW
    The device has already performed checksums in hardware. An example of a hardware checksum is the SPARC HME interface.

`CHECKSUM_NONE`
>    Checksums are still to be verified, and the task must be accomplished by system software. This is the default in newly allocated buffers.

`CHECKSUM_UNNECESSARY`
>    Don't do any checksums. This is the policy in *snull* and in the loopback interface.

Finally, the driver updates its statistics counter to record that a packet has been received. The statistics structure is made up of several fields; the most important are `rx_packets`, `rx_bytes`, `tx_packets`, and `tx_bytes`, which contain the number of packets received and transmitted and the total number of octets transferred. All the fields are thoroughly described in "Statistical Information" later in this chpater.

The last step in packet reception is performed by *netif_rx*, which hands off the socket buffer to the upper layers.

# The Interrupt Handler

Most hardware interfaces are controlled by means of an interrupt handler. The interface interrupts the processor to signal one of two possible events: a new packet has arrived or transmission of an outgoing packet is complete. This generalization doesn't always apply, but it does account for all the problems related to asynchronous packet transmission. Parallel Line Internet Protocol (PLIP) and Point-to-Point Protocol (PPP) are examples of interfaces that don't fit this generalization. They deal with the same events, but the low-level interrupt handling is slightly different.

The usual interrupt routine can tell the difference between a new-packet-arrived interrupt and a done-transmitting notification by checking a status register found on the physical device. The *snull* interface works similarly, but its status word is implemented in software and lives in `dev->priv`. The interrupt handler for a network interface looks like this:

```
void snull_interrupt(int irq, void *dev_id, struct pt_regs *regs)
{
    int statusword;
    struct snull_priv *priv;
    /*
     * As usual, check the "device" pointer for shared handlers.
     * Then assign "struct device *dev"
     */
    struct net_device *dev = (struct net_device *)dev_id;
    /* ... and check with hw if it's really ours */

    if (!dev /*paranoid*/ ) return;

    /* Lock the device */
```

```
    priv = (struct snull_priv *) dev->priv;
    spin_lock(&priv->lock);

    /* retrieve statusword: real netdevices use I/O instructions */
    statusword = priv->status;
    if (statusword & SNULL_RX_INTR) {
        /* send it to snull_rx for handling */
        snull_rx(dev, priv->rx_packetlen, priv->rx_packetdata);
    }
    if (statusword & SNULL_TX_INTR) {
        /* a transmission is over: free the skb */
        priv->stats.tx_packets++;
        priv->stats.tx_bytes += priv->tx_packetlen;
        dev_kfree_skb(priv->skb);
    }

    /* Unlock the device and we are done */
    spin_unlock(&priv->lock);
    return;
}
```

The handler's first task is to retrieve a pointer to the correct `struct net_device`. This pointer usually comes from the `dev_id` pointer received as an argument.

The interesting part of this handler deals with the "transmission done" situation. In this case, the statistics are updated, and *dev_kfree_skb* is called to return the (no longer needed) socket buffer to the system. If your driver has temporarily stopped the transmission queue, this is the place to restart it with *netif_wake_queue*.

Packet reception, on the other hand, doesn't need any special interrupt handling. Calling *snull_rx* (which we have already seen) is all that's required.

## *Changes in Link State*

Network connections, by definition, deal with the world outside the local system. They are thus often affected by outside events, and they can be transient things. The networking subsystem needs to know when network links go up or down, and it provides a few functions that the driver may use to convey that information.

Most networking technologies involving an actual, physical connection provide a *carrier* state; the presence of the carrier means that the hardware is present and ready to function. Ethernet adapters, for example, sense the carrier signal on the wire; when a user trips over the cable, that carrier vanishes, and the link goes down. By default, network devices are assumed to have a carrier signal present. The driver can change that state explicitly, however, with these functions:

```
void netif_carrier_off(struct net_device *dev);
void netif_carrier_on(struct net_device *dev);
```

If your driver detects a lack of carrier on one of its devices, it should call *netif_carrier_off* to inform the kernel of this change. When the carrier returns, *netif_carrier_on* should be called. Some drivers also call *netif_carrier_off* when making major configuration changes (such as media type); once the adapter has finished resetting itself, the new carrier will be detected and traffic can resume.

An integer function also exsists:

```
int netif_carrier_ok(struct net_device *dev);
```

This can be used to test the current carrier state (as reflected in the device structure).

# *The Socket Buffers*

We've now discussed most of the issues related to network interfaces. What's still missing is some more detailed discussion of the `sk_buff` structure. The structure is at the core of the network subsystem of the Linux kernel, and we now introduce both the main fields of the structure and the functions used to act on it.

Although there is no strict need to understand the internals of `sk_buff`, the ability to look at its contents can be helpful when you are tracking down problems and when you are trying to optimize the code. For example, if you look in *loopback.c*, you'll find an optimization based on knowledge of the `sk_buff` internals. The usual warning applies here: if you write code that takes advantage of knowledge of the `sk_buff` structure, you should be prepared to see it break with future kernel releases. Still, sometimes the performance advantages justify the additional maintenance cost.

We are not going to describe the whole structure here, just the fields that might be used from within a driver. If you want to see more, you can look at `<linux/skbuff.h>`, where the structure is defined and the functions are prototyped. Additional details about how the fields and functions are used can be easily retrieved by grepping in the kernel sources.

## *The Important Fields*

The fields introduced here are the ones a driver might need to access. They are listed in no particular order.

```
struct net_device *rx_dev;
struct net_device *dev;
```
   The devices receiving and sending this buffer, respectively.

```
union { /* ... */ } h;
union { /* ... */ } nh;
union { /* ... */} mac;
```
Pointers to the various levels of headers contained within the packet. Each field of the unions is a pointer to a different type of data structure. `h` hosts pointers to transport layer headers (for example, `struct tcphdr *th`); `nh` includes network layer headers (such as `struct iphdr *iph`); and `mac` collects pointers to link layer headers (such as `struct ethdr *ethernet`).

If your driver needs to look at the source and destination addresses of a TCP packet, it can find them in `skb->h.th`. See the header file for the full set of header types that can be accessed in this way.

Note that network drivers are responsible for setting the `mac` pointer for incoming packets. This task is normally handled by *ether_type_trans*, but non-Ethernet drivers will have to set `skb->mac.raw` directly, as shown later in "Non-Ethernet Headers."

```
unsigned char *head;
unsigned char *data;
unsigned char *tail;
unsigned char *end;
```
Pointers used to address the data in the packet. `head` points to the beginning of the allocated space, `data` is the beginning of the valid octets (and is usually slightly greater than `head`), `tail` is the end of the valid octets, and `end` points to the maximum address `tail` can reach. Another way to look at it is that the *available* buffer space is `skb->end - skb->head`, and the *currently used* data space is `skb->tail - skb->data`.

```
unsigned long len;
```
The length of the data itself (`skb->tail - skb->data`).

```
unsigned char ip_summed;
```
The checksum policy for this packet. The field is set by the driver on incoming packets, as was described in "Packet Reception."

```
unsigned char pkt_type;
```
Packet classification used in delivering it. The driver is responsible for setting it to `PACKET_HOST` (this packet is for me), `PACKET_BROADCAST`, `PACKET_MULTICAST`, or `PACKET_OTHERHOST` (no, this packet is not for me). Ethernet drivers don't modify `pkt_type` explicitly because *eth_type_trans* does it for them.

The remaining fields in the structure are not particularly interesting. They are used to maintain lists of buffers, to account for memory belonging to the socket that owns the buffer, and so on.

## Functions Acting on Socket Buffers

Network devices that use a `sock_buff` act on the structure by means of the official interface functions. Many functions operate on socket buffers; here are the most interesting ones:

```
struct sk_buff *alloc_skb(unsigned int len, int priority);
struct sk_buff *dev_alloc_skb(unsigned int len);
```
> Allocate a buffer. The *alloc_skb* function allocates a buffer and initializes both `skb->data` and `skb->tail` to `skb->head`. The *dev_alloc_skb* function is a shortcut that calls *alloc_skb* with `GFP_ATOMIC` priority and reserves some space between `skb->head` and `skb->data`. This data space is used for optimizations within the network layer and should not be touched by the driver.

```
void kfree_skb(struct sk_buff *skb);
void dev_kfree_skb(struct sk_buff *skb);
```
> Free a buffer. The *kfree_skb* call is used internally by the kernel. A driver should use *dev_kfree_skb* instead, which is intended to be safe to call from driver context.

```
unsigned char *skb_put(struct sk_buff *skb, int len);
unsigned char *__skb_put(struct sk_buff *skb, int len);
```
> These inline functions update the `tail` and `len` fields of the `sk_buff` structure; they are used to add data to the end of the buffer. Each function's return value is the previous value of `skb->tail` (in other words, it points to the data space just created). Drivers can use the return value to copy data by invoking `ins(ioaddr, skb_put(...))` or `memcpy(skb_put(...), data, len)`. The difference between the two functions is that *skb_put* checks to be sure that the data will fit in the buffer, whereas *__skb_put* omits the check.

```
unsigned char *skb_push(struct sk_buff *skb, int len);
unsigned char *__skb_push(struct sk_buff *skb, int len);
```
> These functions decrement `skb->data` and increment `skb->len`. They are similar to *skb_put*, except that data is added to the beginning of the packet instead of the end. The return value points to the data space just created. The functions are used to add a hardware header before transmitting a packet. Once again, *__skb_push* differs in that it does not check for adequate available space.

```
int skb_tailroom(struct sk_buff *skb);
```
> This function returns the amount of space available for putting data in the buffer. If a driver puts more data into the buffer than it can hold, the system panics. Although you might object that a *printk* would be sufficient to tag the

error, memory corruption is so harmful to the system that the developers decided to take definitive action. In practice, you shouldn't need to check the available space if the buffer has been correctly allocated. Since drivers usually get the packet size before allocating a buffer, only a severely broken driver will put too much data in the buffer, and a panic might be seen as due punishment.

`int skb_headroom(struct sk_buff *skb);`
Returns the amount of space available in front of `data`, that is, how many octets one can "push" to the buffer.

`void skb_reserve(struct sk_buff *skb, int len);`
This function increments both `data` and `tail`. The function can be used to reserve headroom before filling the buffer. Most Ethernet interfaces reserve 2 bytes in front of the packet; thus, the IP header is aligned on a 16-byte boundary, after a 14-byte Ethernet header. *snull* does this as well, although the instruction was not shown in "Packet Reception" to avoid introducing extra concepts at that point.

`unsigned char *skb_pull(struct sk_buff *skb, int len);`
Removes data from the head of the packet. The driver won't need to use this function, but it is included here for completeness. It decrements `skb->len` and increments `skb->data`; this is how the hardware header (Ethernet or equivalent) is stripped from the beginning of incoming packets.

The kernel defines several other functions that act on socket buffers, but they are meant to be used in higher layers of networking code, and the driver won't need them.

# *MAC Address Resolution*

An interesting issue with Ethernet communication is how to associate the MAC addresses (the interface's unique hardware ID) with the IP number. Most protocols have a similar problem, but we concentrate on the Ethernet-like case here. We'll try to offer a complete description of the issue, so we will show three situations: ARP, Ethernet headers without ARP (like *plip*), and non-Ethernet headers.

## *Using ARP with Ethernet*

The usual way to deal with address resolution is by using ARP, the Address Resolution Protocol. Fortunately, ARP is managed by the kernel, and an Ethernet interface doesn't need to do anything special to support ARP. As long as `dev->addr` and `dev->addr_len` are correctly assigned at open time, the driver doesn't need to worry about resolving IP numbers to physical addresses; *ether_setup* assigns the correct device methods to `dev->hard_header` and `dev->rebuild_header`.

Although the kernel normally handles the details of address resolution (and caching of the results), it calls upon the interface driver to help in the building of the packet. After all, the driver knows about the details of the physical layer header, while the authors of the networking code have tried to insulate the rest of the kernel from that knowledge. To this end, the kernel calls the driver's *hard_header* method to lay out the packet with the results of the ARP query. Normally, Ethernet driver writers need not know about this process—the common Ethernet code takes care of everything.

## Overriding ARP

Simple point-to-point network interfaces such as *plip* might benefit from using Ethernet headers, while avoiding the overhead of sending ARP packets back and forth. The sample code in *snull* also falls into this class of network devices. *snull* cannot use ARP because the driver changes IP addresses in packets being transmitted, and ARP packets exchange IP addresses as well. Although we could have implemented a simple ARP reply generator with little trouble, it is more illustrative to show how to handle physical-layer headers directly.

If your device wants to use the usual hardware header without running ARP, you need to override the default `dev->hard_header` method. This is how *snull* implements it, as a very short function.

```
int snull_header(struct sk_buff *skb, struct net_device *dev,
                 unsigned short type, void *daddr, void *saddr,
                 unsigned int len)
{
    struct ethhdr *eth = (struct ethhdr *)skb_push(skb,ETH_HLEN);

    eth->h_proto = htons(type);
    memcpy(eth->h_source, saddr ? saddr : dev->dev_addr, dev->addr_len);
    memcpy(eth->h_dest,   daddr ? daddr : dev->dev_addr, dev->addr_len);
    eth->h_dest[ETH_ALEN-1]   ^= 0x01;   /* dest is us xor 1 */
    return (dev->hard_header_len);
}
```

The function simply takes the information provided by the kernel and formats it into a standard Ethernet header. It also toggles a bit in the destination Ethernet address, for reasons described later.

When a packet is received by the interface, the hardware header is used in a couple of ways by *eth_type_trans*. We have already seen this call in *snull_rx*:

```
skb->protocol = eth_type_trans(skb, dev);
```

The function extracts the protocol identifier (`ETH_P_IP` in this case) from the Ethernet header; it also assigns `skb->mac.raw`, removes the hardware header from

packet data (with *skb_pull*), and sets `skb->pkt_type`. This last item defaults to `PACKET_HOST` at `skb` allocation (which indicates that the packet is directed to this host), and *eth_type_trans* changes it according to the Ethernet destination address. If that address does not match the address of the interface that received it, the `pkt_type` field will be set to `PACKET_OTHERHOST`. Subsequently, unless the interface is in promiscuous mode, *netif_rx* will drop any packet of type `PACKET_OTHERHOST`. For this reason, *snull_header* is careful to make the destination hardware address match that of the "receiving" interface.

If your interface is a point-to-point link, you won't want to receive unexpected multicast packets. To avoid this problem, remember that a destination address whose first octet has 0 as the least significant bit (LSB) is directed to a single host (i.e., it is either `PACKET_HOST` or `PACKET_OTHERHOST`). The *plip* driver uses 0xfc as the first octet of its hardware address, while *snull* uses 0x00. Both addresses result in a working Ethernet-like point-to-point link.

## Non-Ethernet Headers

We have just seen that the hardware header contains some information in addition to the destination address, the most important being the communication protocol. We now describe how hardware headers can be used to encapsulate relevant information. If you need to know the details, you can extract them from the kernel sources or the technical documentation for the particular transmission medium. Most driver writers will be able to ignore this discussion and just use the Ethernet implementation.

It's worth noting that not all information has to be provided by every protocol. A point-to-point link such as *plip* or *snull* could avoid transferring the whole Ethernet header without losing generality. The *hard_header* device method, shown earlier as implemented by *snull_header*, receives the delivery information—both protocol-level and hardware addresses—from the kernel. It also receives the 16-bit protocol number in the `type` argument; IP, for example, is identified by `ETH_P_IP`. The driver is expected to correctly deliver both the packet data and the protocol number to the receiving host. A point-to-point link could omit addresses from its hardware header, transferring only the protocol number, because delivery is guaranteed independent of the source and destination addresses. An IP-only link could even avoid transmitting any hardware header whatsoever.

When the packet is picked up at the other end of the link, the receiving function in the driver should correctly set the fields `skb->protocol`, `skb->pkt_type`, and `skb->mac.raw`.

`skb->mac.raw` is a char pointer used by the address-resolution mechanism implemented in higher layers of the networking code (for instance, *net/ipv4/arp.c*).

It must point to a machine address that matches `dev->type`. The possible values for the device type are defined in `<linux/if_arp.h>`; Ethernet interfaces use `ARPHRD_ETHER`. For example, here is how *eth_type_trans* deals with the Ethernet header for received packets:

```
skb->mac.raw = skb->data;
skb_pull(skb, dev->hard_header_len);
```

In the simplest case (a point-to-point link with no headers), `skb->mac.raw` can point to a static buffer containing the hardware address of this interface, `proto-col` can be set to `ETH_P_IP`, and `packet_type` can be left with its default value of `PACKET_HOST`.

Because every hardware type is unique, it is hard to give more specific advice than already discussed. The kernel is full of examples, however. See, for example, the AppleTalk driver (*drivers/net/appletalk/cops.c*), the infrared drivers (such as *drivers/net/irda/smc_ircc.c*), or the PPP driver (*drivers/net/ppp_generic.c*).

# *Custom ioctl Commands*

We have seen that the *ioctl* system call is implemented for sockets; `SIOCSIFADDR` and `SIOCSIFMAP` are examples of "socket *ioctls*." Now let's see how the third argument of the system call is used by networking code.

When the *ioctl* system call is invoked on a socket, the command number is one of the symbols defined in `<linux/sockios.h>`, and the function *sock_ioctl* directly invokes a protocol-specific function (where "protocol" refers to the main network protocol being used, for example, IP or AppleTalk).

Any *ioctl* command that is not recognized by the protocol layer is passed to the device layer. These device-related *ioctl* commands accept a third argument from user space, a `struct ifreq *`. This structure is defined in `<linux/if.h>`. The `SIOCSIFADDR` and `SIOCSIFMAP` commands actually work on the `ifreq` structure. The extra argument to `SIOCSIFMAP`, although defined as `ifmap`, is just a field of `ifreq`.

In addition to using the standardized calls, each interface can define its own *ioctl* commands. The *plip* interface, for example, allows the interface to modify its internal timeout values via *ioctl*. The *ioctl* implementation for sockets recognizes 16 commands as private to the interface: `SIOCDEVPRIVATE` through `SIOCDEVPRI-VATE+15`.

When one of these commands is recognized, `dev->do_ioctl` is called in the relevant interface driver. The function receives the same `struct ifreq *` pointer that the general-purpose *ioctl* function uses:

```
int (*do_ioctl)(struct net_device *dev, struct ifreq *ifr, int cmd);
```

The `ifr` pointer points to a kernel-space address that holds a copy of the structure passed by the user. After *do_ioctl* returns, the structure is copied back to user space; the driver can thus use the private commands to both receive and return data.

The device-specific commands can choose to use the fields in `struct ifreq`, but they already convey a standardized meaning, and it's unlikely that the driver can adapt the structure to its needs. The field `ifr_data` is a `caddr_t` item (a pointer) that is meant to be used for device-specific needs. The driver and the program used to invoke its *ioctl* commands should agree about the use of `ifr_data`. For example, *pppstats* uses device-specific commands to retrieve information from the *ppp* interface driver.

It's not worth showing an implementation of *do_ioctl* here, but with the information in this chapter and the kernel examples, you should be able to write one when you need it. Note, however, that the *plip* implementation uses `ifr_data` incorrectly and should not be used as an example for an *ioctl* implementation.

# *Statistical Information*

The last method a driver needs is *get_stats*. This method returns a pointer to the statistics for the device. Its implementation is pretty easy; the one shown works even when several interfaces are managed by the same driver, because the statistics are hosted within the device data structure.

```
struct net_device_stats *snull_stats(struct net_device *dev)
{
    struct snull_priv *priv = (struct snull_priv *) dev->priv;
    return &priv->stats;
}
```

The real work needed to return meaningful statistics is distributed throughout the driver, where the various fields are updated. The following list shows the most interesting fields in `struct net_device_stats`.

`unsigned long rx_packets;`
`unsigned long tx_packets;`
   These fields hold the total number of incoming and outgoing packets successfully transferred by the interface.

`unsigned long rx_bytes;`
`unsigned long tx_bytes;`
   The number of bytes received and transmitted by the interface. These fields were added in the 2.2 kernel.

```
unsigned long rx_errors;
unsigned long tx_errors;
```
The number of erroneous receptions and transmissions. There's no end of things that can go wrong with packet transmission, and the `net_device_stats` structure includes six counters for specific receive errors and five for transmit errors. See `<linux/netdevice.h>` for the full list. If possible, your driver should maintain detailed error statistics, because they can be most helpful to system administrators trying to track down a problem.

```
unsigned long rx_dropped;
unsigned long tx_dropped;
```
The number of packets dropped during reception and transmission. Packets are dropped when there's no memory available for packet data. `tx_dropped` is rarely used.

```
unsigned long collisions;
```
The number of collisions due to congestion on the medium.

```
unsigned long multicast;
```
The number of multicast packets received.

It is worth repeating that the *get_stats* method can be called at any time—even when the interface is down—so the driver should not release statistic information when running the *stop* method.

## *Multicasting*

A *multicast* packet is a network packet meant to be received by more than one host, but not by all hosts. This functionality is obtained by assigning special hardware addresses to groups of hosts. Packets directed to one of the special addresses should be received by all the hosts in that group. In the case of Ethernet, a multicast address has the least significant bit of the first address octet set in the destination address, while every device board has that bit clear in its own hardware address.

The tricky part of dealing with host groups and hardware addresses is performed by applications and the kernel, and the interface driver doesn't need to deal with these problems.

Transmission of multicast packets is a simple problem because they look exactly like any other packets. The interface transmits them over the communication medium without looking at the destination address. It's the kernel that has to assign a correct hardware destination address; the *hard_header* device method, if defined, doesn't need to look in the data it arranges.

The kernel handles the job of tracking which multicast addresses are of interest at any given time. The list can change frequently, since it is a function of the applications that are running at any given time and the user's interest. It is the driver's job to accept the list of interesting multicast addresses and deliver to the kernel any packets sent to those addresses. How the driver implements the multicast list is somewhat dependent on how the underlying hardware works. Typically, hardware belongs to one of three classes, as far as multicast is concerned:

- Interfaces that cannot deal with multicast. These interfaces either receive packets directed specifically to their hardware address (plus broadcast packets), or they receive every packet. They can receive multicast packets only by receiving every packet, thus potentially overwhelming the operating system with a huge number of "uninteresting" packets. You don't usually count these interfaces as multicast capable, and the driver won't set `IFF_MULTICAST` in `dev->flags`.

  Point-to-point interfaces are a special case, because they always receive every packet without performing any hardware filtering.

- Interfaces that can tell multicast packets from other packets (host-to-host or broadcast). These interfaces can be instructed to receive every multicast packet and let the software determine if this host is a valid recipient. The overhead introduced in this case is acceptable, because the number of multicast packets on a typical network is low.

- Interfaces that can perform hardware detection of multicast addresses. These interfaces can be passed a list of multicast addresses for which packets are to be received, and they will ignore other multicast packets. This is the optimum case for the kernel, because it doesn't waste processor time dropping "uninteresting" packets received by the interface.

The kernel tries to exploit the capabilities of high-level interfaces by supporting at its best the third device class, which is the most versatile. Therefore, the kernel notifies the driver whenever the list of valid multicast addresses is changed, and it passes the new list to the driver so it can update the hardware filter according to the new information.

## Kernel Support for Multicasting

Support for multicast packets is made up of several items: a device method, a data structure and device flags.

```
void (*dev->set_multicast_list)(struct net_device *dev);
```
   This device method is called whenever the list of machine addresses associated with the device changes. It is also called when `dev->flags` is modified, because some flags (e.g., `IFF_PROMISC`) may also require you to reprogram the hardware filter. The method receives a pointer to `struct net_device` as an argument and returns `void`. A driver not interested in implementing this

method can leave the field set to `NULL`.

`struct dev_mc_list *dev->mc_list;`
> This is a linked list of all the multicast addresses associated with the device. The actual definition of the structure is introduced at the end of this section.

`int dev->mc_count;`
> The number of items in the linked list. This information is somewhat redundant, but checking `mc_count` against 0 is a useful shortcut for checking the list.

`IFF_MULTICAST`
> Unless the driver sets this flag in `dev->flags`, the interface won't be asked to handle multicast packets. The *set_multicast_list* method will nonetheless be called when `dev->flags` changes, because the multicast list may have changed while the interface was not active.

`IFF_ALLMULTI`
> This flag is set in `dev->flags` by the networking software to tell the driver to retrieve all multicast packets from the network. This happens when multicast routing is enabled. If the flag is set, `dev->mc_list` shouldn't be used to filter multicast packets.

`IFF_PROMISC`
> This flag is set in `dev->flags` when the interface is put into promiscuous mode. Every packet should be received by the interface, independent of `dev->mc_list`.

The last bit of information needed by the driver developer is the definition of `struct dev_mc_list`, which lives in *<linux/netdevice.h>*.

```
struct dev_mc_list {
    struct dev_mc_list   *next;           /* Next address in list */
    __u8                 dmi_addr[MAX_ADDR_LEN]; /* Hardware address */
    unsigned char        dmi_addrlen;     /* Address length */
    int                  dmi_users;       /* Number of users */
    int                  dmi_gusers;      /* Number of groups */
};
```

Because multicasting and hardware addresses are independent of the actual transmission of packets, this structure is portable across network implementations, and each address is identified by a string of octets and a length, just like `dev->dev_addr`.

## A Typical Implementation

The best way to describe the design of *set_multicast_list* is to show you some pseudocode.

The following function is a typical implementation of the function in a full-featured (`ff`) driver. The driver is full featured in that the interface it controls has a complex hardware packet filter, which can hold a table of multicast addresses to be received by this host. The maximum size of the table is `FF_TABLE_SIZE`.

All the functions prefixed with `ff_` are placeholders for hardware-specific operations.

```
void ff_set_multicast_list(struct net_device *dev)
{
    struct dev_mc_list *mcptr;

    if (dev->flags & IFF_PROMISC) {
        ff_get_all_packets();
        return;
    }
    /* If there's more addresses than we handle, get all multicast
    packets and sort them out in software. */
    if (dev->flags & IFF_ALLMULTI || dev->mc_count > FF_TABLE_SIZE) {
        ff_get_all_multicast_packets();
        return;
    }
    /* No multicast?  Just get our own stuff */
    if (dev->mc_count == 0) {
        ff_get_only_own_packets();
        return;
    }
    /* Store all of the multicast addresses in the hardware filter */
    ff_clear_mc_list();
    for (mc_ptr = dev->mc_list; mc_ptr; mc_ptr = mc_ptr->next)
        ff_store_mc_address(mc_ptr->dmi_addr);
    ff_get_packets_in_multicast_list();
}
```

This implementation can be simplified if the interface cannot store a multicast table in the hardware filter for incoming packets. In that case, `FF_TABLE_SIZE` reduces to 0 and the last four lines of code are not needed.

As was mentioned earlier, even interfaces that can't deal with multicast packets need to implement the *set_multicast_list* method to be notified about changes in `dev->flags`. This approach could be called a "nonfeatured" (`nf`) implementation. The implementation is very simple, as shown by the following code:

```
void nf_set_multicast_list(struct net_device *dev)
{
    if (dev->flags & IFF_PROMISC)
        nf_get_all_packets();
    else
        nf_get_only_own_packets();
}
```

Implementing `IFF_PROMISC` is important, because otherwise the user won't be able to run *tcpdump* or any other network analyzers. If the interface runs a point-to-point link, on the other hand, there's no need to implement *set_multicast_list* at all, because users receive every packet anyway.

# Backward Compatibility

Version 2.3.43 of the kernel saw a major rework of the networking subsystem. The new "softnet" implementation was a great improvement in terms of performance and clean design. It also, of course, brought changes to the network driver interface—though fewer than one might have expected.

## Differences in Linux 2.2

First of all, Linux 2.3.14 renamed the network device structure, which had always been `struct device`, to `struct net_device`. The new name is certainly more appropriate, since the structure was never meant to describe devices in general.

Prior to version 2.3.43, the functions *netif_start_queue*, *netif_stop_queue*, and *netif_wake_queue* did not exist. Packet transmission was, instead, controlled by three fields in the `device` structure, and *sysdep.h* implements the three functions using the three fields when compiling for 2.2 or 2.0.

`unsigned char start;`
> This variable indicated that the interface was ready for operations; it was normally set to 1 in the driver's *open* method. The current implementation is to call *netif_start_queue* instead.

`unsigned long interrupt;`
> `interrupt` was used to indicate that the device was servicing an interrupt—accordingly, it was set to 1 at the beginning of the interrupt handler and to 0 before returning. It was never a substitute for proper locking, and its use has been replaced with internal spinlocks.

`unsigned long tbusy;`
> When nonzero, this variable indicated that the device could handle no more outgoing packets. Where a 2.4 driver will call *netif_stop_queue*, older drivers would set `tbusy` to 1. Restarting the queue required setting `tbusy` back to 0 and calling `mark_bh(NET_BH)`.

Normally, setting `tbusy` was sufficient to ensure that the driver's *hard_start_xmit* method would not be called. However, if the networking system decided that a transmitter lockup must have occurred, it would call that method anyway. There was no *tx_timeout* method before softnet was integrated. Thus, pre-softnet drivers had to explicitly check for a call to *hard_start_xmit* when `tbusy` was set and react accordingly.

The type of the `name` field in `struct device` was different. The 2.2 version was simply

```
char *name;
```

Thus, the storage for the interface name had to be allocated separately, and `name` assigned to point to that storage. Usually the device name was stored in a static variable within the driver. The `%d` notation for dynamically assigned interface names was not present in 2.2; instead, if the name began with a null byte or a space character, the kernel would allocate the next `eth` name. The 2.4 kernel still implements this behavior, but its use is deprecated. Starting with 2.5, only the `%d` format is likely to be recognized.

The `owner` field (and the `SET_MODULE_OWNER` macro) were added in kernel 2.4.0-test11, just before the official stable release. Previously, network driver modules had to maintain their own use counts. *sysdep.h* defines an empty `SET_MODULE_OWNER` for kernels that do not have it; portable code should also continue to manage its use count manually (in addition to letting the networking system do it).

The link state functions (*netif_carrier_on* and *netif_carrier_off*) did not exist in the 2.2 kernel. The kernel simply did without that information in those days.

## *Further Differences in Linux 2.0*

The 2.1 development series also saw its share of changes to the network driver interface. Most took the form of small changes to function prototypes, rather than sweeping changes to the network code as a whole.

Interface statistics were kept in a structure called `struct 1enet_statistics`, defined in `<linux/if_ether.h>`. Even non-Ethernet drivers used this structure. The field names were all the same as the current `struct net_device_stats`, but the `rx_bytes` and `tx_bytes` fields were not present.

The 2.0 kernel handled transmitter lockups in the same way as 2.2 did. There was, however, an additional function:

```
void dev_tint(struct device *dev);
```

This function would be called by the driver after a lockup had been cleared to restart the transmission of packets.

A couple of functions had different prototypes. *dev_kfree_skb* had a second, integer argument that was either `FREE_READ` for incoming packets (i.e., `skb`s allocated by the driver) or `FREE_WRITE` for outgoing packets (`skb`s allocated by the networking code). Almost all calls to *dev_kfree_skb* in network driver code used `FREE_WRITE`. The nonchecking versions of the `skb` functions (such as *__skb_push*) did not exist; *sysdep.h* in the sample code provides emulation for these functions under 2.0.

The *rebuild_header* method had a different set of arguments:

```
int (*rebuild_header) (void *eth, struct device *dev,
unsigned long raddr, struct sk_buff *skb);
```

The Linux kernel also made heavier use of *rebuild_header*; it did most of the work that *hard_header* does now. When *snull* is compiled under Linux 2.0, it builds hardware headers as follows:

```
int snull_rebuild_header(void *buff, struct net_device *dev, unsigned long dst,
                   struct sk_buff *skb)
{
    struct ethhdr *eth = (struct ethhdr *)buff;

    memcpy(eth->h_source, dev->dev_addr, dev->addr_len);
    memcpy(eth->h_dest, dev->dev_addr, dev->addr_len);
    eth->h_dest[ETH_ALEN-1]   ^= 0x01;   /* dest is us xor 1 */
    return 0;
}
```

The device methods for header caching were also significantly different in this kernel. If your driver needs to implement these functions directly (very few do), and it also needs to work with the 2.0 kernel, see the definitions in `<linux/netdevice.h>` to see how things were done in those days.

## Probing and HAVE_DEVLIST

If you look at the source for almost any network driver in the kernel, you will find some boilerplate that looks like this:

```
#ifdef HAVE_DEVLIST
/*
 * Support for an alternate probe manager,
 * which will eliminate the boilerplate below.
 */
struct netdev_entry netcard_drv =
{cardname, netcard_probe1, NETCARD_IO_EXTENT, netcard_portlist};
#else
/* Regular probe routine defined here */
```

Interestingly, this code has been around since the 1.1 development series, but we are still waiting for the promised alternate probe manager. It is probably safe to not worry about being prepared for this great change, especially since ideas for how to implement it will likely have changed in the intervening years.

# *Quick Reference*

This section provides a reference for the concepts introduced in this chapter. It also explains the role of each header file that a driver needs to include. The lists of fields in the `net_device` and `sk_buff` structures, however, are not repeated here.

`#include <linux/netdevice.h>`
> This header hosts the definitions of `struct net_device` and `struct net_device_stats`, and includes a few other headers that are needed by network drivers.

`int register_netdev(struct net_device *dev);`
`void unregister_netdev(struct net_device *dev);`
> Register and unregister a network device.

`SET_MODULE_OWNER(struct net_device *dev);`
> This macro will store a pointer to the current module in the device structure (or in any structure with an `owner` field, actually); it is used to enable the networking subsystem to manage the module's use count.

`netif_start_queue(struct net_device *dev);`
`netif_stop_queue(struct net_device *dev);`
`netif_wake_queue(struct net_device *dev);`
> These functions control the passing of packets to the driver for transmission. No packets will be transmitted until *netif_start_queue* has been called. *netif_stop_queue* suspends transmission, and *netif_wake_queue* restarts the queue and pokes the network layer to restart transmitting packets.

`void netif_rx(struct sk_buff *skb);`
> This function can be called (including at interrupt time) to notify the kernel that a packet has been received and encapsulated into a socket buffer.

`#include <linux/if.h>`
> Included by *netdevice.h*, this file declares the interface flags (`IFF_` macros) and `struct ifmap`, which has a major role in the *ioctl* implementation for network drivers.

`void netif_carrier_off(struct net_device *dev);`
`void netif_carrier_on(struct net_device *dev);`
`int netif_carrier_ok(struct net_device *dev);`
> The first two functions may be used to tell the kernel whether a carrier signal is currently present on the given interface. *netif_carrier_ok* will test the carrier state as reflected in the device structure.

```
#include <linux/if_ether.h>
ETH_ALEN
ETH_P_IP
struct ethhdr;
```
Included by *netdevice.h*, *if_ether.h* defines all the `ETH_` macros used to represent octet lengths (such as the address length) and network protocols (such as IP). It also defines the `ethhdr` structure.

```
#include <linux/skbuff.h>
```
The definition of `struct sk_buff` and related structures, as well as several inline functions to act on the buffers. This header is included by *netdevice.h*.

```
struct sk_buff *alloc_skb(unsigned int len, int priority);
struct sk_buff *dev_alloc_skb(unsigned int len);
void kfree_skb(struct sk_buff *skb);
void dev_kfree_skb(struct sk_buff *skb);
```
These functions handle the allocation and freeing of socket buffers. Drivers should normally use the `dev_` variants, which are intended for that purpose.

```
unsigned char *skb_put(struct sk_buff *skb, int len);
unsigned char *__skb_put(struct sk_buff *skb, int len);
unsigned char *skb_push(struct sk_buff *skb, int len);
unsigned char *__skb_push(struct sk_buff *skb, int len);
```
These functions add data to an `skb`; *skb_put* puts the data at the end of the `skb`, while *skb_push* puts it at the beginning. The regular versions perform checking to ensure that adequate space is available; double-underscore versions leave those tests out.

```
int skb_headroom(struct sk_buff *skb);
int skb_tailroom(struct sk_buff *skb);
void skb_reserve(struct sk_buff *skb, int len);
```
These functions perform management of space within an `skb`. *skb_headroom* and *skb_tailroom* tell how much space is available at the beginning and end, respectively, of an `skb`. *skb_reserve* may be used to reserve space at the beginning of an `skb`, which must be empty.

```
unsigned char *skb_pull(struct sk_buff *skb, int len);
```
*skb_pull* will "remove" data from an `skb` by adjusting the internal pointers.

```
#include <linux/etherdevice.h>
void ether_setup(struct net_device *dev);
```
This function sets most device methods to the general-purpose implementation for Ethernet drivers. It also sets `dev->flags` and assigns the next available `ethx` name to `dev->name` if the first character in the name is a blank space or the null character.

```
unsigned short eth_type_trans(struct sk_buff *skb, struct
     net_device *dev);
```
When an Ethernet interface receives a packet, this function can be called to set `skb->pkt_type`. The return value is a protocol number that is usually stored in `skb->protocol`.

```
#include <linux/sockios.h>
```
SIOCDEVPRIVATE
This is the first of 16 *ioctl* commands that can be implemented by each driver for its own private use. All the network *ioctl* commands are defined in *sockios.h*.