# Exercise 8. Write on how a SNULL(Simple Network Utility for Loading Localities) works, need not execute.

---

The network interfaces fit in with the rest of the Linux kernel and provides examples in the form of a memory-based modularized network interface, which is called *snull*.

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 different only in tiny details related to the actual network protocol.

## 1. Entry and exit functions

The Linux network device driver basically makes a fuss about the net_device structure. The entry function of the program is mainly to allocate and register the net_device structure. Correspondingly, the exit function needs to cancel the logout and release the net_device structure. This experiment will eventually use the ping command to test the sending and receiving of two network devices, so first define two pointers to net_device:

struct net_device *snull_devs[2];

### 1.1 Entry function:

```
1.  static __init int snull_init_module(void)
2.  {
3.     int result, i, ret = -ENOMEM;
4.     /* Allocate the devices */
5.          snull_devs[0] = alloc_netdev(sizeof(struct snull_priv), "sn%d", snull_init);
6.          snull_devs[1] = alloc_netdev(sizeof(struct snull_priv), "sn%d", snull_init);
7.          if (snull_devs[0] == NULL || snull_devs[1] == NULL)
8.                  goto out;
9.     ret = -ENODEV;
10.         for (i = 0; i < 2;  i++)
11.                 if ((result = register_netdev(snull_devs[i])))
12.                         printk("snull: error %i registering device \"%s\"\n",
13.                                         result, snull_devs[i]->name);
14.                 else
15.                         ret = 0;
16.  out:
```

```
17.        if (ret)
18.               snull_cleanup();
19.        return ret;
20.}
21. module_init(snull_init_module);
```

Alloc_netdev() function is used to allocate net_device structure:

```
1.  struct net_device *alloc_netdev(int sizeof_priv, const char *name,
2.                         void (*setup)(struct net_device *));
```

The first parameter of the function is the size of the private data, the second parameter is the device name, and the third parameter setup points to a function that takes struct net_device* as a parameter and is used to initialize some members of net_device.

1.2 The private data of the program defines snull_priv (the usage of spin lock members is not repeated):

```
1.  struct snull_priv {
2.         struct net_device_stats stats;
3.         int status;
4.         struct snull_packet *ppool;
5.         struct snull_packet *rx_queue;
6.         int rx_int_enabled;
7.         int tx_packetlen;
8.         u8 *tx_packetdata;
9.         struct sk_buff *skb;
10.        spinlock_t lock;
11. };
```

1.2.1 The snull_packet structure is used to save the sent and received data:

```
1.  struct snull_packet {
2.         struct snull_packet *next;
3.         struct net_device *dev;
4.         int     datalen;
5.         u8 data[ETH_DATA_LEN];
6.  };
```

The series of operation functions defined for struct snull_packet *pool are as follows:

```
1.  int pool_size = 8;
2.  module_param(pool_size, int, 0);
```

```
3.  /*
4.     * Allocate 8 snull_packets and add them to a linked list
5.  */
6.  void snull_setup_pool(struct net_device *dev)
7.  {
8.          struct snull_priv *priv = netdev_priv(dev);
9.          int i;
10.         struct snull_packet *pkt;
11.         priv->ppool = NULL;
12.         for (i = 0; i < pool_size; i++) {
13.                 pkt = kmalloc (sizeof (struct snull_packet), GFP_KERNEL);
14.                 if (pkt == NULL) {
15.                         printk (KERN_NOTICE "Ran out of memory allocating
    packet pool\n");
16.                         return;
17.                 }
18.                 pkt->dev = dev;
19.                 pkt->next = priv->ppool;
20.                 priv->ppool = pkt;
21.         }
22. }
23. /*
24.    * Release the linked list composed of snull_packet
25. */
26. void snull_teardown_pool(struct net_device *dev)
27. {
28.         struct snull_priv *priv = netdev_priv(dev);
29.         struct snull_packet *pkt;
30.
31.         while ((pkt = priv->ppool)) {
32.                 priv->ppool = pkt->next;
33.                 kfree (pkt);
34.         }
35. }
36. /*
37.    * Take a snull packet from the linked list for sending
38. */
39. struct snull_packet *snull_get_tx_buffer(struct net_device *dev)
40. {
41.         struct snull_priv *priv = netdev_priv(dev);
```

```
42.        unsigned long flags;
43.        struct snull_packet *pkt;
44.
45.        spin_lock_irqsave(&priv->lock, flags);
46.        pkt = priv->ppool;
47.        priv->ppool = pkt->next;
48.        if (priv->ppool == NULL) {
49.                printk (KERN_INFO "Pool empty\n");
50.                netif_stop_queue(dev);
51.        }
52.        spin_unlock_irqrestore(&priv->lock, flags);
53.        return pkt;
54. }
55. /*
56.  * Put a snull packet back to the linked list
57. */
58. void snull_release_buffer(struct snull_packet *pkt)
59. {
60.        unsigned long flags;
61.        struct snull_priv *priv = netdev_priv(pkt->dev);
62.
63.        spin_lock_irqsave(&priv->lock, flags);
64.        pkt->next = priv->ppool;
65.        priv->ppool = pkt;
66.        spin_unlock_irqrestore(&priv->lock, flags);
67.        if (netif_queue_stopped(pkt->dev) && pkt->next == NULL)
68.                netif_wake_queue(pkt->dev);
69. }
```

1.2.2 rx_queue is the receiving queue of the device, the operation functions are as follows:

```
1.  void snull_enqueue_buf(struct net_device *dev, struct snull_packet *pkt)
2.  {
3.        unsigned long flags;
4.        struct snull_priv *priv = netdev_priv(dev);
5.        spin_lock_irqsave(&priv->lock, flags);
6.        pkt->next = priv->rx_queue;  /* FIXME - misorders packets */
7.        priv->rx_queue = pkt;
8.        spin_unlock_irqrestore(&priv->lock, flags);
9.  }
```

```
10. struct snull_packet *snull_dequeue_buf(struct net_device *dev)
11. {
12.         struct snull_priv *priv = netdev_priv(dev);
13.         struct snull_packet *pkt;
14.         unsigned long flags;
15.         spin_lock_irqsave(&priv->lock, flags);
16.         pkt = priv->rx_queue;
17.         if (pkt != NULL)
18.                 priv->rx_queue = pkt->next;
19.         spin_unlock_irqrestore(&priv->lock, flags);
20.         return pkt;
21. }
```

1.2.3 rx_int_enabled is used to enable or disable the receiving interrupt (the interrupt is not a real hardware interrupt, but is simulated by software), the operation function is as follows:

```
1.  /*
2.  * Enable and disable receive interrupts.
3.  */
4.  static void snull_rx_ints(struct net_device *dev, int enable)
5.  {
6.          struct snull_priv *priv = netdev_priv(dev);
7.          priv->rx_int_enabled = enable;
8.  }
```

1.2.4 struct skb_buff socket buffer, this structure is critical, used to transfer data between different layers of the network protocol, the skb_buff structure will be further explained later;

1.3 The setup parameter of alloc_netdevice() points to a function that is used to initialize other members of net_device. The snull_init function is used in the program:

```
1.  static int timeout = SNULL_TIMEOUT; //SNULL_TIMEOUT is defined as 5 in
    snull.h
2.  module_param(timeout,int,0);
3.  static const struct net_device_ops snull_dev_ops = {
4.    .ndo_open = snull_open,
5.    .ndo_stop = snull_release,
6.    .ndo_start_xmit = snull_tx,
7.    .ndo_do_ioctl = snull_ioctl,
8.    .ndo_get_stats = snull_stats,
```

```
9.     .ndo_tx_timeout = snull_tx_timeout,
10. };
11. static const struct header_ops snull_header_ops= {
12.   .create   = snull_header,
13.   .rebuild = snull_rebuild_header,
14.   .cache = NULL,
15. };
16. void snull_init(struct net_device *dev)
17. {
18.   struct snull_priv *priv = NULL;
19.   ether_setup(dev);
20.
21.   dev->netdev_ops = &snull_dev_ops;
22.   dev->header_ops = &snull_header_ops;
23.   dev->watchdog_timeo = timeout;
24.      dev->flags |= IFF_NOARP; //Forbid ARP
25.   dev->features |= NETIF_F_NO_CSUM;
26.   priv = netdev_priv(dev);
27.   memset(priv, 0, sizeof(struct snull_priv));
28.   spin_lock_init(&((struct snull_priv *)priv)->lock);
29.      snull_rx_ints(dev, 1);                /* enable receive interrupts */
30.      snull_setup_pool(dev);
31.
32.   return;
33. }
```

snull_init() first initializes the members in the net_device structure: call the ether_setup() function to assign default values to many members in the net_device; Next, set the device method, because the kernel version is different, the device method is directly included in the net_device structure, which is required here Modified, from the code, you can see that the operations such as open and release are encapsulated in the netdev_ops, header_ops structure, and the specific implementation of the device methods such as open, send, and receive will be further explained later. The watchdog_timeo member sets the transmission timeout period, followed by some flags, indicating that ARP is forbidden and force packet verification.

After setting up the net_device member, use netdev_priv() to obtain the private data of the device and set the spin lock, receive enable and data buffer. netdev_priv() is dedicated to access private data of network devices:

*void *netdev_priv(struct net_device *dev);*

1.4 After allocating and initializing the net_device structure, we also need to register this device with register_netdev, and register the function prototype for unregistering the network device:

1. *int register_netdev(struct net_device *dev);*
2. *void unregister_netdev(struct net_device *dev);*

1.5 **Export function:**

1. *static __exit void snull_cleanup(void)*
2. *{*
3. *    int i;*
4. 
5. *    for (i=0; i<2;i++)*
6. *    {*
7. *        if(snull_devs[i])*
8. *        {*
9. *            unregister_netdev(snull_devs[i]);*
10. *                        snull_teardown_pool(snull_devs[i]);*
11. *            free_netdev(snull_devs[i]);*
12. *        }*
13. *    }*
14. *    return;*
15. *}*
16. *module_exit(snull_cleanup);*

The operation of the exit function is the opposite of the entry, mainly to unregister the device, release the data buffer and release the memory occupied by the device.

The explanation for the parameters used in 1.3:

Initializing the net_device member is mainly to implement some device operation methods, as follows:

1. static const struct net_device_ops snull_dev_ops = {
2.     .ndo_open = snull_open,
3.     .ndo_stop = snull_release,
4.     .ndo_start_xmit = snull_tx,
5.     .ndo_do_ioctl = snull_ioctl,
6.     .ndo_get_stats = snull_stats,
7.     .ndo_tx_timeout = snull_tx_timeout,
8. };

```
9.  static const struct header_ops snull_header_ops= {
10.     .create   = snull_header,
11.     .rebuild = snull_rebuild_header,
12.     .cache = NULL,
13. };
```

The following specifically describes how to implement these functions:

1. snull_open: open the device

```
1.  int snull_open(struct net_device *dev)
2.  {
3.     memcpy(dev->dev_addr, "\0SNUL0", ETH_ALEN);
4.     if(dev == snull_devs[1])
5.        dev->dev_addr[ETH_ALEN-1] ++;
6.     netif_start_queue(dev);
7.     return 0;
8.  }
```

The open function assigns snull_devs[0] and snull_devs[1] addresses, and starts the transmission queue of the interface. The transmission queue operations include:

1. netif_start_queue(struct net_device *dev); //Start transmission queue
2. netif_stop_queue(struct net_device *dev); //Stop transmission queue
3. netif_wake_queue(struct net_device *dev); //Restart the transmission queue

2. snull_release: Called when the interface is closed, it simply stops the transmission queue:

```
1.  int snull_release(struct net_device *dev)
2.  {
3.     netif_stop_queue(dev);
4.     return 0;
5.  }
```

3. snull_tx: transmit data packets

```
1.  int snull_tx(struct sk_buff *skb, struct net_device *dev)
2.  {
3.         int len;
4.         char *data, shortpkt[ETH_ZLEN];
5.         struct snull_priv *priv = netdev_priv(dev);
6.
```

```
7.          data = skb->data;
8.          len = skb->len;
9.          if (len < ETH_ZLEN) {
10.                 memset(shortpkt, 0, ETH_ZLEN);
11.                 memcpy(shortpkt, skb->data, skb->len);
12.                 len = ETH_ZLEN;
13.                 data = shortpkt;
14.         }
15.         dev->trans_start = jiffies; /* save the timestamp */
16.         /* Remember the skb, so we can free it at interrupt time */
17.         priv->skb = skb;
18.         /* actual deliver of data is device-specific, and not shown here */
19.         snull_hw_tx(data, len, dev);
20.         return 0; /* Our simple device can not fail */
21.}
```

skb is a pointer to sk_buff, get the data and data length to be sent from skb, and then call snull_hw_tx function to send, snull_hw_tx code is as follows:

```
1.  void snull_hw_tx(char *buf, int len, struct net_device *dev)
2.  {
3.    struct iphdr *ih;
4.    struct net_device *dest;
5.    struct snull_priv *priv = netdev_priv(dev);
6.    u32 *saddr, *daddr;
7.        struct snull_packet *tx_buffer;
8.    if (len < sizeof(struct ethhdr) + sizeof(struct iphdr))
9.    {
10.     printk("snull: Hmm... packet too short (%i octets)\n",len);
11.     return;
12.  }
13.  ih = (struct iphdr *)(buf+sizeof(struct ethhdr));
14.  saddr = &ih->saddr;
15.  daddr = &ih->daddr;
16.
17.       if (dev == snull_devs[0])
18.     printk(KERN_INFO"before change: %08x:%05i --> %08x:%05i\n",
19.         ntohl(ih->saddr),ntohs(((struct tcphdr *)(ih+1))->source),
20.         ntohl(ih->daddr),ntohs(((struct tcphdr *)(ih+1))->dest));
21.  else
22.     printk(KERN_INFO"before change: %08x:%05i <-- %08x:%05i\n",
```

```
23.          ntohl(ih->daddr),ntohs(((struct tcphdr *)(ih+1))->dest),
24.          ntohl(ih->saddr),ntohs(((struct tcphdr *)(ih+1))->source));
25.
26.   ((u8 *)saddr)[2] ^= 1;
27.   ((u8 *)daddr)[2] ^= 1;
28.
29.   ih->check = 0;
30.   ih->check = ip_fast_csum((unsigned char *)ih,ih->ihl);
31.   if (dev == snull_devs[0])
32.     printk(KERN_INFO"after change: %08x:%05i --> %08x:%05i\n",
33.          ntohl(ih->saddr),ntohs(((struct tcphdr *)(ih+1))->source),
34.          ntohl(ih->daddr),ntohs(((struct tcphdr *)(ih+1))->dest));
35.   else
36.     printk(KERN_INFO"after change: %08x:%05i <-- %08x:%05i\n",
37.          ntohl(ih->daddr),ntohs(((struct tcphdr *)(ih+1))->dest),
38.          ntohl(ih->saddr),ntohs(((struct tcphdr *)(ih+1))->source));
39.
40.       dest = snull_devs[dev == snull_devs[0] ? 1 : 0];
41.       priv = netdev_priv(dest);
42.       tx_buffer = snull_get_tx_buffer(dev);
43.       tx_buffer->datalen = len;
44.       memcpy(tx_buffer->data, buf, len);
45.       snull_enqueue_buf(dest, tx_buffer);
46.       if (priv->rx_int_enabled) {
47.              priv->status |= SNULL_RX_INTR;
48.              snull_interrupt(0, dest, NULL);
49.       }
50.       priv = netdev_priv(dev);
51.   priv->status = SNULL_TX_INTR;
52.   priv->tx_packetlen = len;
53.   priv->tx_packetdata = buf;
54.   if (lockup && ((priv->stats.tx_packets + 1) % lockup) == 0)
55.   {
56.     netif_stop_queue(dev);
57.     printk(KERN_INFO"Simulate lockup at %ld, txp %ld\n", jiffies,(unsigned long)
   priv->stats.tx_packets);
58.   }
59.   else
60.     snull_interrupt(0, dev, NULL);
61.   return;
```

62. }

This function modifies the source address and destination address of the data packet to produce the effect of pinging two network interfaces. For example, setting several IP addresses is as follows:

192.168.0.1 : local0

192.168.0.2 : remote0

192.168.1.2 : local1

192.168.1.1 : remote1

When sending ping -c 1 remote0, the source address is local0, the destination address is remote0, and the sending function modifies it to remote1 and local1. After local1 receives the data, it sends a response to remote1 The sending function changes the source and destination addresses to remote0 and local0, and finally local0 receives the data, causing the illusion of pinging. Here add the data packet to the receive queue, call snull_interrupt to imitate a hardware interrupt, and set the state of the two network interfaces. The snull_interrupt function is:

1.  *void snull_interrupt(int irq, void *dev_id, struct pt_regs *regs)*
2.  *{*
3.  *int statusword;*
4.  *struct snull_priv *priv;*
5.  *struct snull_packet *pkt = NULL;*
6.  
7.  *struct net_device *dev = (struct net_device *)dev_id;*
8.  
9.  *if (!dev)*
10.  *return;*
11.  */* Lock the device */*
12.  *priv = netdev_priv(dev);*
13.  *spin_lock(&priv->lock);*
14.  
15.  *statusword = priv->status;*
16.  *priv->status = 0;*
17.  *if (statusword & SNULL_RX_INTR)*
18.  *{*
19.  *if (dev_id == snull_devs[0])*
20.  *printk(KERN_INFO"sn0 rx\n");*

```
21.              if (dev_id == snull_devs[1])
22.                      printk(KERN_INFO"sn1 rx\n");
23.     /* send it to snull_rx for handling */
24.              pkt = priv->rx_queue;
25.              if (pkt) {
26.                      priv->rx_queue = pkt->next;
27.                      snull_rx(dev, pkt);
28.              }
29.  }
30.  if (statusword & SNULL_TX_INTR)
31.  {
32.              if (dev_id == snull_devs[0])
33.                      printk(KERN_INFO"sn0 tx\n");
34.              if (dev_id == snull_devs[1])
35.                      printk(KERN_INFO"sn1 tx\n");
36.     priv->stats.tx_packets++;
37.     priv->stats.tx_bytes += priv->tx_packetlen;
38.     dev_kfree_skb(priv->skb);
39.  }
40.  spin_unlock(&priv->lock);
41.       if (pkt)
42.              snull_release_buffer(pkt); /* Do this outside the lock! */
43.
44.  return;
45.}
```

If it is to receive an interrupt, call the snull_rx function: the main process is to allocate a sk_bull socket buffer, write data and set some members, and submit it to the upper layer software for processing with the netif_rx function. For a detailed description, refer to the description in the book.

```
1.  void snull_rx(struct net_device *dev, struct snull_packet *pkt)
2.  {
3.    struct sk_buff *skb;
4.    struct snull_priv *priv = netdev_priv(dev);
5.    skb = dev_alloc_skb(pkt->datalen + 2);
6.    if (!skb)
7.    {
8.      printk("snull rx: low on mem - packet dropped\n");
9.      priv->stats.rx_dropped++;
10.     return;
```

```
11.  }
12.  skb_reserve(skb, 2);
13.      memcpy(skb_put(skb, pkt->datalen), pkt->data, pkt->datalen);
14.  /* Write metadata, and then pass to the receive level */
15.      skb->dev = dev;
16.      skb->protocol = eth_type_trans(skb, dev);
17.      skb->ip_summed = CHECKSUM_UNNECESSARY; /* don't check it */
18.      priv->stats.rx_packets++;
19.      priv->stats.rx_bytes += pkt->datalen;
20.      netif_rx(skb);
21.  return;
22.}
```

There are two ways to receive, based on interrupt or polling, here only retain the use of interrupt receiving, based on the description of the polling method reference book.

If it is a sending interrupt, release the socket buffer after setting the relevant state.

4. snull_ioctl: it will be called when using the ifconfig command, but here to simplify the procedure, it simply returns 0 and does nothing:

```
1.  int snull_ioctl(struct net_device *dev, struct ifreq *rq, int cmd)
2.  {
3.    return 0;
4.  }
```

5. Snull_stats: returns the current status of the device

```
1.  struct net_device_stats *snull_stats(struct net_device *dev)
2.  {
3.    struct snull_priv *priv = netdev_priv(dev);
4.    return &priv->stats;
5.  }
```

6. snull_tx_timeout: send timeout processing, just simulate sending interrupt again and restart the transmission queue

```
1.  void snull_tx_timeout (struct net_device *dev)
2.  {
3.    struct snull_priv *priv = netdev_priv(dev);
4.    printk(KERN_INFO"Transmit timeout at %ld, latency %ld\n", jiffies,jiffies -
       dev->trans_start);
```

```
5.    priv->status = SNULL_TX_INTR;
6.    snull_interrupt(0, dev, NULL);
7.    priv->stats.tx_errors++;
8.    netif_wake_queue(dev);
9.    return;
10.}
```

7. Snull_header: Called before hard_start_xmit, it is used to build a hardware header based on the retrieved source and target hardware addresses.

```
1.  int snull_header(struct sk_buff *skb, struct net_device *dev,
2.            unsigned short type, const void *daddr, const void *saddr,
3.            unsigned int len)
4.  {
5.    struct ethhdr *eth = (struct ethhdr *)skb_push(skb,ETH_HLEN);
6.    eth->h_proto = htons(type);
7.    memcpy(eth->h_source, saddr ? saddr : dev->dev_addr, dev->addr_len);
8.    memcpy(eth->h_dest, daddr ? daddr : dev->dev_addr, dev->addr_len);
9.    eth->h_dest[ETH_ALEN-1] ^= 0x01;
10.   return (dev->hard_header_len);
11.}
```

8. snull_rebuild_header: used to re-establish the hardware header after completing the ARP parsing before transmitting the data packet

```
1.  int snull_rebuild_header(struct sk_buff *skb)
2.  {
3.    struct ethhdr *eth = (struct ethhdr *) skb->data;
4.    struct net_device *dev = skb->dev;
5.
6.    memcpy(eth->h_source, dev->dev_addr, dev->addr_len);
7.    memcpy(eth->h_dest, dev->dev_addr, dev->addr_len);
8.    eth->h_dest[ETH_ALEN-1] ^= 0x01;
9.    return 0;
10.}
```

After implementing these device methods, basically the driver is complete.