# Network drivers

**Team Emertxe** 



# **Topics**



- Linux Network stack and Driver
- Interrupt
- Memory access
- RX/TX
- SKBuffs and NAPI
- Mellanox Adapters Programmer's Reference Manual



# PCI Driver Registration



- int pci\_register\_driver(struct pci\_driver \*drv);
- int pci\_unregister\_driver(struct pci\_driver \*drv);
- struct pci\_driver
- -const char \*name
- -const struct pci\_dev\_id \*id\_table;
- .PCI\_DEVICE(vendor, device);
- .PCI\_DEVICE\_CLASS(dev\_class, dev\_class\_mask);
- -int (\*probe)(pci\_dev, id\_table);
- -void (\*remove)(pci\_dev);



# The 'probe' Function



```
int probe(struct pci dev *d, struct pci dev id *id)
  /* Initialize the PCI Device */
  /* Enable the PCI Device */
  pci_enable_device(d);
  return 0; /* Claimed. Negative for not Claimed */
```



### Network interface Overview

**Linux Network Stack** 

**Network Driver** 

PCI interface (Hardware)

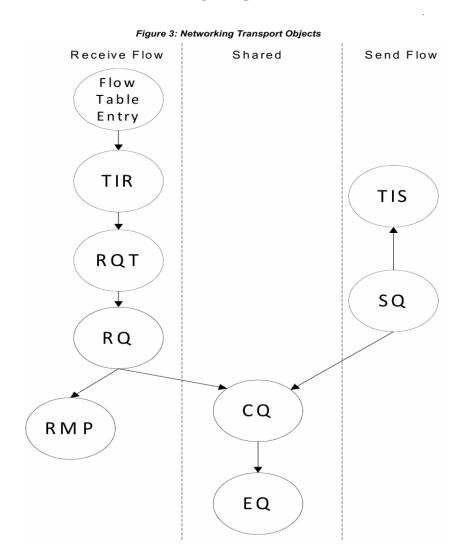
Firmware (Proprietary Software)

ASIC (Proprietary Hardware)

Connectors (SFP)



# Network interface Overview -Firmware



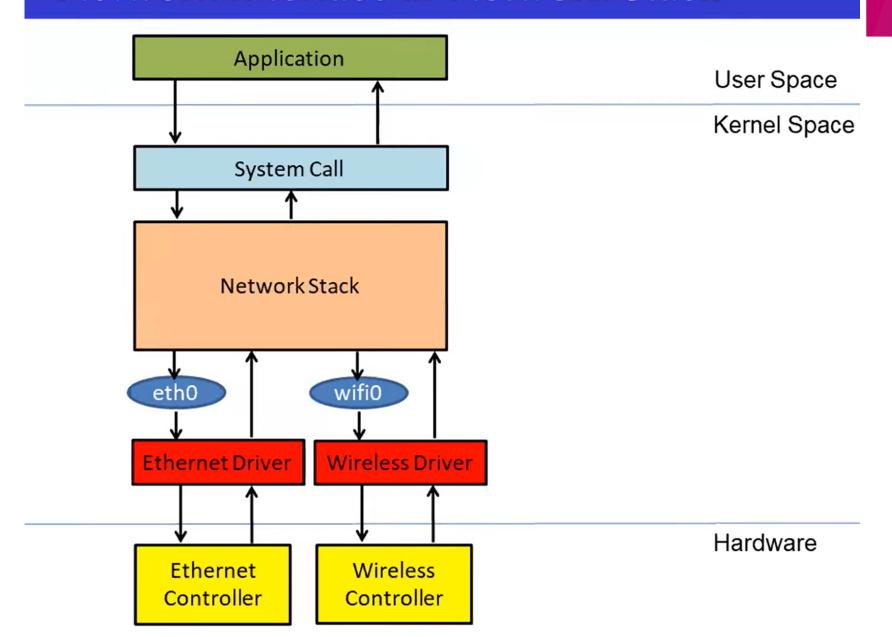




- PCI: Probe, BAR, MSIX
- **IEEE 803 standards**
- Linux kernel concept: IRQ, threads, memory, etc.
- Linux network stack: sk\_buff, net\_device\_ops, ethtool
- Tools: RSS settings, Wireshark, iPerf, etc.



#### Network Interface in Network Stack





#### Device registration with network stack

```
/**
         register netdev
                             - register a network device
         @dev: device to register
*
         Take a completed network device structure and add it to the kernel
         interfaces. A %NETDEV REGISTER message is sent to the netdev notifier
         chain. 0 is returned on success. A negative errno code is returned
         on a failure to set up the device, or if the name is a duplicate.
         This is a wrapper around register netdevice that takes the rtnl semaphore
         and expands the device name if you passed a format string to
         alloc netdev.
*/
int register_netdev(struct net_device *dev)
```



#### Driver registration

```
static int __init dummy_init_module(void)
        int i, err = 0;
        down_write(&pernet_ops_rwsem);
        rtnl_lock();
        err = __rtnl_link_register(&dummy link ops);
        if (err < 0)
                goto out;
        for (i = 0; i < numdummies && !err; i++) {
                err = dummy_init_one();
                cond resched();
        if (err < 0)
                __rtnl_link_unregister(&dummy link ops);
out:
        rtnl_unlock();
        up_write(&pernet_ops_rwsem);
        return err;
```





```
static int init dummy init one (void)
        struct net_device *dev_dummy;
        int err;
        dev_dummy = alloc_netdev(0, "dummy%d", NET_NAME_ENUM, dummy_setup);
        if (!dev dummy)
                return -ENOMEM;
        dev dummy->rtnl link ops = &dummy link ops;
        err = register_netdevice(dev_dummy);
        if (err < 0)
                goto err;
        return 0:
err:
        free_netdev(dev_dummy);
        return err;
```



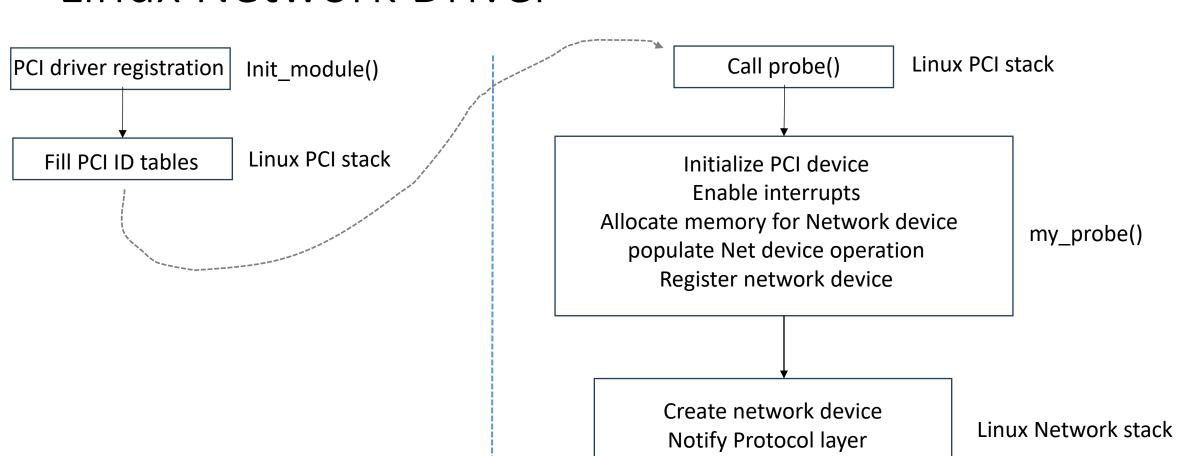
```
static void dummy setup(struct net device *dev)
       ether_setup(dev);
       /* Initialize the device structure. */
       dev->netdev_ops = &dummy_netdev_ops;
       dev->ethtool_ops = &dummy_ethtool_ops;
       dev->needs free netdev = true;
       /* Fill in device structure with ethernet-generic values. */
       dev->flags |= IFF_NOARP;
       dev->flags &= ~IFF_MULTICAST;
       dev->priv_flags |= IFF_LIVE_ADDR_CHANGE | IFF_NO_QUEUE;
       dev->lltx = true:
       dev->features |= NETIF_F_SG | NETIF_F_FRAGLIST;
       dev->features |= NETIF_F_GSO_SOFTWARE;
       dev->features |= NETIF_F_HW_CSUM | NETIF_F_HIGHDMA;
       dev->features |= NETIF_F_GSO_ENCAP_ALL;
       dev->hw features |= dev->features;
       dev->hw_enc_features |= dev->features;
       eth_hw_addr_random(dev);
       dev->min_mtu = 0;
       dev -> max mtu = 0;
```



Device registration with network stack

```
static const struct net_device_ops dummy_netdev_ops = {
                                = dummy_dev_init,
        .ndo_init
        .ndo_start_xmit
                                = dummy_xmit,
        .ndo_validate_addr
                                = eth_validate_addr,
        .ndo_set_rx_mode
                                = set_multicast_list,
        .ndo_set_mac_address
                                = eth_mac_addr,
                                = dummy_get_stats64,
        .ndo_get_stats64
        .ndo_change_carrier
                                = dummy_change_carrier,
};
```

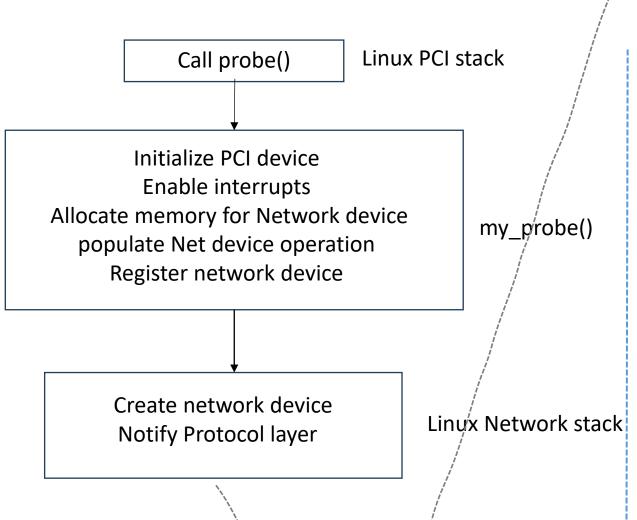






#### Network device operations

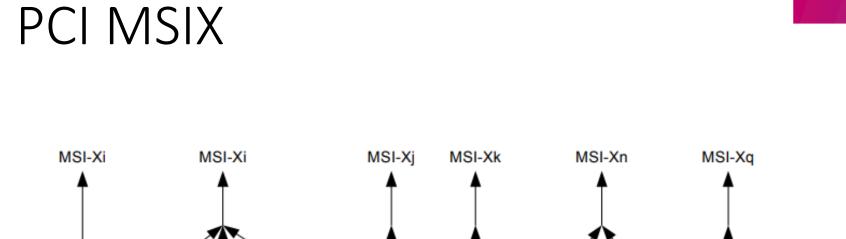
```
824
      static const struct net_device_ops e1000_netdev_ops = {
825
              .ndo open
                                      = e1000 open,
826
              .ndo stop
                                      = e1000 close,
              .ndo start xmit
                                      = e1000 xmit frame,
827
              .ndo set rx mode
                                      = e1000 set rx mode,
828
              .ndo_set_mac_address
                                      = e1000 set_mac,
829
              .ndo tx timeout
                                      = e1000 tx timeout,
830
              .ndo change mtu
                                      = e1000 change mtu,
831
              .ndo eth ioctl
                                      = e1000 ioctl,
832
833
              .ndo validate addr
                                      = eth validate addr,
                                      = e1000_vlan_rx_add_vid,
              .ndo_vlan_rx_add_vid
834
                                      = e1000_vlan_rx_kill_vid,
              .ndo vlan rx kill vid
835
      #ifdef CONFIG_NET_POLL_CONTROLLER
836
837
              .ndo poll controller
                                      = e1000 netpoll,
      #endif
838
839
              .ndo fix features
                                      = e1000 fix features,
                                      = e1000 set features,
840
              .ndo set features
      };
841
```

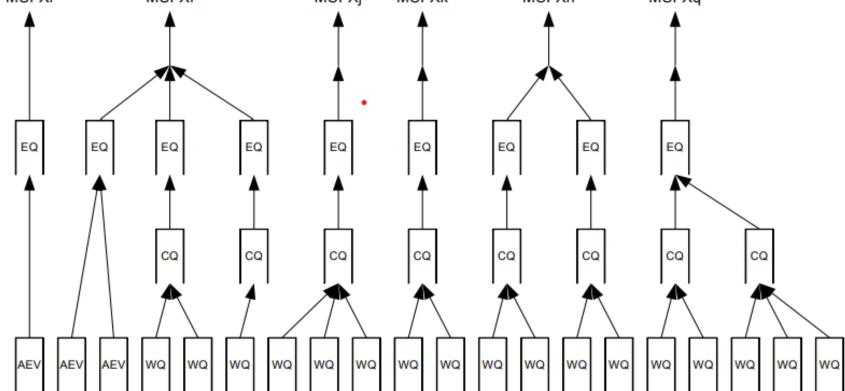


Allocate transmit descriptors
Allocate receive descriptors
Power up Net adapter
Allocate an interrupt
Net Polling napi\_enable()
netif\_start\_queue - allow transmit

Driver my\_open or e1000\_open









# Linux Network Driver interrupt

- request\_irq()
- request\_threaded\_irq()
- free\_irq()
- disable\_irq()
- enable\_irq()







interrupt

```
static int e1000_request_irq(struct e1000_adapter *adapter)
250
251
               struct net_device *netdev = adapter->netdev;
252
               irq_handler_t handler = e1000_intr;
253
               int irq flags = IRQF SHARED;
254
255
               int err;
256
               err = request_irq(adapter->pdev->irq, handler, irq_flags, netdev->name,
257
258
                                 netdev);
259
               if (err) {
                       e_err(probe, "Unable to allocate interrupt Error: %d\n", err);
260
261
262
263
               return err;
264
265
```





#### **Memory Access**

- What is the Memory type [ Port mapped or Memory mapped ]
- Size of memory access [Register read(bytes or words) or DMA]

Port mapped (IO mapped)
Include asm/io.h
readb(), writeb(), readw(), writew() & readl(), writel()

Memory mapped
Normal memory access

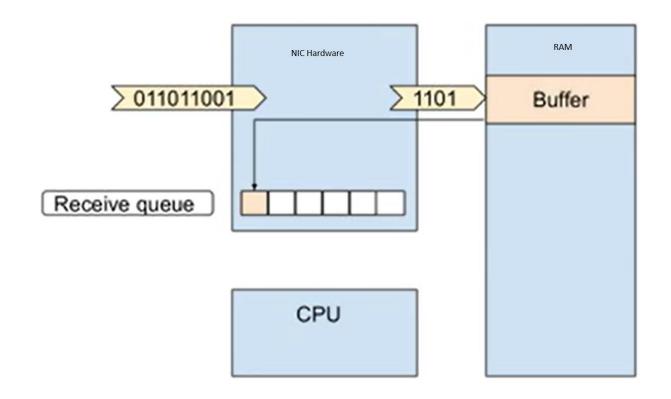


**Memory Access** 

```
Port mapped (IO mapped)
Include asm/io.h
readb(), writeb(), readw(), writew() & readl(), writel()
```



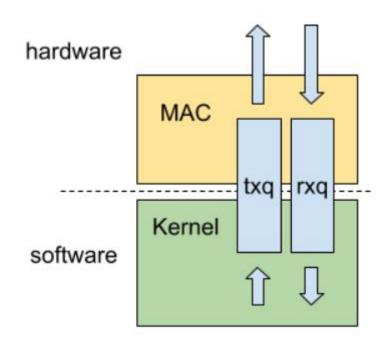




- ► The MAC received data and writes it to RAM using DMA
- A descriptor is created
- It's address is put in a queue



RX and TX





```
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
```

```
/* allocate transmit descriptors */
err = e1000_setup_all_tx_resources(adapter);
if (err)
        goto err setup tx;
/* allocate receive descriptors */
err = e1000_setup_all_rx_resources(adapter);
if (err)
        goto err_setup_rx;
```



#### Struct skbuff

```
* DOC: Basic sk buff geometry
* struct sk buff itself is a metadata structure and does not hold any pac
* data. All the data is held in associated buffers.
* &sk buff.head points to the main "head" buffer. The head buffer is divid
* into two parts:
  - data buffer, containing headers and sometimes payload;
    this is the part of the skb operated on by the common helpers
    such as skb put() or skb pull();
  - shared info (struct skb shared info) which holds an array of pointer.
    to read-only data in the (page, offset, length) format.
* Optionally &skb shared info.frag list may point to another skb.
* Basic diagram may look like this::
    headroom | data | tailroom | skb shared info
                                 + [page frag]
                                 + [page frag]
                                 + [page frag]
                                 + [page frag]
                                 + frag list
```



#### **NAPI**

```
napi enable - enable NAPI scheduling
        @n: NAPI context
 * Resume NAPI from being scheduled on this context.
 * Must be paired with napi disable.
 */
void napi enable(struct napi struct *n)
 * netif napi add() - initialize a NAPI context
 * @dev: network device
 * @napi: NAPI context
 * @poll: polling function
 * netif napi add() must be used to initialize a NAPI context prior to cal
 * *any* of the other NAPI-related functions.
static inline void
netif_napi_add(struct net_device *dev, struct napi_struct *napi,
               int (*poll)(struct napi_struct *, int))
        netif_napi_add_weight(dev, napi, poll, NAPI_POLL_WEIGHT);
```

