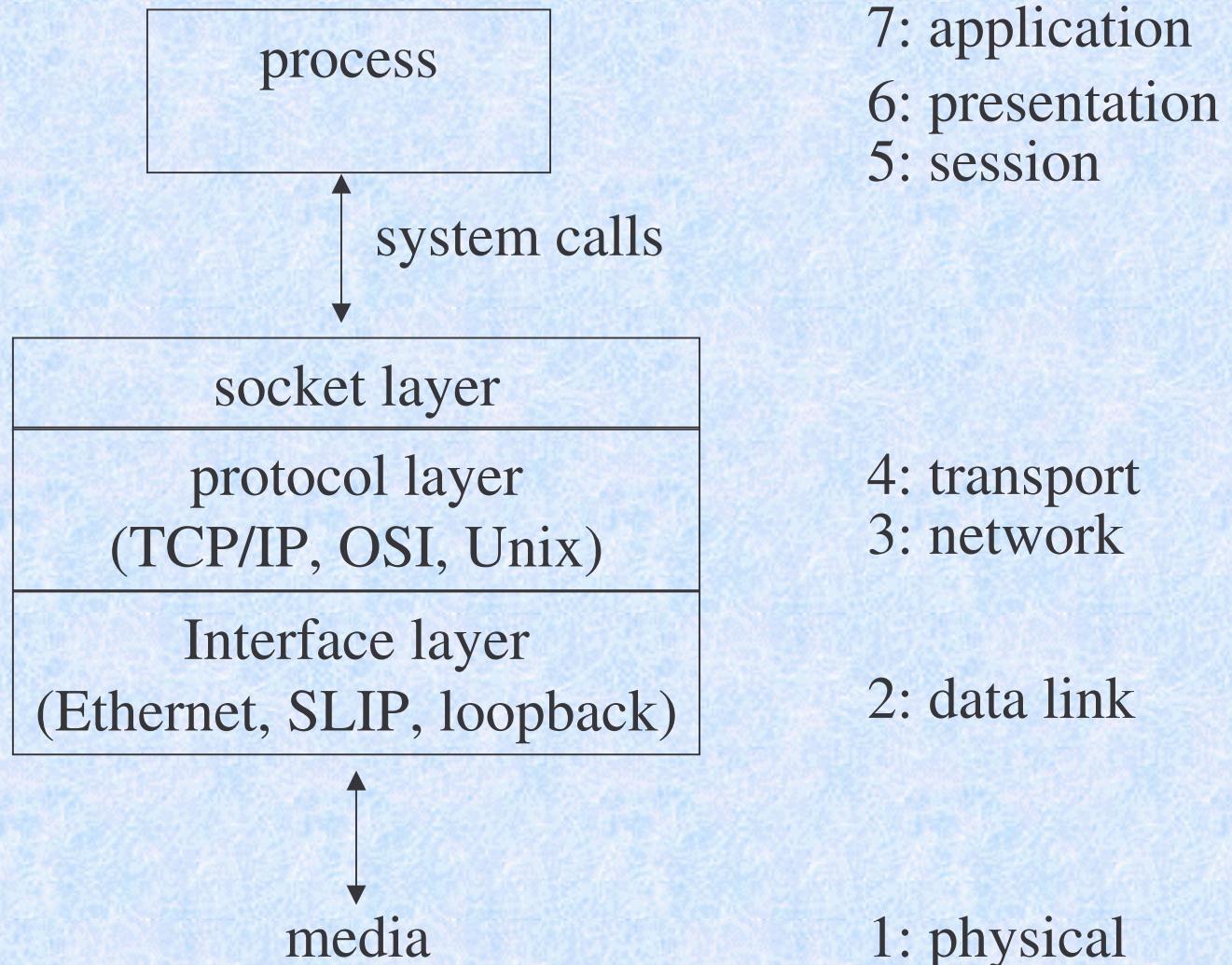


# Kernel Code Directory Structure

Under /usr/src/sys

- **dev/fxp – Intel fxp network card driver**
- i386 – intel specific
- kern – general kernel stuff that does belong else where
- **net – general networking stuff**
- netiso – OSI protocols
- **netinet – TCP/IP protocols**
- pci – general pci support functions
- sys – header files
- nfs – NFS

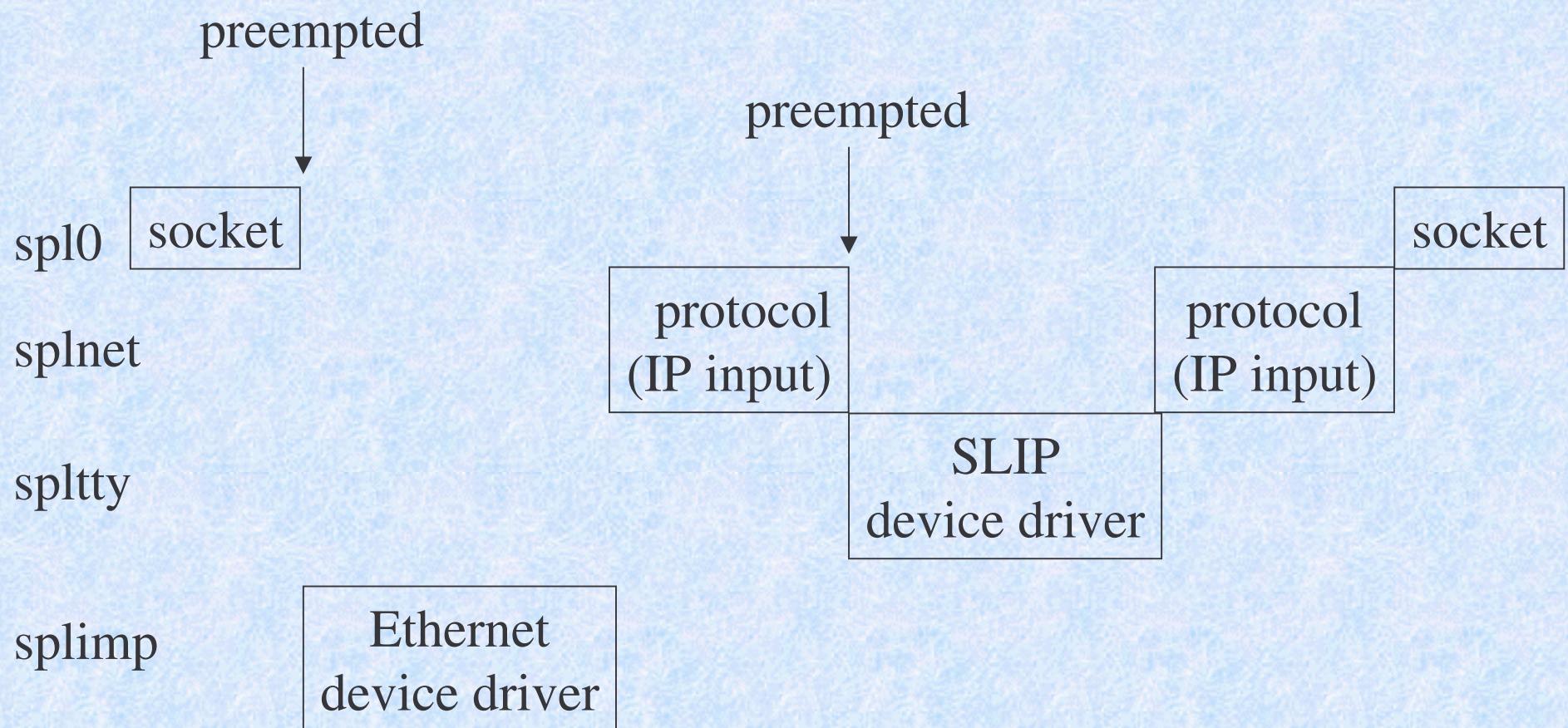
# Network Code Organization



# Set Priority Level Commands

spl0	normal operating mode, nothing blocked (lowest priority)
splsoftclock	low-priority clock processing
splnet	network protocol processing
spltty	terminal I/O
splbio	disk and tape I/O
splimp	network device I/O
splclock	high priority clock processing
splhigh	all interrupts blocked (highest priority)
splx(s)	return to level s

# SPL Examples



# packet queue data structure

mbuf 1

m_next
m_nextpkt
m_len
m_data
m_type
m_flags
m_pkthdr.len
p_pkthdr.rcvif
100 bytes of data

mbuf 2

m_next
m_nextpkt
m_len
m_data
m_type
m_flags
4 bytes of data
104 bytes (unused)

100

MT\_DATA

M\_PKTHDR

104

ptr to ifnet

NULL

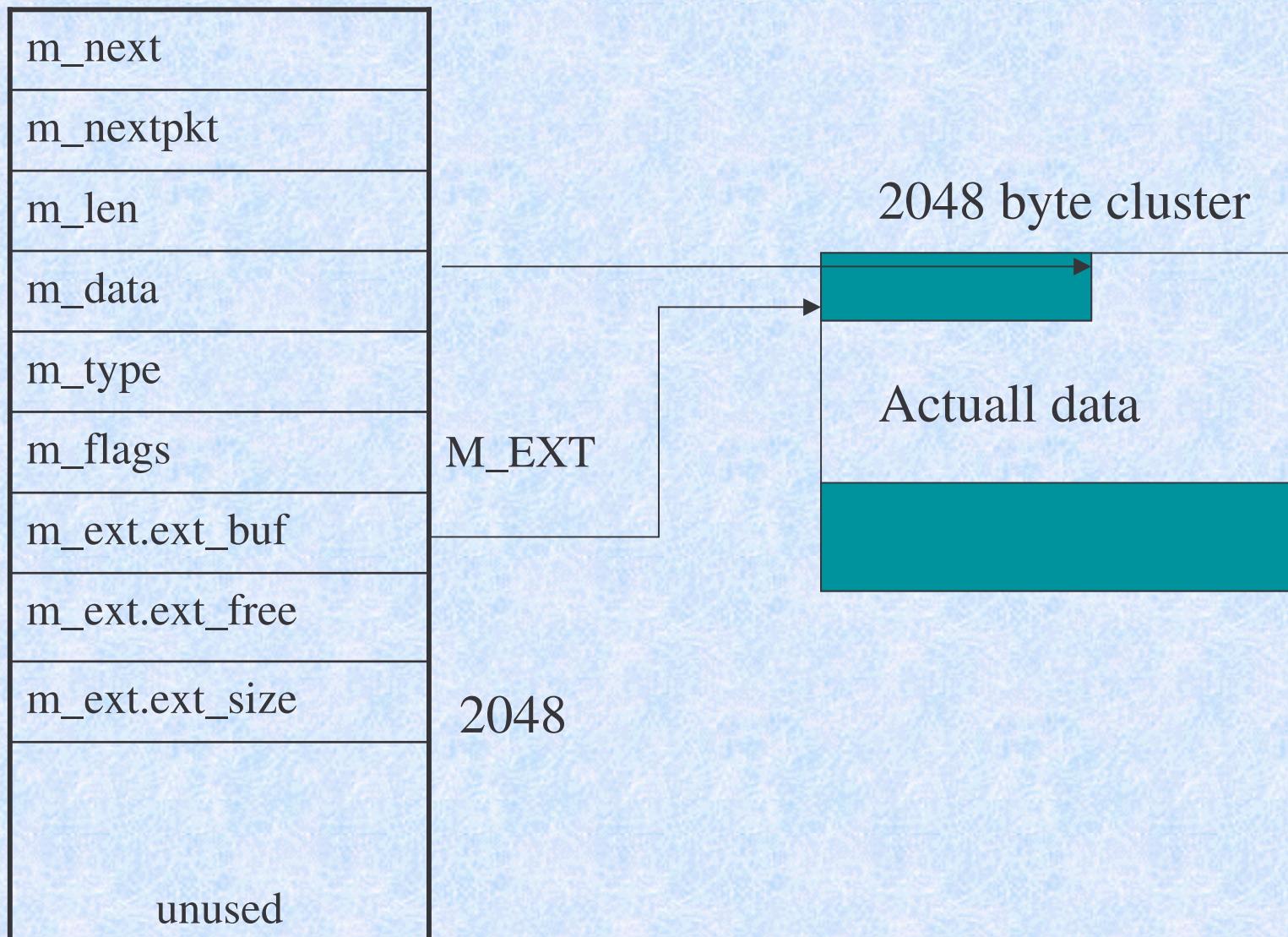
NULL

4

MT\_DATA

0

# External mbuf structure



# Mbuf functions

m_adj(*m, len)	remove len bytes from start of m if len>0 or abs(len) bytes from end of m if len<0
m_cat(*m, n)	concatenate mbuf n to mbuf m
mbuf *m_copy(*m, offset, len)	copy len bytes from offset in mbuf m and return new mbuf
m_free, m_freem(*m)	free one mbuf or mbuf chain
mbuf *m_get(wait, type)	get new mbuf

by \*m I mean struct mbuf \*

# Mbuf functions cont'

<code>m_pullup(*m, len)</code>	arrange data contiguously in memory
<code>M_PREPEND(*m, len, wait)</code>	prepend len bytes to start of m
<code>mtod(*m, type)</code>	cast the mbuf * to new type

# ifnet structure

- Each interface has a struct ifnet holding its data.
- All struct ifnets are connected in a linked list.
- ifnet functions:
  - if\_init – initialize the interface
  - if\_output – queue outgoing packets for transmission.
  - if\_start – initiate transmission of packets.
  - if\_ioctl – the interface ioctl function
- Also hold interface name, flags, statistics, etc...

# Incoming Ethernet Packet Handling

- Driver interrupt routine (fxp\_intr) getting the packet from the hardware into an mbuf chain. calling ethernet driver
- Ethernet driver validate the packet. Look at layer 3 type, put packet on appropriate queue, and initiate appropriate software interrupt
- Layer 3 interrupt routine get packet from queue. Do the layer 3 job.

# Incoming Packet – Network Card Driver

- Uses DMA to get incoming packet data into kernel memory.
- `fpx_intr`:
  - Uses preallocated mbuf pool to build an mbuf chain holding the packet.
  - Doing checksum validation.
  - Calling the input function of the right interface.

`src/sys/dev/fpx/if_fpx.c`

# Incoming Packet - Ethernet

```
void  
ether_input(struct ifnet *ifp, struct ether_header *eh, struct mbuf *m)  
    m->m_pkthdr.rcvif = ifp;  
  
    eh = mtod(m, struct ether_header *);  
  
    m->m_data += sizeof(struct ether_header);  
  
    m->m_len -= sizeof(struct ether_header);  
  
    m->m_pkthdr.len = m->m_len;  
  
    ether_demux(ifp, eh, m);
```

src/sys/net/if\_ETHERSUBR.C

# Incoming Packet – Ethernet (cont')

```
void ether_demux(ifp, eh, m)
    ether_type = ntohs(eh>ether_type);
    switch (ether_type) {
        case ETHERTYPE_IP:
            schednetisr(NETISR_IP);
            inq = &ipintrq;
            break;
        .
        .
        .
        (void) IF_HANDOFF(inq, m, NULL);
```

# Incoming Packet - IP

ipintr – the IP interrupt routine

```
while(1) {  
    s = splimp();  
    IF_DEQUEUE(&ipintrq, m);  
    splx(s);  
    if (m == 0)  
        return;  
    ip_input(m);  
}
```

[src/sys/netinet/ip\\_input.c](src/sys/netinet/ip_input.c)

# Incoming Packet – IP (cont’)

```
void ip_input(struct mbuf *m)
```

IPFIREWALL and DUMMYNET hooks

```
if (m->m_len < sizeof (struct ip) &&
    (m = m_pullup(m, sizeof (struct ip))) == 0) {
```

```
    ipstat.ips_toosmall++;
    return;
```

```
ip = mtod(m, struct ip *);
```

validation checks on IP header

# Incoming Packet – IP (cont’)

IP routing handling

IP fragments reassembly

m\_adj to remove the IP header

inetsw[ip\_protox[ip->ip\_p]].pr\_input)(m, off, nh);\*)

# Outgoing packet - UDP

```
static int udp_output(struct inpcb *inp, struct mbuf *m
, struct sockaddr *addr struct mbuf *control, struct thread *td) {

    M_PREPEND(m, sizeof(struct udppiphdr) + max_linkhdr, M_DONTWAIT);
    ui = mtod(m, struct udppip_hdr *);

    Filling UDP header fields, pseudo IP header fields

    error = ip_output(m, inp->inp_options, NULL, ipflags,
                      inp->inp_m_options, inp)

}
```

src/sys/netinet/udp\_usrreq.c

# Outgoing packet - IP

```
int ip_output(struct mbuf *m, struct mbuf *opt, struct route *ro,
              int flags, struct ip_moptions *imo, struct inpcb *inp) {

    ip = mtod(m, struct ip *);

    Multicast handling

    Routing handling – finding out which interface we go out from
    Cut to pieces smaller than the interfaces's MTU

    For (; m; m = m0) {

        m0 = m->m_nextpkt
        m->m_nextpkt = 0;
        error = *ifp->if_output(ifp, m, dst, ro->ro_rt);
    }
}
```

src/sys/netinet/ip\_output.c

# Outgoing packet - Ethernet

```
int ether_output(struct ifnet *ifp, struct mbuf *m,
                 struct sockaddr *dst, struct rtentry *rt0) {

    switch (dst->sa_family) {

        case AF_INET:

            error = arpresolve(ifp, rt0, m, dst, edst);

            ...
    }

    M_PREPEND(m, ETHER_HDR_LEN, M_DONTWAIT);

    eh = mtod(m, struct ether_header *);

    IFQ_HANDOFF(ifp, m, error);

}
```

`src/sys/net/if_ETHERSUBR.C`

```
#define IFQ_HANDOFF(ifp, m, err)
do {
    int len;
    short mflags;
    len = (m)->m_pkthdr.len;
    mflags = (m)->m_flags;
    IFQ_ENQUEUE(&(ifp)->if_snd, m, err);
    if ((err) == 0) {
        (ifp)->if_obytes += len;
        if (mflags & M_MCAST)
            (ifp)->if_omcasts++;
        if (((ifp)->if_flags & IFF_OACTIVE) == 0)
            if_start(ifp);
    }
} while (0)
```

src/sys/net/if\_var.h

# Outgoing Packet – Network Card Driver

- fxp\_start found in src/sys/dev/fxp/if\_fxp.c
- While there are packets in the send queue do:
  - get a packet off the queue
  - hand it to the hardware to transmit

# References

- The Design and Implementation of the 4.4  
BSD Operating System  
McKusick, Bostic, Karels, Quarterman
- TCP/IP Illustrated Volume 2  
Stevens and Wright
- FreeBSD source code