



1.

Linux Kernel : 2.4.13

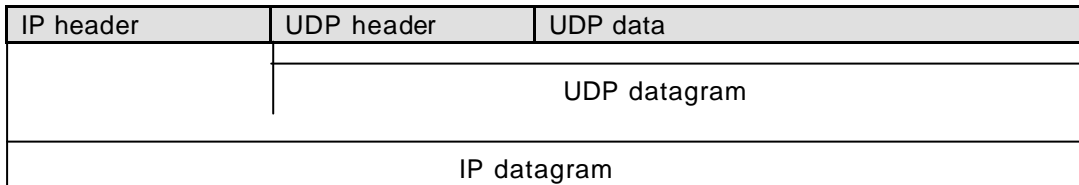
Processor: Intel x86 family

Ethernet Device : 3COM 3c509

## 2. Base Knowledge

UDP(User Datagram Protocol) ?

UDP transport layer      IP layer      OSI 7 layer



가

4가



checksum

).

(

(flow control)

가

(ex> Optic cable)  
layer      TCP

- 가 ( 가 ).
- (
- half duplexing ).
- 가
-



Local Network

Device ↔ Kernel Area ↔ User Area

NIC

system call )

가

가

( )가

가

가

가

가

가

가

가 system call

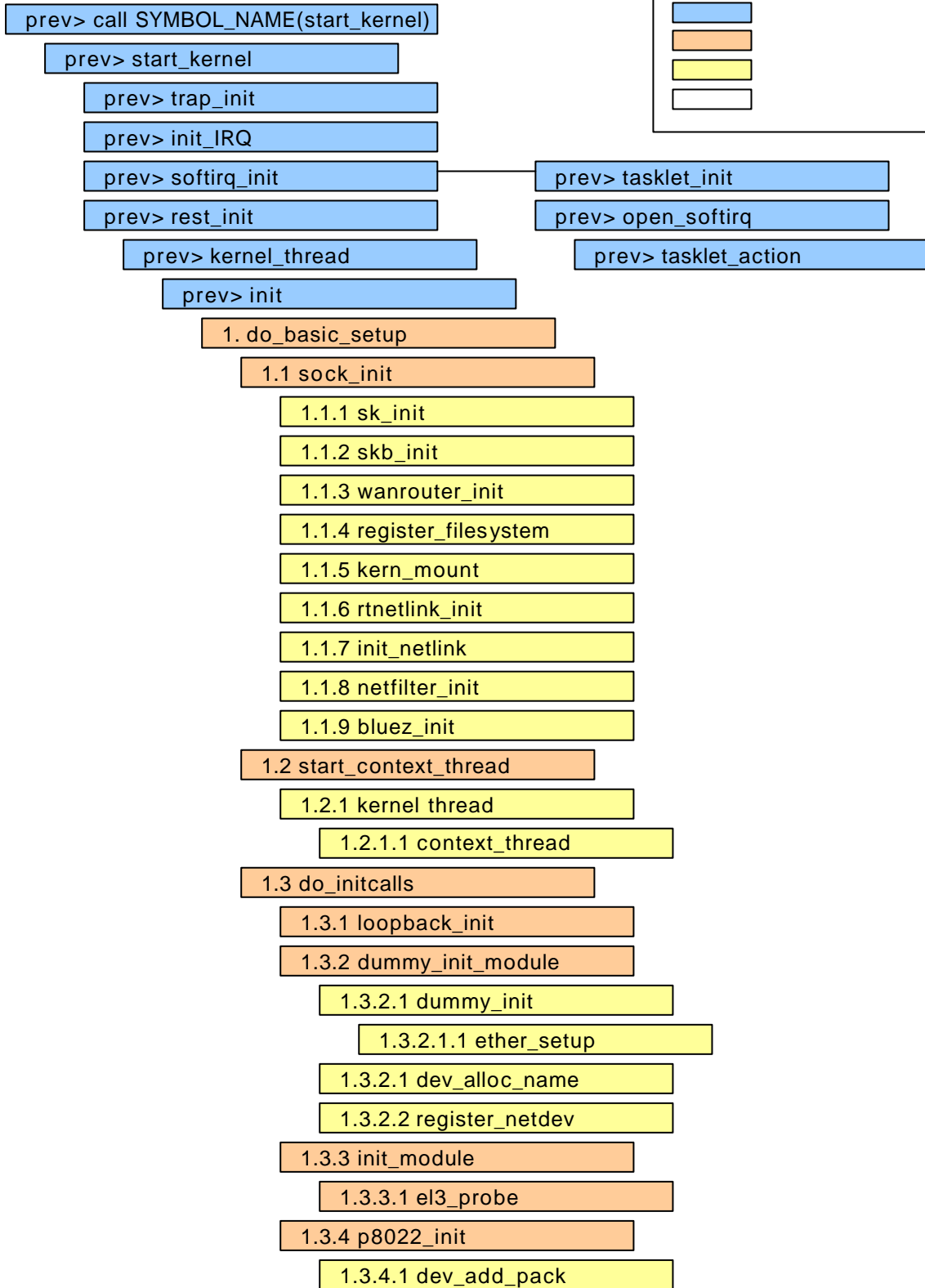
Comment

가



### 3. Network Initialization

Function Call Tree



### 1.3.5 snap\_init

1.3.5.1 register\_8022\_client

### 1.3.6 inet\_init

1.3.6.1 sock\_register

1.3.6.2 inet\_add\_protocol

1.3.6.3 arp\_init

1.3.6.4 ip\_init

1.3.6.4.1 dev\_add\_pack

1.3.6.4.2 ip\_rt\_init

1.3.6.4.3 proc\_net\_create

1.3.6.5 tcp\_v4\_init

1.3.6.5.1 init\_waitqueue\_head

1.3.6.6 tcp\_init

1.3.6.7 icmp\_init

1.3.6.7.1 init\_waitqueue\_head

1.3.6.8 ipip\_init

1.3.6.9 ipgre\_init

1.3.6.10 ip\_mr\_init

1.3.6.11 proc\_net\_create

1.3.6.11.1 create\_proc\_info\_entry

1.3.6.11.1.1 create\_proc\_entry

1.3.6.11.1.2 ip\_mc\_procinfo

### 1.3.7 af\_unix\_init

1.3.7.1 sock\_register

1.3.7.2 create\_proc\_read\_entry

1.3.7.3 unix\_sysctl\_register

### 1.3.8 netlink\_proto\_init

1.3.8.1 sock\_register

1.3.8.2 create\_proc\_read\_entry

### 1.3.9 packet\_init

1.3.9.1 sock\_register

1.3.9.2 register\_netdevice\_notifier

1.3.9.3 create\_proc\_read\_entry

### 1.3.10 atalk\_init

1.3.10.1 register\_netdevice\_notifier

1.3.10.2 aarp\_proto\_init



### Tree Analysis

#### Comment

#### Kernel

32 x86 arch/i386/kernel/head.S  
processor , page table directory/page table , paging , SMP

arch/i386/kernel/head.S 가

```
#ifdef CONFIG_SMP
    movb ready, %cl
    cmpb $1,%cl
    je 1f                # the first CPU calls start_kernel
                        # all other CPUs call initialize_secondary
    call SYMBOL_NAME(initialize_secondary)
    jmp L6
1:
#endif
call SYMBOL_NAME(start_kernel)
```

CPU start\_kernel , SMP CPU

initialize\_secondary  
initialize\_secondary arch/i386/kernel/smpboot.c  
stack pointer instruction pointer

init/main.c start\_kernel [asmlinkage void \_\_init start\_kernel(void) in  
init/main.c] 가 architecture setup, trap[trap\_init() in  
arch/i386/kernel/traps.c], IRQ, scheduler, softirq, time, console, memory, fork,  
cache, signal, ipc, smp rest\_init [static void rest\_init(void)  
in init/main.c]

#### Help Comment ? trap init

trap\_init() divide\_error, overflow, invalid\_op, device\_not\_available, stack\_segment,  
page\_fault system\_call interrupt vector  
system\_call index 80 (SYSCALL\_VECTOR)

softirq\_init() 32 tasklet softirq vector TASKLET\_SOFTIRQ  
HI\_SOFTIRQ softirq action tasklet\_action

#### Help Comment ? softirq & tasklet

softirq	tasklet	Linux Kernel	2.4	가	network stack
		Bottom Half			
		Bottom Half	softirq가		
Bottom Half		CPU가		CPU	
		softirq		CPU	SMP
					32
tasklet			32 가	CPU	
tasklet		가			32

rest\_init kernel\_thread init [static int init(void \* unused) in init/main.c]  
, init do\_basic\_setup [static void \_\_init do\_basic\_setup(void) in



## Network Initialization

init/main.c]

sock\_init [void \_\_init sock\_init(void) in net/socket.c]

1. do\_basic\_setup in init/main.c
  - architecture dependant CPU PCI, SBUS, ISAPNP
  - device
  - sock\_init(), start\_context\_thread(), do\_initcalls()
  - process context가 start\_context\_thread() 가 가
  - protocol do\_initcalls()
  - IRDA, PCMCIA

Comment

### Network

1.1 sock\_init

in net/socket.c

<pre>for (l = 0; l &lt; NPROTO; l++)     net_families[l] = NULL; sk_init(); skb_init(); wanrouter_init(); register_filesystem(&amp;sock_fs_type); sock_mnt = kern_mount(&amp;sock_fs_type); rtnetlink_init(); init_netlink(); retfilter_init(); bluez_init();</pre>	<p>← protocol family table NULL (kernel protocol )</p> <p>← sock SLAB cache</p> <p>← skbuff SLAB cache</p> <p>← WAN router</p> <p>← protocol module</p> <p>← network link</p> <p>← network filter</p> <p>← Bluetooth</p>
---	--

1.2 start\_context\_thread

in kernel/context.c

```
context_thread (context_thread kernel_thread
)
(wait_for_completion(&startup)).
```

1.3 do\_initcalls

in init/main.c

CPU UID cache, Memory(CPU cache/swap/shared), File System, block/character device  
 Network Device  
 IDE, CD-ROM, PCI  
 Network  
 debug System.map

1.3.1 loopback\_init

in drivers/net/loopback.c

loopback device device list 가

Help Comment ? loopback device

loopback

network device



## Network Initialization

가 loopback

1.3.2 dummy\_init\_module in drivers/net/dummy.c  
 dummy device dummy\_init  
 dev\_dummy.init = dummy\_init; ← device  
 SET\_MODULE\_OWNER(&dev\_dummy); ← device  
 err=dev\_alloc\_name(&dev\_dummy, ← dummy%d device  
 "dummy%d ");  
 err = register\_netdev(&dev\_dummy); ← dev\_dummy net\_device list

Help Comment – dummy device

dummy device loopback device 가 device  
 dummy 가  
 가 standard output 가  
 dummy device

1.3.3 init\_module in drivers/net/3c509.c

가 3c509 가  
 drivers/net/3c509.c init\_module 가  
 while (el3\_probe(0) == 0) { ← device  
 el3\_root\_dev->irq = irq[el3\_cards]; ← IRQ  
 el3\_root\_dev->if\_port = xcvr[el3\_cards]; ← Interface port number  
 }

1.3.3.1 el3\_probe in drivers/net/3c509.c  
 device가 bus(EISA/MCA ) 가  
 device

1.3.4 p8022\_init in net/802/p8022.c  
 ETH\_P\_802\_2 (dev\_add\_pack)

1.3.5 snap\_init in net/802/psnap.c  
 SNAP layer IP SNAP layer

1.3.6 inet\_init in net/ipv4/af\_inet.c  
 TCP/IP

(void) sock\_register(&inet\_family\_ops); ← SOCKET module inet\_family\_ops  
 link  
 for (p=inet\_protocol\_base; p != NULL; ) { ← protocol  
 inet\_add\_protocol(p);  
 }  
 arp\_init(); ← ARP (2 )  
 ip\_init(); ← IP  
 tcp\_v4\_init(); ← TCP v4  
 tcp\_init(); ← TCP  
 icmp\_init(); ← ICMP



## Network Initialization

```


    ipip_init();
    ipgre_init();
    ip_mr_init();
    proc_net_create ( "... ", 0, ..._get_info);
    proc_net_create ( "... ", 0, ..._get_info);


```

← IPIP  
 ← IPGRE  
 ← IP multicast router  
 ← raw, netstat, smp, sockstat, tcp, udp  
 /proc entry

### 1.3.6.6. ip\_init

in net/ipv4/ip\_output.c

```


    dev_add_pack(&ip_packet_type);
    ip_rt_init();
    proc_net_create( "igmp ", 0,
    ip_mc_procinfo);


```

← ip packet  
 ← ip routing table  
 ← IP multicast  
 igmp  
 /proc entry

### 1.3.7 af\_unix\_init

in net/unix/af\_unix.c

BSD Unix socket  
 unix\_family\_ops SOCKET module

"net/unix" proc entry

### 1.3.8 netlink\_proto\_init

in net/netlink/af\_netlink.c

netlink protocol  
 netlink\_family\_ops SOCKET module

"net/netlink" proc entry

### 1.3.9 packet\_init

in net/packet/af\_packet.c

raw packet socket  
 packet\_family\_ops SOCKET module  
 "net/packet" proc entry  
 packet\_netdev\_notifier notifier\_chain

### 1.3.10 atalk\_init

in net/appletalk/ddp.c

Appletalk DDP protocol for Ethernet

```


    sock_register(&atalk_family_ops);

    dev_add_pack(Italk_packet_type);
    dev_add_pack(ppptalk_packet_type);
    register_netdevice_notifier(&ddp_notifier);
    aarp_proto_init();
    proc_net_create( "... ", 0, ..._get_info);


```

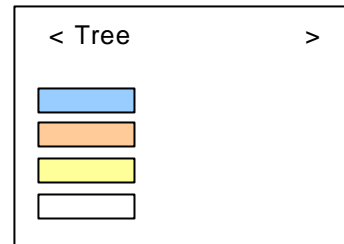
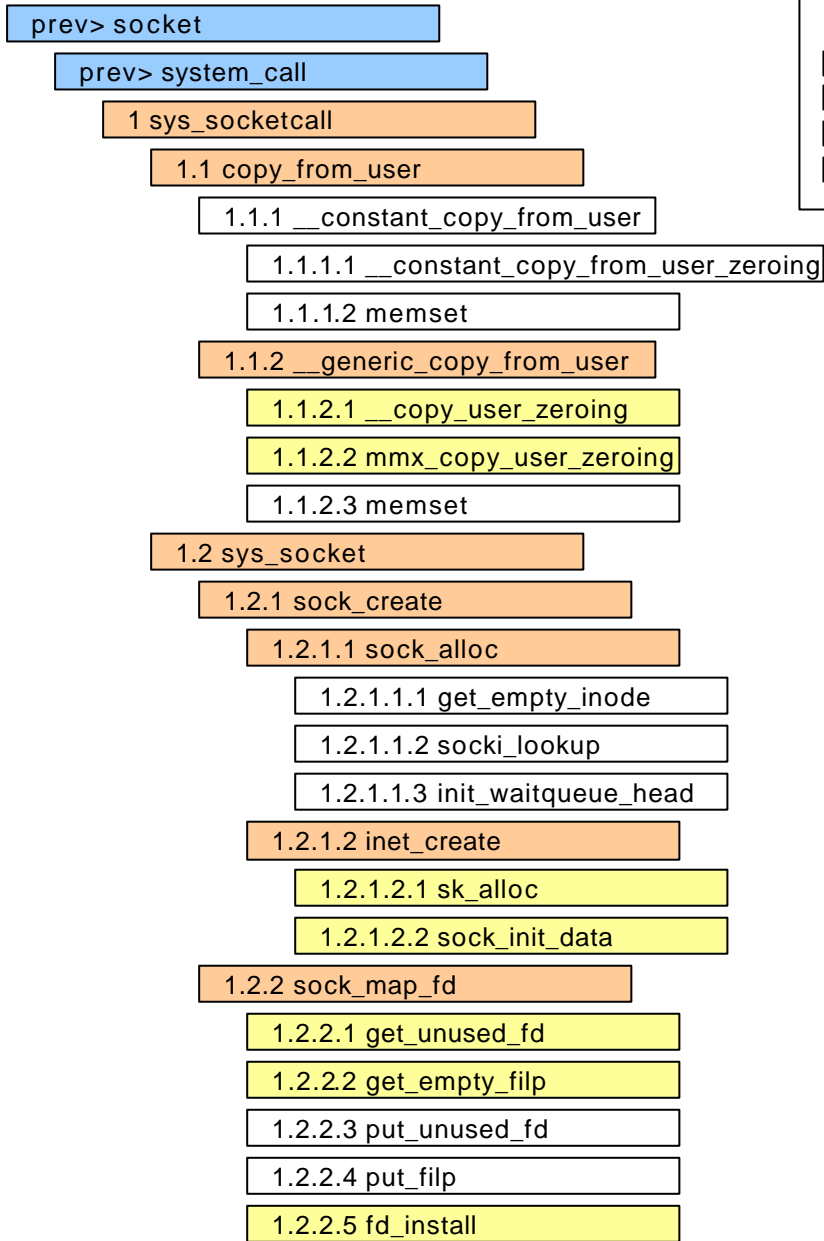
← atalk\_family\_ops SOCKET module  
 ← Italk/ppptalk\_packet\_type  
 ddp\_notifier notifier\_chain  
 ← appletalk ARP  
 ← "appletalk", "atalk\_route",  
 "atalk\_iface" proc entry





### 4. Socket Creation

#### Function Call Tree





### Tree Analysis

```

socket [socket(int domain, int type, int protocol)]
system call (interrupt vector index 0x80, "#define SYSCALL_VECTOR
0x80" in include/asm-i386/hw_irq.h) . ( UDP
domain AF_INET, type SOCK_DGRAM, protocol 0(default) )
) system_call sys_socketcall 가
    
```

1. sys\_socketcall in net/socket.c  
 → asmlinkage long sys\_socketcall(int call, unsigned long \*args)

```

if (copy_from_user(a, args, nargs[call]))
    return -EFAULT;
switch(call)
{
    case SYS_SOCKET:
        err = sys_socket(a0, a1, a[2]);
        break;
    case SYS_BIND:
        ...
    case SYS_CONNECT:
        ...
    ...
}
return err;
    
```

← User (args) Kernel (a)  
 nargs[call]  
 ← call action  
 ←  
 ←  
 ←  
 ← action (LISTEN, ACCEPT, SEND, SENDTO, RECV, RECVFROM )

call SYS\_SOCKET sys\_socket  
 가 long array a family, type,  
 protocol 가 sys\_socket

### Help Comment ? system call table

System call table arch/i386/kernel/entry.S



1.1 copy\_from\_user in include/asm-i386/uaccess.h

User Memory Area Kernel Memory Area  
 copy\_from\_user(to, from, n), from to n(size)  
 n built-in  
 \_\_constant\_copy\_from\_user, \_\_generic\_copy\_from\_user  
 \_\_constant\_copy\_from\_user inline

1.1.2 \_\_generic\_copy\_from\_user in arch/i386/lib/usercopy.c  
 → unsigned long \_\_generic\_copy\_from\_user(void \*to, const void \*from, unsigned long n)

inline AMD 3DNow  
 3DNow \_\_copy\_user\_zeroing  
 mmx\_copy\_user\_zeroing 3DNow

### Help Comment ? 3DNow

3DNow AMD 가 SIMD floating point number  
 가



floating point					3D
3DNow			가	Intel AMD	SSE SSE2 가 (Athlon XP) Hammer
Intel SSE2	SSE		AMD		
		3DNow			

1.2 sys\_socket in net/socket.c  
 → asmlinkage long sys\_socket(int family, int type, int protocol)  
 family, type, protocol socket struct  
 file descriptor 가

1.2.1 sock\_create in net/socket.c  
 → int sock\_create(int family, int type, int protocol, struct socket \*\*res)

```

struct socket *sock; ← socket
if (family < 0 || family >= NPROTO) ← family type
    return ?EAFNOSUPPORT;
if (type < 0 || type >= SOCK_MAX)
    return ?EINVAL;
if (net_families[family] == NULL) { ← family가
    l = -EAFNOSUPPORT;
    goto out;
}
if (!(sock = sock_alloc())) { ←
    ...; goto out;
}
if ((l = net_families[family] ->create(sock, ← net_families family
protocol) < 0) { net_proto_family create
    ...; goto out;
} inet_family_ops가 (
    가 ). , inet_create
    AF_INET, SOCK_DGRAM, default
    
```

1.2.1.1 sock\_alloc in net/socket.c  
 inode socket bind socket  
 inode가 NULL

1.2.1.2 inet\_create in net/ipv4/af\_inet.c  
 → static int inet\_create(struct socket \*sock, int protocol)  
 INET socket  
 sock [struct sock \*sk\_alloc(int family, int priority, int zero\_it) in  
 net/core/socket.c] type/protocol pair

1.2.2 sock\_map\_fd in net/socket.c

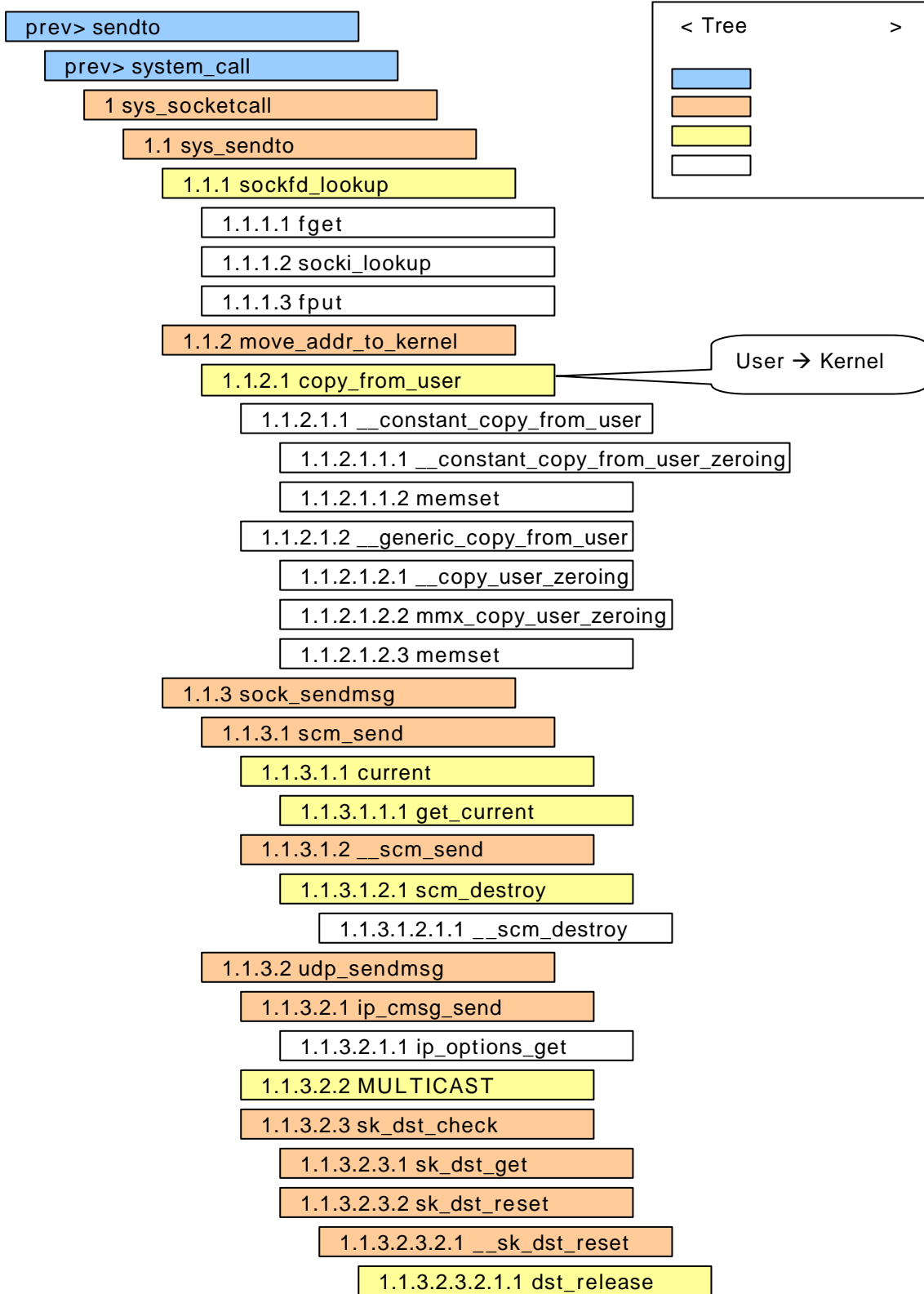
```

fd = get_unused_fd (); ← file descriptor
struct file *file = get_empty_filp (); ← file pointer
...
...
fd_install(fd, file); ← file fd array install
return fd; ← fd
    
```



## 5. Data Sending

### Function Call Tree







### Tree Analysis

sendto(sockfd, buff, len, dst\_addr, addrlen) system call  
 (interrupt vector index 0x80, "#define SYSCALL\_VECTOR 0x80" in include/asm-i386/hw\_irq.h) . system call interrupt handler sys\_socketcall  
 가 UDP .

1. sys\_socketcall in net/socket.c  
 → asmlinkage lon sys\_socketcall(int call, unsigned long \*args)

```
if (copy_from_user(a, args, nargs[call]))
    return ?EFAULT;
switch(call)
{
    case SYS_SOCKET:
        ...
    case SYS_SENDTO:
        err = sys_sendto(a0,(void *)a1,
            a[2], a[3], (struct sockaddr *)
            a[4], a[5]);
        break;
    ...
}
return err;
```

← User (args) Kernel (a)  
 nargs[call]  
 ← call action .

←

←

UDP call SYS\_SENDTO  
 sys\_sendto 가 long array a  
 fd(file descriptor), \*buff(pointer of data buffer), length, flags, addr(sockaddr type struct), addr\_len sys\_sendto

1.1 sys\_sendto in net/socket.c  
 → asmlinkage long sys\_sendto(int fd, void \* buff, size\_t len, unsigned flags, struct sockaddr \*addr, int addr\_len)

#### Datagram

```
struct socket *sock;
struct msghdr msg;
sock = sockfd_lookup(fd, &err);
...
if(addr)
{
    err = move_addr_to_kernel(addr,
        addr_len, address);
}
err = sock_sendmsg(sock, &msg, len);
```

← socket  
 ← message header  
 ← file descriptor(file handle)  
 socket slot  
 ← msg

← socket address User  
 Kernel

←



1.1.2 move\_addr\_to\_kernel in net/socket.c  
 → int move\_addr\_to\_kernel(void \*uaddr, int ulen, void \*kaddr)

```
if(ulen<0||ulen>MAX_SOCKET_ADD)
    return ?EINVAL;
```

← length 가



## Data Sending

```

if(ulen==0)
    return 0;
if(copy_from_user(kaddr,uaddr,ulen))
    return ?EFAULT;
return 0;
    
```

← length 0 ' 가  
 ← User address  
 Kernel .  
 copy\_from\_user kaddr, uaddr, ulen copy\_from\_user  
 to, from, size uaddr ulen 가 kaddr  
 copy\_from\_user Socket Creation

1.1.3 sock\_sendmsg in net/socket.c  
 → int sock\_sendmsg(struct socket \*sock, struct msghdr \*msg, int size)

```

struct scm_cookie scm;
err = scm_send(sock, msg, &scm);
if (err >= 0) {
    err = sock->ops->sendmsg(sock,
        msg, size, &scm);
    scm_destroy(&scm);
}
    
```

← scm\_cookie  
 ← scm\_send ( )  
 ← sock->ops sendmsg  
 , sock UDP  
 sendmsg udp\_sendmsg 가  
 sock UDP sendmsg udp\_sendmsg 가

1.1.3.1 scm\_send in include/net/scm.h  
 → static \_\_inline\_\_ int scm\_send(struct socket \*sock, struct msghdr \*msg, struct scm\_cookie \*scm)

```

scm->creds.uid = current->uid;
scm->creds.gid = current->gid;
scm->creds.pid = current->pid;
return __scm_send(sock, msg, scm);
    
```

← scm uid(user ID), gid(group ID),  
 pid(process ID) task uid, gid, pid  
 ← \_\_scm\_send  
 current include/asm\_i386/curreht.h  
 get\_current(void) 가 task task\_struct

1.1.3.1.2 \_\_scm\_send in net/core/scm.c  
 → int \_\_scm\_send(struct socket \*sock, struct msghdr \*msg, struct scm\_cookie \*p)

```

for (cmsg = CMSG_FIRSTHDR(msg);
    cmsg;
    cmsg=CMSG_NXTHDR(msg, cmsg))
{
    if (cmsg->cmsg_len < ... )
        goto error;
    if (cmsg->cmsg_level!=SOL_SOCKET)
        continue;
    switch (cmsg->cmsg_type) {
        case SCM_RIGHTS:
            ...
        case SCM_CREDENTIALS:
            ...
        default:
            goto error;
    }
}
    
```

← message header(msg)  
 control message header  
 ← message length 가  
 ← message level  
 ← message type  
 ← rw: access rights (array of int)  
 ← rw: struct ucred  
 ← error



```

return 0;
error:
    scm_destroy(p);
}
    
```

socket level control message type  
SCM\_CONNECT 가

← 가  
←  
\_\_scm\_destory  
include/linux/socket.h

### 1.1.3.2 udp\_sendmsg

in net/ipv4/udp.c

→ int udp\_sendmsg(struct sock \*sk, struct msghdr \*msg, int len)

```

struct ipcm_cookie ipc;
struct udpfakehdr ufh;
struct rtable *rt = NULL;
...

if (msg->msg_controllen) {
    err = ip_cmsg_send(msg, &ipc);
    ...
}
...
if (MULTICAST(daddr))
    ...

if (connected)
    rt = (struct rtable*)sk_dst_check(sk,0);
...
if (rt == NULL)
    err = ip_route_output(&rt, daddr,
        ufh.saddr, tos, ipc.oif);
    ...
    if (rt->rt_flags&RTCF_BROADCAST &&
        !sk->broadcast)
        goto out;
    if (connected)
        sk_dst_set(sk,
            dst_clone(&rt->u.dst));

if (msg->msg_flags&MSG_CONFIRM)
    goto do_confirm;
back_from_confirm:
...
err = ip_build_xmit(sk,
    (sk->no_check==UDP_CSUM_NOXMIT
    ? udp_getfrag_nosum
    : udp_getfrag),
    &ufh,ulen, &ipc, rt, msg->msg_flags);

out:
    
```

← len, msg  
← ip control message가  
← 가 multicast  
ipc.oif ufh.saddr ,  
← UDP  
← sk destination entry rtable  
cast rt  
← rtable  
← rt 가  
← (broadcast)  
← UDP routing table  
socket  
← Network path  
path  
← UDP header(ufh)  
←  
← socket no\_check 가  
UDP\_CSUM\_NOXMIT  
check sum  
check sum





```

ip_rt_put(rt);
...
do_confirm:
dst_confirm(&rt->u.dst);
...

```

← routing table rt  
 ←  
 ← path confirm  
 ← rt path  
 ← msg->msg\_flags out  
 back\_from\_confirm goto

1.1.3.2.1 ip\_cmsg\_send in net/ipv4/ip\_sockglue.c

```

→ int ip_cmsg_send(struct msghdr *msg, struct ipcm_cookie *ipc)
    __scm_send (1.1.3.1.2)
    msg control message length level type
control message type include/linux/in.h IP_RETOPTS
IP_PKTINFO

```

1.1.3.2.3 sk\_dst\_check in include/net/sock.h

```

→ static inline struct dst_entry * sk_dst_check(struct sock *sk, u32 cookie)
    sk destination 가 sk dst_cache
dst_entry
sk destination field NULL

```

1.1.3.2.3.1 sk\_dst\_get in include/net/sock.h

```

→ static inline struct dst_entry * sk_dst_get(struct sock *sk)
    sk dst_cache dst_entry

```

1.1.3.2.3.2 sk\_dst\_reset in include/net/sock.h

```

→ static inline void sk_dst_reset(struct sock *sk)
    __sk_dst_reset(sk)

```

1.1.3.2.3.2.1 \_\_sk\_dst\_reset in include/net/sock.h

```

→ static inline void __sk_dst_reset(struct sock *sk)
    sk dst_cache NULL dst_cache
(dst_release).

```

1.1.3.2.4 ip\_route\_output in include/net/route.h

```

→ static inline int ip_route_output(struct rtable **rp, u32 daddr, u32 saddr, u32 tos,
int oif)
{ dst :daddr, src:saddr, oif:oif, tos:tos } rt_key routing table 가
(ip_route_output_key ).

```

1.1.3.2.4.1 ip\_route\_output\_key in net/ipv4/route.c

```

→ int ip_route_output_key(struct rtable **rp, const struct rt_key *key)
    key hash code [rt_hash_code() in net/ipv4/route.c]
rt_hash_table chain entry 가
ip_route_output_slow

```

1.1.3.2.6 sk\_dst\_set in include/net/sock.h

```

→ static inline void sk_dst_set(struct sock *sk, struct dst_entry *dst)
    __sk_dst_set(sk, dst)

```



1.1.3.2.6.1 \_\_sk\_dst\_set in include/net/sock.h

```
→ static inline void __sk_dst_set(struct sock *sk, struct dst_entry *dst)
__sk_dst_reset
sk          dst_cache          (dst_release)
dst
```

1.1.3.2.7 ip\_build\_xmit in net/ipv4/ip\_output.c

```
→ int ip_build_xmit(struct sock *sk,
                    int getfrag( const void *, char *, unsigned int, unsigned int),
                    const void *frag, unsigned length, struct ipcm_cookie *ipc,
                    struct rtable *rt, int flags )
(fragmentation)
```

```
struct sk_buff *skb;
struct iphdr *iph;
if (!sk->protinfo.af_inet.hdrincl) {
    ...
    if (length > rt->u.dst.pmtu
        || ipc->opt != NULL)
        return ip_build_xmit_slow(...);

} else {
    if (length > rt->u.dst.dev->mtu)
        ...
}
if (ip_dont_fragment(sk, &rt->u.dst))
    df = htons(IP_DF);
...
...
...
skb = sock_alloc_send_skb (...);
skb_reserve(skb, hh_len);

skb->priority = sk->priority;
skb->dst = dst_clone(&rt->u.dst);
skb->nh.iph = iph = (struct iphdr
                    *)skb_put(skb, length);
if (!sk->protinfo.af_inet.hdrincl) {
    iph->version=4;
    iph->ihl=5;
    ...
    iph->check = ip_fast_csum(...);
}
...
err = NF_HOOK(PF_INET,
              NF_IP_LOCAL_OUT, skb, NULL,
              rt->u.dst.dev, output_maybe_reroute);
```

← ip header  
 ← socket af\_inet header가  
 ← ( )가 HOST mtu control message option ip\_build\_xmit\_slow(...)  
 ← ( ) 가 device mtu  
 ← path  
 가 IP\_DF  
 include/net/ip.h IP  
 flags 'Don t 'Fragment"  
 df iph->frag\_off  
 ← buffer  
 ← skb data, tail 가  
 ←  
 ←  
 ← ip header  
 ← iph IP header  
 ← hooking, output\_maybe\_reroute  
 skb Kernel  
 가 Device



...

1.1.3.2.7.2 sock\_alloc\_send\_skb in net/core/sock.c  
 → struct sk\_buff \*sock\_alloc\_send\_skb(struct sock \*sk, unsigned long size, int noblock, int \*errcode)  
 send/receive buffer  
 sk\_buffer (alloc\_skb)  
 (skb\_set\_owner\_w)

1.1.3.2.7.2.2 skb\_set\_owner\_w in include/net/sock.h  
 → static inline void skb\_set\_owner\_w(struct sk\_buff \*skb, struct sock \*sk)  
 skb sk sk destructor

```

skb->sk = sk;
skb->destructor = sock_wfree;
atomic_add(skb->truesize,
            &sk->wmem_alloc);
    
```

← sk  
 ←  
 ← sk wmem\_alloc skb->truesize

sock\_wfree net/core/sock.c kfree\_skb

1.1.3.2.7.3 skb\_reserve in include/linux/skbuff.h  
 → static inline int skb\_skb\_reserve(struct sk\_buff \*skb, unsigned int len)  
 headroom  
 skb data tail len 가

1.1.3.2.7.5 ip\_select\_ident in include/net/ip.h  
 → static inline void ip\_select\_ident(struct iphdr \*iph, struct dst\_entry \*dst, struct sock \*sk)  
 ip header id (identification) \_\_ip\_select\_ident

1.1.3.2.7.5.1 \_\_ip\_select\_ident in net/ipv4/route.c  
 → void \_\_ip\_select\_ident(struct iphdr \*iph, struct dst\_entry \*dst)  
 routing table( dst rtable ) peer [extern inline \_\_u16  
 inet\_getid(struct inet\_peer \*p) in include/net/inetpeer.h] id id

1.1.3.2.7.5.1.1 inet\_getid in include/net/inetpeer.h  
 → extern inline \_\_u16 inet\_getid(struct inet\_peer \*p)  
 p ip\_id\_count  
 1 가 ip\_id\_count

### Help Comment

Kernel	가	Device
--------	---	--------



1.1.3.2.7.7 NF\_HOOK in include/linux/netfilter.h  
 → #define NF\_HOOK(pf, hook, skb, indev, outdev, okfn) (okfn)(skb)  
 output\_maybe\_reroute skb  
 debug debug NF\_HOOK가  
 nf\_hook\_slow [in net/core/netfilter.c]

1.1.3.2.7.7.1 output\_maybe\_reroute in net/ipv4/ip\_output.c



→ static inline int output\_maybe\_reroute(struct sk\_buff \*skb) [skb->dst ->outout(skb)] . ip\_output 가 .

1.1.3.2.7.7.1.1 ip\_output in net/ipv4/ip\_output.c

→ int ip\_output(struct \*skb)  
ip\_finish\_output(skb)

1.1.3.2.7.7.1.1.2 ip\_finish\_output in net/ipv4/ip\_output.c

→ \_\_inline\_\_ int ip\_finish\_output(struct sk\_buff \*skb)  
skb dev protocol network filter hooking( , NF\_HOOK)  
ip\_finish\_output2 .

1.1.3.2.7.7.1.1.2.1 ip\_finish\_output2 in net/ipv4/ip\_output.c

→ static inline int ip\_finish\_output2(struct sk\_buff \*skb)  
skb dst(destination) hh(hh\_cache )가 hh output [ dev\_queue\_xmit()] skb .

1.1.3.2.7.7.1.1.2.1.1 dev\_queue\_xmit in net/core/dev.c

→ int dev\_queue\_xmit(struct sk\_buff \*skb)  
(skb ->dev)

```

struct Qdisc *q;
...
if (skb ->ip_summed == CHECKSUM_HW
    && ... ) {
    if ((skb = skb_checksum_help(skb))
        == NULL )
        return ?ENOMEM;
}
q = dev->qdisc;
if (q ->enqueue) {
    int ret = q ->enqueue(skb, q);
    ...
    return ...;
}
...
    
```

← checksum 가  
가  
checksum  
checksum  
←  
← 가  
← . ( 3COM 3c509  
el3\_start\_xmit 가  
)  
← 가 ( loopback software  
device가 ).  
el3\_start\_xmit

3COM 3c509

1.1.3.2.7.7.1.1.2.1.1.1 el3\_start\_xmit in drivers/net/3c509.c

→ start int el3\_start\_xmit(struct sk\_buff \*skb, struct net\_device \*dev)

```

struct el3_private *lp
    = (struct el3_private *)dev->priv;
...
netif_stop_queue (dev);
lp ->stats.tx_bytes += skb ->len;
...
    
```

← driver specific  
← device queue  
(race condition )  
← ( )



```

outw(skb ->len, ioaddr + TX_FIFO); ←
...
outsl(ioaddr + TX_FIFO, skb ->data, ←
      (skb ->len + 3) >> 2); ←

```

port, src, count

1.1.3.2.7.7.1.1.2.1.1.1.1.1.1 outsl in include/asm\_sh/io.h  
 → void outsl(unsigned long port, void \*src, unsigned long count)  
 Kernel Device

```

outsl(ioaddr + TX_FIFO, skb ->data, (skb ->len + 3) >> 2);
, "skb ->data" "(skb ->len+3) >> 2" "ioaddr + TX_FIFO"

```

1.1.3.2.7.8 ip\_build\_xmit\_slow in net/ipv4/ip\_output.c  
 → int ip\_build\_xmit(struct sock \*sk,  
 int getfrag( const void \*, char \*, unsigned int, unsigned int),  
 const void \*frag, unsigned length, struct ipcm\_cookie \*ipc,  
 struct rtable \*rt, int flags )  
 ip\_build\_xmit  
 ip\_build\_xmit 가

```

...
fraglen = length - offset + fragheaderlen; ←
...
id = sk->protinfo.af_inet.id++; ←
id ←
do { ←
    char *data;
    struct sk_buff *skb;
    skb = sock(alloc_send_skb( ... ) ); ← buffer
    ...
    data = skb_put(skb, fraglen); ←
    ...
} while (offset >= 0); ←
...

```

1.1.3.2.8 ip\_rt\_put in include/net/route.h  
 → static inline void ip\_rt\_put(struct rtable \* rt)  
 rt dst (dst\_entry) (dst\_release(&rt->dst))

1.1.3.2.9 dst\_confirm in include/net/dst.h  
 → static inline void dst\_confirm(struct dst\_entry \*dst)  
 entry가 dst neighbour neigh\_confirm

1.1.3.2.9.1 neigh\_confirm in include/net/neighbour.h  
 → static inline int neigh\_confirm(struct neighbour \*neigh)  
 neigh가 confirmed jiffies



## 6. Data Receiving

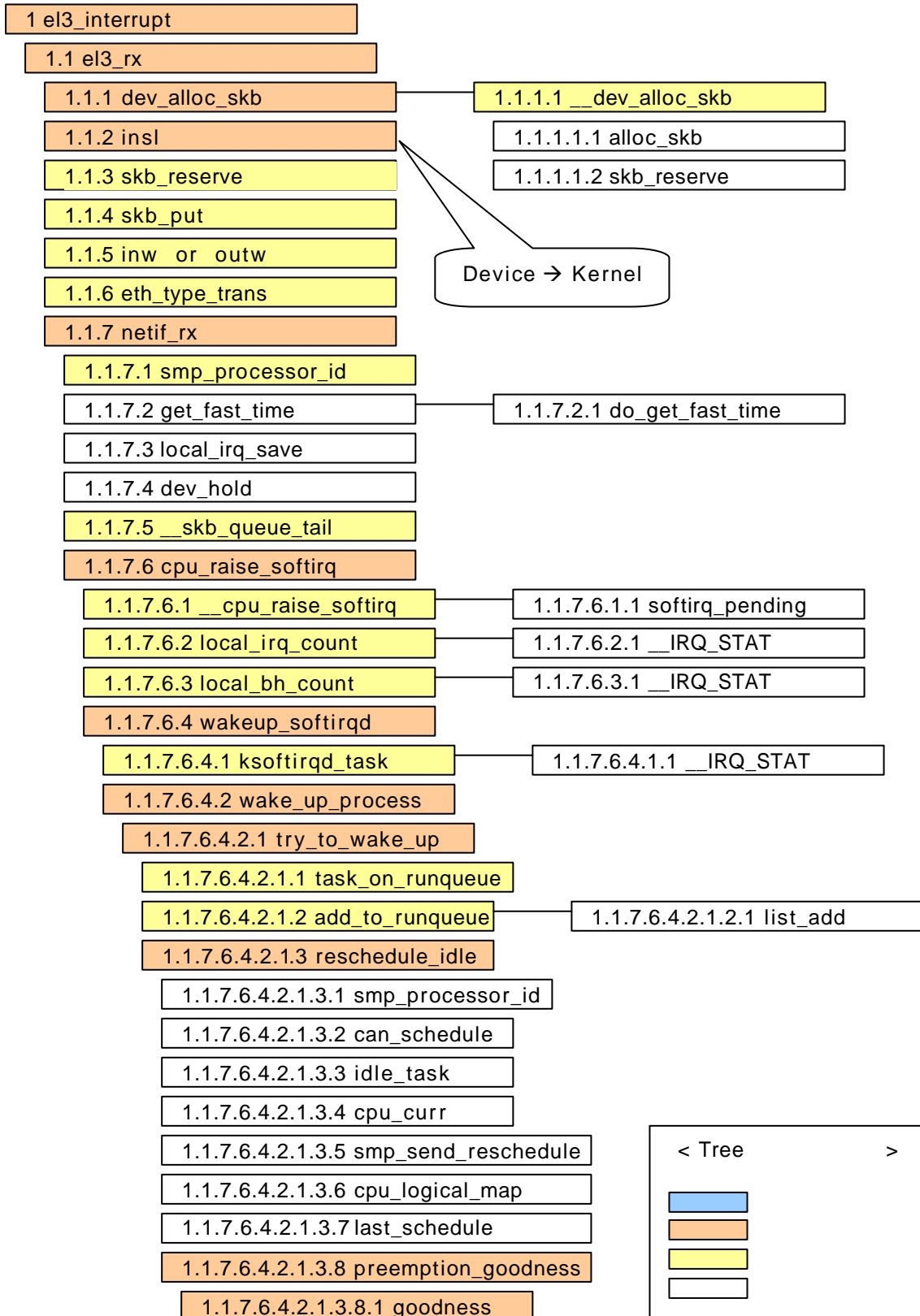
Comment

UDP

1. NIC가 CPU softirq
2. softirq Sub System NET\_RX\_SOFTIRQ raise  
net\_rx\_action CPU layer sock  
IP layer UDP layer
3. sock ( )가 system call(recvfrom )  
가 가  
가  
Data Receiving Part I, II, III  
Part



## Data Receiving - Part I Function Call Tree





## Tree Analysis

### Help Comment - interrupt

NIC가 layer broadcast	local MAC address 가	. NIC DMA PIO	link
	RAM NIC	sk buffer	Tree
1 el3_interrupt '	1.1 el3_probe '		
	1.1.7 netif_rx(skb) '		

1. el3\_interrupt in drivers/net/3c509.c  
 → static void el3\_interrupt(int irq, void \*dev\_id, struct pt\_regs \*regs)

1.1 el3\_rx in drivers/net/3c509.c  
 → static int el3\_rx(struct net\_device \*dev)

```

struct el3_private *lp = (struct el3_private *)dev->priv; ← device specific , device
int ioaddr = dev->base_addr; ← device

short rx_status; ←
while ((rx_status ← device
        = inw(ioaddr+RX_STATUS)) > 0)
{
    if (rx_status & 0x4000) { ← error
        outw(RxDiscard,ioaddr+EL3_CMD); ←
        lp->stats.rx_errors++; ←
        switch (error) { ←
            case 0x0000: ←
                lp->stats.rx_over_error++; break; ←
            case 0x0800: ←
                lp->stats.rx_length_error++; break; ←
            ... ←
        } ← error
    } else { ←
        short pkt_len = rx_status & 0x7ff; ← (rx_status 12bit)
        struct sk_buff *skb; ←
        skb = dev_alloc_skb (pkt_len+5); ← ( + 5,
        5' )
        lp->stats.rx_bytes += pkt_len; ←
        if (skb != NULL) { ←
            insl(ioaddr+RX_FIFO, skb_put(skb, ←
                pkt_len), (pkt_len+3) >> 2); ←
            outw(RxDiscard,ioaddr+EL3_CMD); ← Kernel device
            skb->protocol = ←
                eth_type_trans(skb, dev); ←
            netif_rx(skb); ← device layer
        }
    }
}
    
```





Data Receiving

```

lp ->stats.rx_packets++;
continue;
}

outw(RxDiscard, ioadr + EL3_CMD);
lp ->stats.rx_dropped++;
}
...
}
    
```

← 가 .  
 ←  
 ←  
 ← device  
 ← 가 .

1.1.1 dev\_alloc\_skb in include/linux/skbuff.h

```

→ static inline struct sk_buff *dev_alloc_skb(unsigned int length)
sk_buff
__dev_alloc_skb GFP_ATOMIC gfp_mask
    
```

1.1.2 insl in include/asm\_sh/io.h

```

→ void insl(unsigned long port, void *dst, unsigned long count)
Device Kernel
    
```



```

insl(ioaddr+RX_FIFO, skb_put(skb, pkt_len), (pkt_len+3) >> 2);
, 'ioaddr+RX_FIFO " "(pkt_len+3) >> 2 " 'skb_put(skb,
pkt_len) "
    
```

1.1.7 netif\_rx in /net/core/dev.c

```

→ int netif_rx(struct sk_buff *skb)
device layer
    
```

```

int this_cpu = smp_processor_id();
struct softnet_data *queue;
queue = &softnet_data[this_cpu];

netdev_rx_stat[this_cpu].total++;
if (queue->input_pkt_queue qlen
    <= netdev_max_backlog)
{
    if (queue->input_pkt_queue qlen)
    {
        if (queue->throttle)
            goto drop;
        enqueue:
        __skb_queue_tail(
            &queue->input_pkt_queue, skb);
        cpu_raise_softirq(this_cpu,
            NET_RX_SOFTIRQ);
    }
    return
    
```

← CPU ID  
 ← queue  
 ← softnet\_data  
 CPU  
 ← CPU 가  
 ← 가 backlog  
 ← 가 0'  
 ← 가  
 ← softirq Bottom Half  
 NET\_RX\_SOFTIRQ raise CPU softirq  
 )  
 ← CPU congestion level



```

    softnet_data[this_cpu].cng_level;
}
if (queue->throttle) {
    queue->throttle = 0;
#ifdef CONFIG_NET_HW_FLOWCONTROL
    if (atomic_dec_and_test(
        &netdev_dropping))
        netdev_wakeup();
#endif
}
goto enqueue;
}

if (queue->throttle == 0) {
    queue->throttle = 1;
    netdev_rx_stat[this_cpu].throttled++;
#ifdef CONFIG_NET_HW_FLOWCONTROL
    if (atomic_inc(&netdev_dropping))
#endif
}
drop:
netdev_rx_stat[this_cpu].dropped++;
...

```

Important Data Structure softnet\_data  
 netdev\_dropping Network Flow Control

Help Comment - softirq

Kernel	2.4	softirq	가	network stack
Bottom Half	softirq가			
Bottom Half	CPU가			CPU
	softirq	CPU		SMP

1.1.7.6 cpu\_raise\_softirq in kernel/softirq.c  
 → inline void cpu\_raise\_softirq(unsigned int cpu, unsigned int nr)  
 CPU softirq nr raise(pending) [\_\_cpu\_raise\_softirq(cpu, nr) in  
 include/linux/interrupt.h], CPU local irq local bottom half 가  
 0' softirq daemon [wakeup\_softirqd() in kernel/softirq.c]  
 NET\_RX\_SOFTIRQ softirq  
 do\_softirq() 가 net\_rx\_action() 가  
 do\_softirq(), net\_rx\_action() Data Receiving Part II

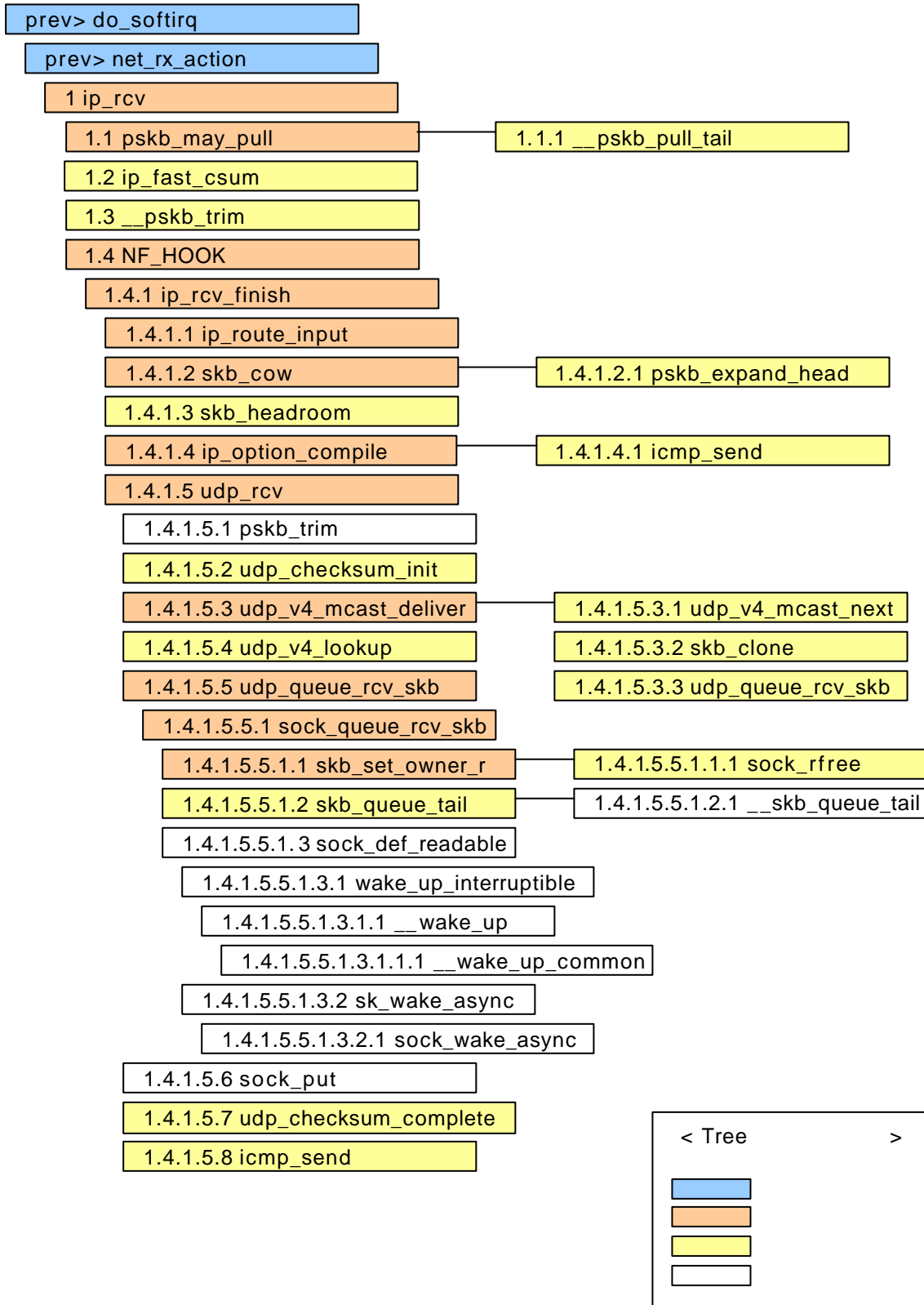




weight = -1;	←	
if (p->policy & SCHED_YIELD)	←	Real-Time
goto out;		
if (p->policy == SCHED_OTHER) {	←	가
weight = p->counter;	←	clock tick
if (!weight)	←	clock tick 0 ( , )
goto out;		
#ifdef CONFIG_SMP	←	SMP
if (p->processor == this_cpu)	←	가
weight += PROC_CHANGE_PENALTY;		
#endif		( )가 가
if (p->mm == this_mm    !p->mm)	←	가
weight += 1;		
...		가
}		
weight = 1000 + p->rt_priority;	←	Real-Time 가
out:		( )
return weight;	←	가



## Data Receiving - Part II Function Call Tree





### Tree Analysis

가 [schedule() in kernel/sched.c] softirq 가  
do\_softirq [asmlinkage void do\_softirq() in  
kernel/softirq.c]가  
UDP Data Receiving Part I  
CPU 1.1.7.6.4 wakeup\_softirqd '  
softirq 가 Running Queue 가 .  
do\_softirq .  
do\_softirq CPU softirq pending  
1.1.7.6 cpu\_raise\_softirq ' NET\_RX\_SOFTIRQ  
net\_rx\_action [static void net\_rx\_action(struct softirq\_action \*h) in net/core/dev.c]가

net\_rx\_action CPU (struct  
packet\_type) [packet\_type.func(struct sk\_buff \*, struct  
net\_device \*, struct packet\_type\*)  
UDP , IP ip\_rcv 가 .

1. ip\_rcv in net/ipv4/ip\_input.c  
→ int ip\_rcv(struct sk\_buff \*skb, struct net\_device \*dev, struct packet\_type \*pt)  
IP 가 .

```

if( skb ->pkt_type ← 가
    == PACKET_OTHERHOST)
    goto drop;

if ((skb = skb_share_check(skb, ← 가
    GFP_ATOMIC)) == NULL)
    goto out;
if (!pskb_may_pull(skb, ← 가 IP
    sizeof(struct iphdr)))
    goto inhdr_error;
...
if (!pskb_may_pull(skb, iph ->ipl*4) ← 가 IP length
    goto inhdr_error;
if (ip_fast_csum((u8 *)iph, iph ->ihl) != 0) ← check sum 가
    goto inhdr_error;
...
if (skb ->len > len) { ← 가
    __pskb_trim(skb, len); ( )
    ...
}
return NF_HOOK(..., ip_rcv_finish); ← network filter hooking
... ( )
    
```

1.1 pskb\_may\_pull in include/linux/skbuff.h  
→ static inline int pskb\_may\_pull(struct sk\_buff \*skb, unsigned int len)  
skb pull( ) 가  
skb 가 len true, skb 가 len



```
false
...
skb len-skb
[__pskb_pull_tail(struct sk_buff *skb, int delta) in net/core/skbuff.c].
```

```
1.4 NF_HOOK in include/linux/netfilter.h
→ #define NF_HOOK(pf, hook, skb, indev, outdev, okfn) (okfn)(skb)
    ip_rcv_finish      skb
    debug              debug          NF_HOOK가
nf_hook_slow [in net/core/netfilter.c]
```

```
1.4.1 ip_rcv_finish in net/ipv4/ip_input.c
→ static inline int ip_rcv_finish(struct sk_buff *skb)
    IP
```

```
...
if (iph->ihl > 5 ) {
    struct ip_options *opt;
    if (skb_cow(skb, skb_headroom(skb)) ←
        goto drop;

    ip (ip_options_compile(NULL, skb)) ← IP option
        goto inhdr_error;
    opt = &(IPCB(skb)->opt); ← compile IP
    ...
}
return skb->dst->input(skb); ← input          udp_rcv
...
가
```

```
1.4.1.2 skb_cow in include/linux/skbuff.h
→ static inline int skb_cow(struct sk_buff *skb, unsigned int headroom)
    skb [int
pskb_expand_head(struct sk_buff *skb, int nhead, int ntail, int gfp_mask) in
net/core/skbuff.c]
```

```
1.4.1.4 ip_option_compile in net/ipv4/ip_options.c
→ int ip_options_compile(struct ip_options * opt, struct sk_buff *skb)
    skb          opt가 NULL
    skb->data가 IP 가
    skb          parameter
error가 ICMP [void icmp_send(struct sk_buff *skb_in, int
type, int code, u32 info) in net/ipv4/icmp.c]
```

Comment

IP layer 가 , TCP layer 가

```
1.4.1.5 udp_rcv in net/ipv4/udp.c
→ int udp_rcv(struct sk_buff *skb)
```

```
struct sock *sk;
struct udphdr *uh; ← UDP
struct rtable *rt =(struct rtable*)skb->dst; ← routing table
```





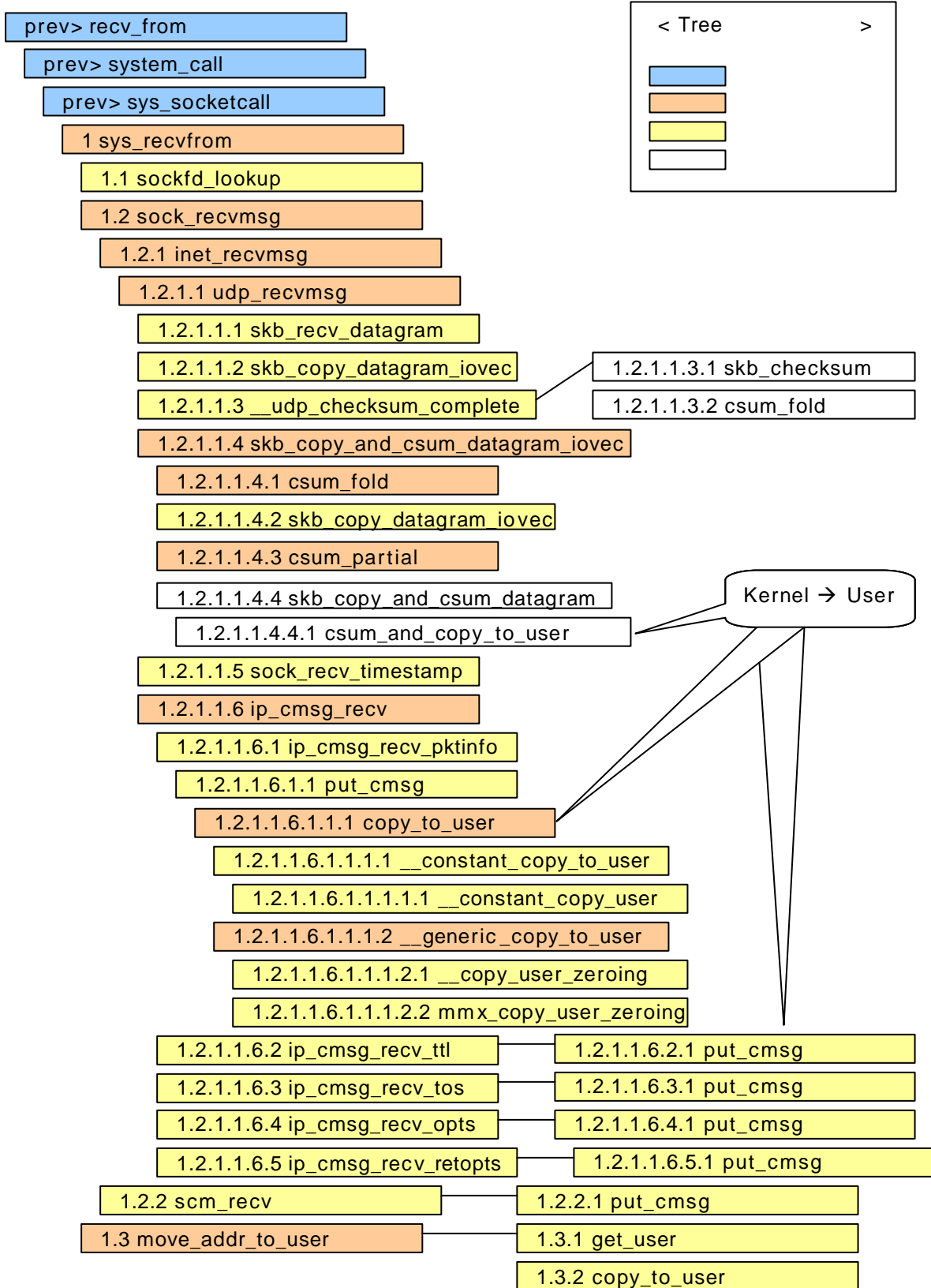


1.4.1.5.5.1 sock\_queue\_rcv\_skb in include/net/sock.h  
→ static inline void skb\_set\_rcv\_skb(struct sock \*sk, struct sk\_buff \*skb)  
    skb                          [skb\_set\_owner\_r in include/net/sock.h]    sk  
skb

1.4.1.5.5.1.1 skb\_set\_owner\_r in include/net/sock.h  
→ static inline void skb\_set\_owner\_r(struct sk\_buff \*skb, struct sock \*sk)  
    skb                          가 sk  
    skb                          [void sock\_rfree(struct sk\_buff \*skb) in net/core/sock.c]  
                                  가                                  가                                  ( )



## Data Receiving - Part III Function Call Tree





## Tree Analysis

( )가 recvfrom system call  
 (interrupt vector index 0x80, "#define SYSCALL\_VECTOR 0x80" in include/asm-i386/hw\_irq.h) . system call interrupt handler sys\_socketcall  
 가 UDP call SYS\_RECVFROM  
 sys\_recvfrom 가 . (sys\_socketcall Socket Creation Data Sending )

1. sys\_recvfrom in net/socket.c  
 → asmlinkage long sys\_recvfrom(int fd, void \* ubuf, size\_t size, unsigned flags, struct sockaddr \*addr, int \*addr\_len)

```

UDP 가
struct socket *sock;
struct msghdr msg;
...
sock = sockfd_lookup(fd, &err);
...
err=sock_recvmsg(sock, &msg, size,
                  flags);
if ( err >= 0 && addr != NULL
    && msg.msg_namelen)
{
    err2=move_addr_to_user(address,
                          msg.msg_namelen, addr, addr_len);
    ...
}
...
    
```

← fd socket  
 (fd : file handle)  
 ← msg  
 ←  
 ← 가  
 ← User address  
 (Kernel -> User)

1.2 sockfd\_recvmsg in net/socket.c  
 → int sock\_recvmsg(struct socket \*sock, struct msghdr \*msg, int size, int flags)

```

struct scm_cookie scm;
size = sock->ops->recvmsg(sock, msg,
                          size, flags, &scm);

if (size >= 0)
    scm_recv(sock, msg, &scm, flags);
return size;
    
```

← socket control message cookie  
 ← sock operation  
 inet inet\_recvmsg  
 가  
 ← msg control message  
 ←

1.2.1 inet\_recvmsg in net/ipv4/af\_inet.c  
 → int inet\_recvmsg(struct socket \*sock, struct msghdr \*msg, int size, int flags, struct scm\_cookie \*scm)

sk->prot->recvmsg( layer )  
 UDP udp\_recvmsg 가

1.2.1.1 udp\_recvmsg in net/ipv4/udp.c  
 → static \_\_inline\_\_ int udp\_checksum\_complete(struct sk\_buff \*skb)



```

...
skb = skb_rcv_datagram(sk, flags, noblock, &err); ← sk
...
copied=skb->len ? sizeof(struct udphdr); ← ( )
...
// option
  skb_copy_datagram_iovec(...) ← checksum : User IO Vector
...
  or ( __udp_checksum_complete(skb) ← MSG_TRUNC flag : checksum
    and skb_copy_datagram_iovec(...)
    or skb_copy_and_csum_datagram_iovec(...) ← checksum
  // skb io vector
...
sock_rcv_timestamp(msg, sk, skb); ←
...
if (sk->protinfo.af_inet.cmsg_flags) ← address
  ip_cmsg_rcv(msg, skb); ← control message flag가 ( )
...

```

1.2.1.1.4 `skb_copy_and_csum_datagram_iovec` in `net/core/datagram.c`  
 → `int skb_copy_and_csum_datagram_iovec(const struct sk_buff *skb, int hlen, struct iovec *iovc)`

`skb` `checksum` `User` `iovc`(IO vector) `skb` .

1.2.1.1.4.1 `csum_fold` in `include/asm-i386/checksum.h`  
 → `static inline unsigned int csum_fold(unsigned int sum)`

`checksum` ( ).

1.2.1.1.4.3 `csum_partial` in `include/asm-i386/checksum.h`  
 → `asmlinkage unsigned int csum_partial(const unsigned char *buff, int len, unsigned int sum)`

(buff) (len) sum

1.2.1.1.6 `ip_cmsg_rcv` in `net/ipv4/ip_sockglue.c`  
 → `void ip_cmsg_rcv(struct msghdr *msg, struct sk_buff *skb)`

IP layer control message .  
`skb->sk->protinfo.af_inet.cmsg_flags` .  
`pktinfo > ttl > tos > opts > reopts`

`ip_cmsg_rcv_xxx(struct msghdr *msg, struct sk_buff *skb)` .  
`put_cmsg` [in `net/core/scm.c`]

control message `User` [copy\_to\_user in `include/asm-i386/uaccess.h`]



1.2.1.1.6.1.1.1 `copy_to_user` in `include/asm-i386/uaccess.h`

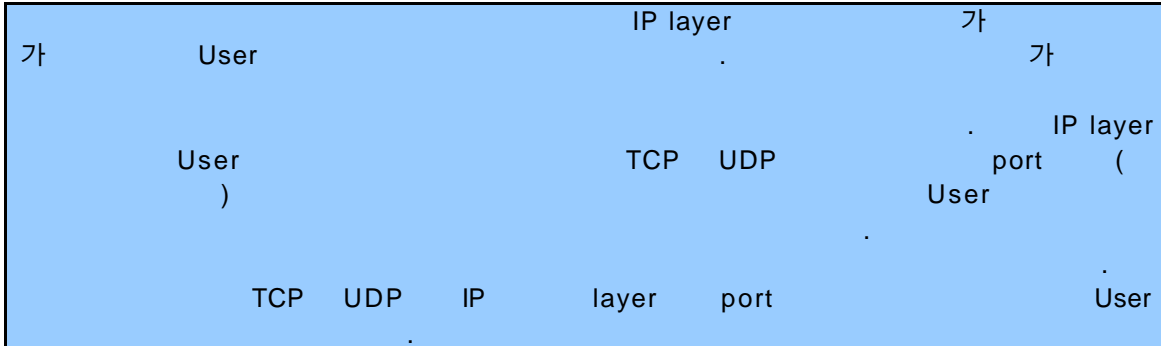
Kernel Memory Area `User` Memory Area .  
 Socket Creation `copy_from_user` .  
`copy_to_user(to, from, n)` , `from` `to` `n`(size)  
`n` `built-in`

`__constant_copy_to_user` , `__generic_copy_to_user` .



\_\_constant\_copy\_to\_user inline

Help Comment ? delayed copy to user



1.2.1.1.6.1.1.1.2 \_\_generic\_copy\_to\_user in arch/i386/lib/usercopy.c  
 → unsigned long \_\_generic\_copy\_to\_user(void \*to, const void \*from, unsigned long n)

inline . 3DNow AMD 3DNow \_\_copy\_\_user mmx\_copy\_user

1.3 move\_addr\_to\_user in net/socket.c

→ int move\_addr\_to\_user(void \*kaddr, int klen, void \*uaddr, int \*ulen)

Kernel address User  
 kaddr: , klen: , uaddr: User  
 ulen: User 가

```

if ((err=get_user(len, ulen))) ← User ulen
    return err; ← error
if (len>klen) ←
    len = klen;
if (len<0 || len> MAX_SOCKET_ADDR) ← 가 가
    return ?EINVAL; ← error
if (len) { ←
    if (copy_to_user(uaddr, kaddr, len)) ← Kernel address User
        return ?EFAULT;
}
return __put_user(klen, ulen); ←
    
```

copy\_to\_user



### Appendix

#### Important Data Structures

struct net\_proto\_family in include/linux/net.h

```

int family;
int (* create)(struct socket *sock, int protocol);
short authentication;
short encryption;
short encrypt_net;

struct inet_family_ops, unix_family_ops xxx_family_ops struct
inet_family_ops
struct net_proto_family inet_family_ops = {
    family: PF_INET,
    create: inet_create
};
, inet_family_ops PF_INET inet_create

```

struct proto\_ops in include/linux/net.h

```

int family;
int (*release) (struct socket *sock);
int (*bind) (struct socket *sock, struct sockaddr *umyaddr, int sockaddr_len);
int (*connect) (struct socket *sock, struct sockaddr *uaddr, int sockaddr_len, int flags);
int (*socketpair) (struct socket *sock1, struct socket *sock2);
int (*accept) (struct socket *sock, struct socket *newsock, int flags);
int (*getname) (struct socket *sock, struct sockaddr *uaddr, int *usockaddr_len, int peer);
unsigned int (*poll) (struct file *file, struct socket *sock, struct poll_table_struct *wait);
int (*ioctl) (struct socket *sock, unsigned int cmd, unsigned long arg);
int (*listen) (struct socket *sock, int len);
int (*shutdown) (struct socket *sock, int flags);
int (*setsockopt) (struct socket *sock, int level, int optname, char *optval, int optlen);
int (*getsockopt) (struct socket *sock, int level, int optname, char *optval, int *optlen);
int (*sendmsg) (struct socket *sock, struct msghdr *m, int total_len, struct scm_cookie *scm);
int (*recvmsg) (struct socket *sock, struct msghdr *m, int total_len, int flags, struct scm_cookie *scm);
int (*mmap) (struct file *file, struct socket *sock, struct vm_area_struct *vma);
ssize_t (*sendpage) (struct socket *sock, struct page *page, int offset, size_t size, int flags);

```

socket

struct socket in include/linux/net.h

```

socket_state state;

unsigned long flags;
struct proto_ops *ops;
struct inode *inode;
struct fasync_struct *fasync_list; /* Asynchronous wake up list */
struct file *file; /* File back pointer for gc */
struct sock *sk;
wait_queue_head_t wait;

```



```
short          type;
unsigned char  passcred;
```

```
static struct net_proto_family *net_families[NPROTO];           in net/socket.c
                                protocol          list          . net_proto_family   Network
Initialization      Important Data Struct
```

```
struct msghdr                                     in include/net/route.h
```

```
union
{
    struct dst_entry  dst;
    struct rtable     *rt_next;
} u;
unsigned            rt_flags;
unsigned            rt_type;
__u32               rt_dst; /* Path destination */
__u32               rt_src; /* Path source */
int                 rt_iif;
/* Info on neighbour */
__u32               rt_gateway;
/* Cache lookup keys */
struct rt_key       key;
/* Miscellaneous cached information */
__u32               rt_spec_dst; /* RFC1122 specific destination */
struct inet_peer   *peer; /* long-living peer info */
#ifdef CONFIG_IP_ROUTE_NAT
__u32               rt_src_map;
__u32               rt_dst_map;
#endif
#endif
```

```
struct iphdr                                     in include/linux/ip.h
```

```
#if defined(__LITTLE_ENDIAN_BITFIELD)
__u8   ihl:4,
       version:4;
#elif defined (__BIG_ENDIAN_BITFIELD)
__u8   version:4,
       ihl:4;
#else
#error "Please fix <asm/byteorder.h>"
#endif
__u8   tos;
__u16  tot_len;
__u16  id;
__u16  frag_off;
__u8   ttl;
__u8   protocol;
__u16  check;
__u32  saddr;
__u32  daddr;
/*The options start here. */
```



struct udpfakehdr in net/ipv4/udp.c

```
struct udphdr uh;
u32 saddr;
u32 daddr;
struct iovec *iov;
u32 wcheck;
```

struct udphdr in include/linux/udp.h

```
__u16 source;
__u16 dest;
__u16 len;
__u16 check;
```

struct scm\_cookie in include/net/scm.h

```
struct ucred creds; /* Skb credentials */
struct scm_fp_list *fp; /* Passed files */
unsigned long seq; /* Connection seqno */
```

→ socket control message

struct ipcm\_cookie in include/net/ip.h

```
u32 addr;
int oif;
struct ip_options *opt;
```

→ ip control message

struct ip\_options in include/linux/ip.h

```
__u32 faddr; /* Saved first hop address */
unsigned char optlen;
unsigned char srr;
unsigned char rr;
unsigned char ts;
unsigned char is_setbyuser:1, /* Set by setsockopt?
 */
is_data:1, /* Options in __data, rather than skb
 */
is_strictroute:1, /* Strict source route */
srr_is_hit:1, /* Packet destination addr was our one
 */
is_changed:1, /* IP checksum more not valid
 */
rr_needaddr:1, /* Need to record addr of outgoing
dev */
ts_needtime:1, /* Need to record timestamp
 */
ts_needaddr:1; /* Need to record addr of outgoing
dev */
```





```
unsigned char router_alert;
unsigned char __pad1;
unsigned char __pad2;
unsigned char __data[0];
```

struct softnet\_data in include/linux/netdevice.h

```
int throttle;
int cng_level;
int avg_blog;
struct sk_buff_head input_pkt_queue;
struct net_device *output_queue;
struct sk_buff *completion_queue;
```

extern struct netif\_rx\_stats netdev\_rx\_stat[] in include/linux/netdevice.h

struct netif\_rx\_stats in include/linux/netdevice.h

```
unsigned total;
unsigned dropped;
unsigned time_squeeze;
unsigned throttled;
unsigned fastroute_hit;
unsigned fastroute_success;
unsigned fastroute_defer;
unsigned fastroute_deferred_out;
unsigned fastroute_latency_reduction;
unsigned cpu_collision;
```

### Comment

가	가	.
---	---	---

struct rtable in include/net/route.h

→ source path

struct dst\_entry in include/net/dst.h

→ source path

struct sk\_buff in include/linux/skbuff.h

→ source path

struct sock in include/net/sock.h

→ source path

struct net\_device in include/linux/netdevice.h or netfilter.h

→ source path



## References

### [Book or Document]

The Linux Kernel, by David A Rusling

Linux Kernel Internals, by Michael Beck, and so on : ADDISON -WESLEY

TCP/IP Illustrated, by Wright Stevens

### [Web Site]

<http://lxr.linux.no/>

<http://www.gnumonks.org/ftp/pub/doc/packet-journey-2.4.html>

<http://linuxkernel.to/module/port-2.4/eng/lkp.html#toc>

### [Source Code]

Linux Kernel version 2.4.13