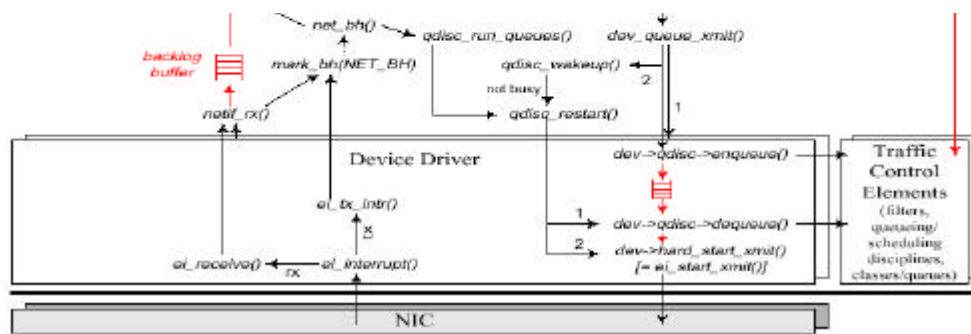


Linux Networking Kernel

Part II

- device driver -



1. Linux Networking Overview

1.1 Introduction

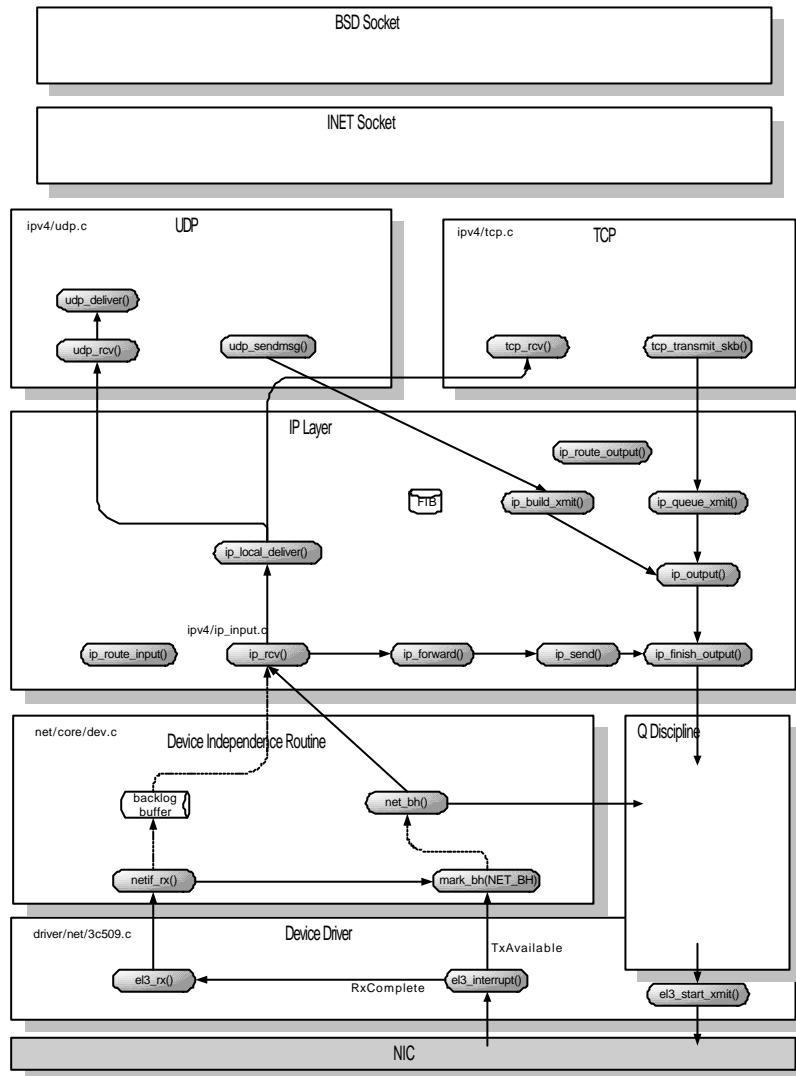
Device driver가 , device
Independence Routine . IP
UDP, TCP Transport layer가 . INET socket
layer가 , BSD socket layer가 가 . BSD socket layer
Application process

/usr/src/linux/net ,
/usr/src/linux/drivers/net .

```
silmaril:/usr/src/linux/net# ls
802/      appletalk/ ethernet/ netlink/  protocols.c  socket.o      x25/
Changes   ax25/      ipv4/      netrom/   protocols.o  sunrpc/
Config.in bridge/    ipv6/      netsyms.c  rose/        sysctl_net.c
Makefile  core/      ipx/       netsyms.o  sched/       sysctl_net.o
README    decnet/    irda/      network.a  sock_n_syms.o  unix/
TUNABLE   econet/    lapb/      packet/    socket.c      wanrouter/
```

```
silmaril:/usr/src/linux/net/core# ls
Makefile  dev.o      filter.c  neighbour.o  scm.o      sysctl_net_core.c
core.o    dev_mcast.c  firewall.c  profile.c    skbuff.c   sysctl_net_core.o
datagram.c  dev_mcast.o  iovec.c    rtnetlink.c  skbuff.o   utils.c
datagram.o  dst.c        iovec.o    rtnetlink.o  sock.c     utils.o
dev.c       dst.o        neighbour.c  scm.c        sock.o
```

Linux Network Overview



1.1

가

1.2

: Data Link Layer

LAN

3c509

driver가

, driver linux/drivers/net/3c509.c .

가

el3_interrupt() 가 2.1 el3_interrupt() 가

가

(RxComplete) el3_rx() 가

(TxAvailable) mark_bh(NET_BH)

mark_bh(NET_BH)

el3_interrupt() el3_rx() el3_rx()

skbuff , netif_rx()

.

netif_rx() net/core/dev.c skbuff backlog

, mark_bh(NET_BH)

mark_bh(NET_BH) 가

Bottom Half Handler 가 ,

, Bottom Half Handler

.

mark_bh(NET_BH) BH

Handler가 mark_bh(NET_BH)

, net_bh()

net_bh() , qdisc_run_queues()

. backlog , dequeue

skbuff 가 skbuff type

. IP , ip_rcv()

1.3 : IP Layer

ip_rcv() net/ipv4/ip_input.c . ip_rcv()

. , ip_local_deliver()

, ip_forward()

ip_local_deliver() port

udp udp_rcv() 가 , tcp tcp_rcv() 가 .

1.4 : UDP Layer

udp_rcv() net/ipv4/udp.c . udp_rcv()

```

        udp
        , udp_deliver()
udp_deliver()
Application

```

1.5

```

        . UDP      udp_sendmsg()   가 socket
        ip_build_xmit()
        . ip_build_xmit()   net/ipv4/ip_output.c   .   skbuff
        .   , ip_output()
ip_output()   ip_finish_out()   . ip_finish_out()   inline
        include/net/ip.h   .   device   output
        .   dev_queue_xmit()
dev_queue_xmit()   net/core/dev.c   .   Queueing
discipline   ,   el3_start_xmit()   가
el3_start_xmit()   drivers/net/3c509.c   .   el3_start_xmit()
skbuff

```

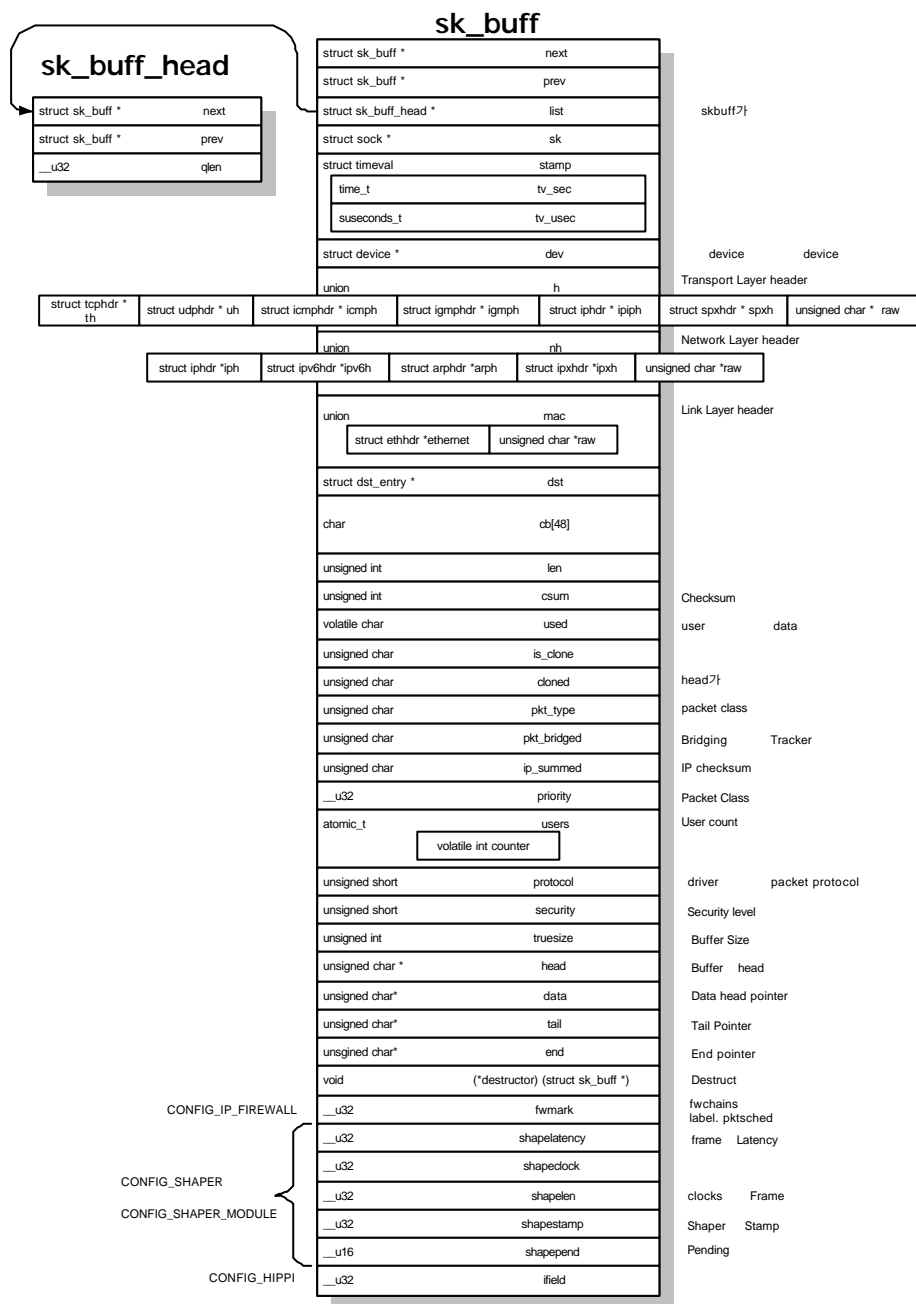
2. skbuff Buffer

2.1 skbuff structure

```

1  #include <linux/skbuff.h>
2
3  struct sk_buff {
4      /* ... */
5      struct sk_buff *next;
6  };
7
8  struct sk_buff_head {
9      /* ... */
10     struct sk_buff *list;
11 };
12
13 #endif

```



2.1 skbuff structure

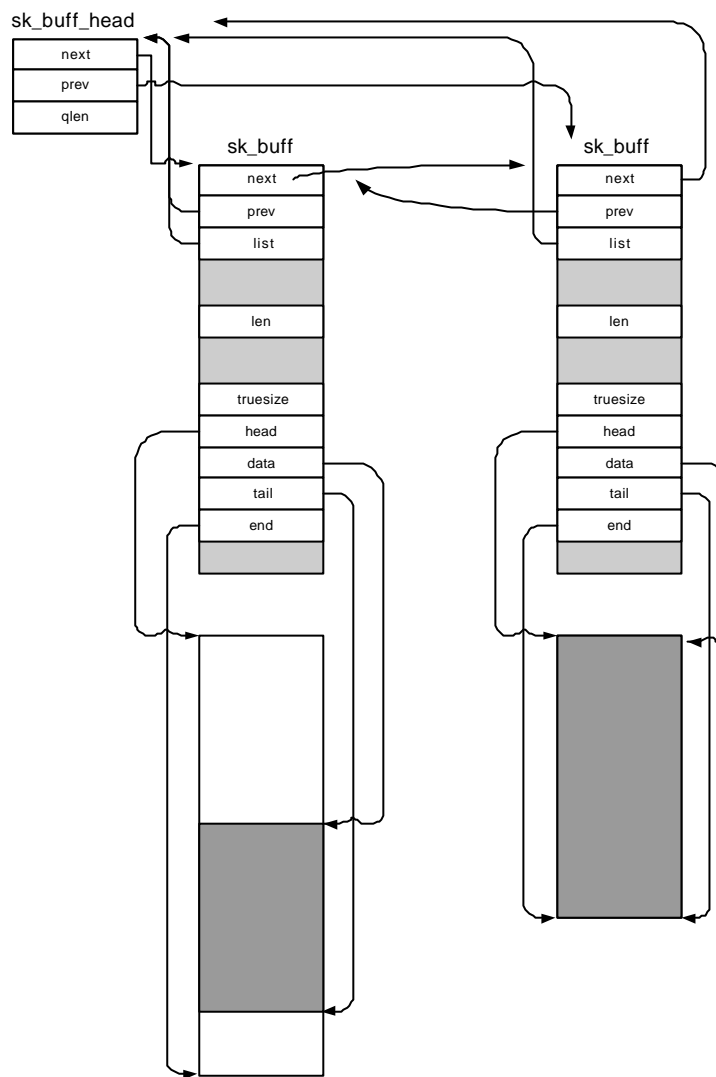
skbuff 가 , 가 inline

가 .

가

가

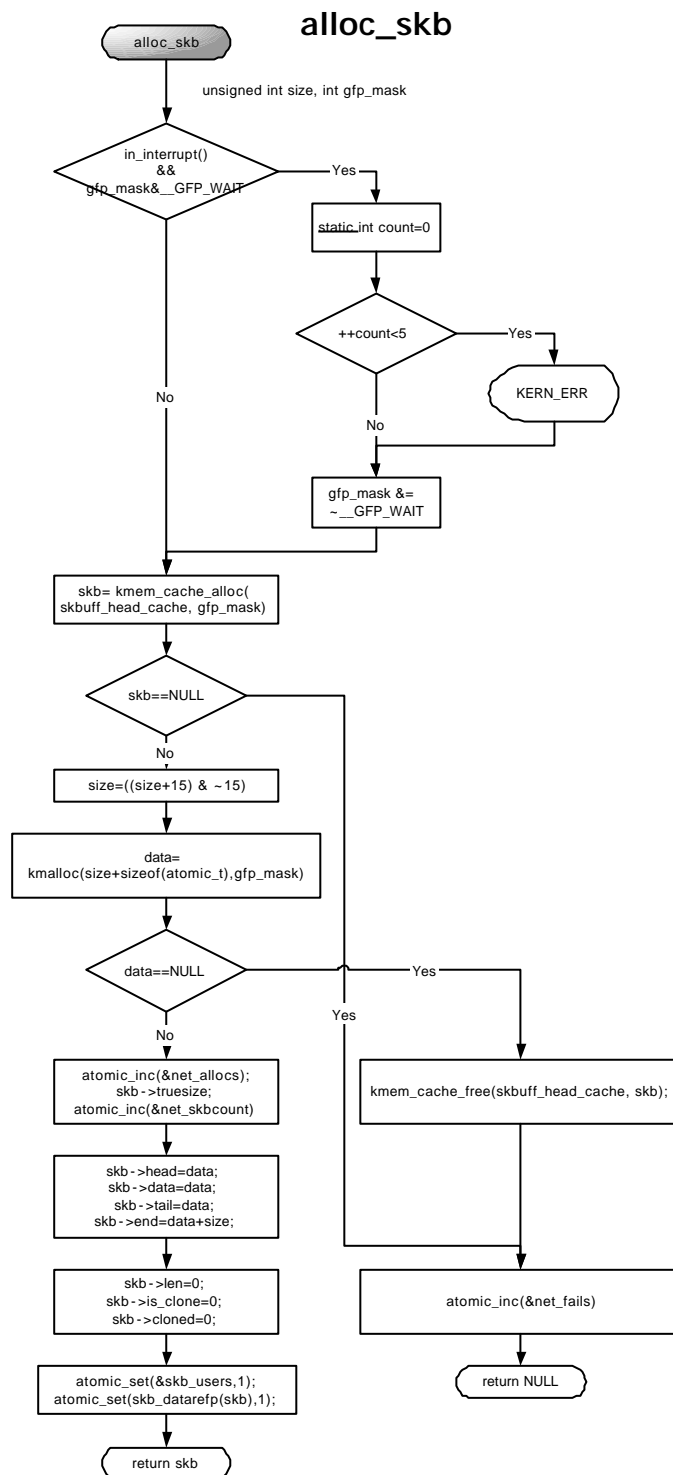
next	prev	skbuff	skbuff	가	.
list	skbuff	header	가	.	.
len				.	.
truesize				.	.
head		가	.	.	.
data		가	.	가	.
tail		가	.	.	.
end		가	.	.	.
skbuff		2.2	.	.	.



2.2 skbuff example

2.2 alloc_skb

alloc_skb() skbuff . size
gfp_mask . net/core/skbuff.c ,



```

119 struct sk_buff *alloc_skb(unsigned int size,int gfp_mask)
120 {
121     struct sk_buff *skb;
122     u8 *data;
123
124     if (in_interrupt() && (gfp_mask & __GFP_WAIT)) {
125         static int count = 0;
126         if (++count < 5) {
127             printk(KERN_ERR "alloc_skb called nonatomically "
128                 "from interrupt %p \n", __builtin_return_address(0));
129         }
130         gfp_mask &= ~__GFP_WAIT;
131     }
132
133     /* Get the HEAD */
134     skb = kmem_cache_alloc(skbuff_head_cache, gfp_mask);
135     if (skb == NULL)
136         goto nohead;

```

124 131 가 ,

.

4 .

. in_interrupt() include/asm-i386/hardirq.h

.

```
#define in_interrupt() ({ int __cpu = smp_processor_id(); \
    (local_irq_count[__cpu] + local_bh_count[__cpu] != 0); })
```

가 in_interrupt()가 true가 . , gfp_mask
__GFP_WAIT .

134 . skbuff_head_cache skbuff.c

kmem_cache_t . kmem_cache_t include/linux/slabb.h struct

kmem_cache_s mm/slabb.c ,

. skbuff_head_cache skb_init()

. skb_init() . 379 struct sk_buff

.

```
376 void __init skb_init(void)
377 {
378     skbuff_head_cache = kmem_cache_create("skbuff_head_cache",

```

```

379                                     sizeof(struct sk_buff),
380                                     0,
381                                     SLAB_HWCACHE_ALIGN,
382                                     skb_headerinit, NULL);
383     if (!skbuff_head_cache)
384         panic("cannot create skbuff cache");
385 }
kmem_cache_alloc                                     , mm/slab.c
.                                     sk_buff                                     skb                                     .

137
138 /* Get the DATA. Size must match skb_add_mtu(). */
139 size = ((size + 15) & ~15);

. 139      4      0      offset      0
.      size      4      가 0
.      10      size      16      가
size가 17      32가      , 1      16      16
.

140 data = kmalloc(size + sizeof(atomic_t), gfp_mask);
141 if (data == NULL)
142     goto noda;
143

140 data      .      atomic_t      가
. atomic_t      include/asm-i386/atomic.h      ,
.

22 #ifdef __SMP__
23 typedef struct { volatile int counter; } atomic_t;
24 #else
25 typedef struct { int counter; } atomic_t;
26 #endif
atomic      .      가

가

144 /* Note that this counter is useless now - you can just look in the
145 * skbuff_head entry in /proc/slabinfo. We keep it only for emergency

```

```

146  * cases.
147  */
148  atomic_inc(&net_allocs);
149

```

```

atomic_inc                                가                                . net_allocs  atomic_t
                                net \ core \ skbuff.c                                . net_allocs
                                가                                , alloc_skb()                                skb_clone()
                                가                                .

```

```

150  skb->truesize = size;
151
152  atomic_inc(&net_skbcount);
153
154  /* Load the data pointers. */
155  skb->head = data;
156  skb->data = data;
157  skb->tail = data;
158  skb->end = data + size;
159

```

```

150  skb                                truesize                                가                                .
152  net_skbcount                                가                                ,                                net_skbcount  sk_buff                                가                                ,
kfree_skb()                                .                                ,                                sk_buff                                .
155  skb                                head,data,tail,end                                .
                                ,                                , tail
                                가                                .

```

```

160  /* Set up other state */
161  skb->len = 0;
162  skb->is_clone = 0;
163  skb->cloned = 0;
164
165  atomic_set(&skb->users, 1);
166  atomic_set(skb_datarefp(skb), 1);
167  return skb;

```

```

                                가                                len=0                                . is_clone                                skbuffer가

```

```

        0 . cloned
        0 . 165 skbuffer users 1
        . skbuffer 가 가,
skb_datarefp(skb) 가 . skbuff.h

extern __inline__ atomic_t *skb_datarefp(struct sk_buff *skb)
{
    return (atomic_t *) (skb->end);
}

skb_datarefp(skb) , 140 atomic_t 가
, 1 . 167 skb

168
169 nodata:
170 kmem_cache_free(skbuff_head_cache, skb);

140 ,

171 nohead:
172 atomic_inc(&net_fails);
173 return NULL;
174 }

net_fails 가 ,
alloc_skb . sk_buffer

```

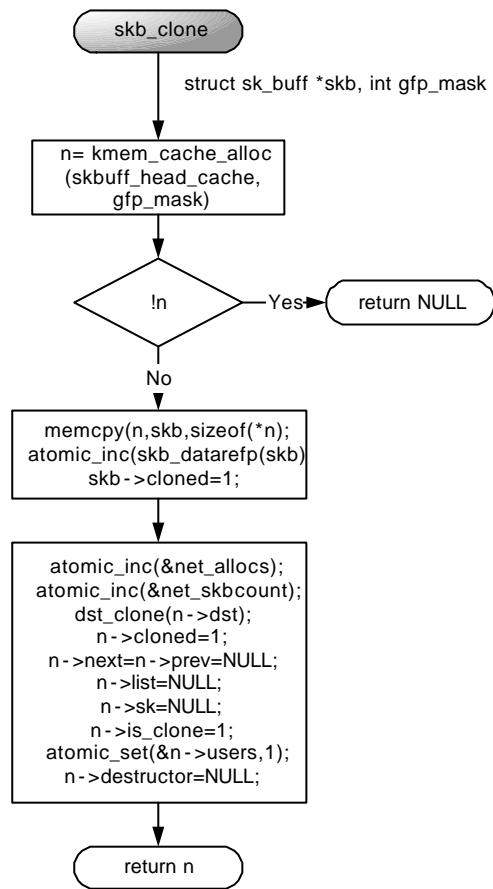
2.3 skb_clone

```

        skb_clone .
, skbuff .

```

net/core/skbuff.c .



2.4 skb_clone

```

235 struct sk_buff *skb_clone(struct sk_buff *skb, int gfp_mask)
236 {
237     struct sk_buff *n;
238
239     n = kmem_cache_alloc(skbuff_head_cache, gfp_mask);
240     if (!n)
241         return NULL;
242
243     memcpy(n, skb, sizeof(*n));
244     atomic_inc(skb->datarefp(skb));
245     skb->cloned = 1;
246
247     atomic_inc(&net_allocs);
248     atomic_inc(&net_skbcount);
  
```

```

239         skbbuffer . 243
sizeof(*n) sizeof(sk_buff)가 , sk_buff n sk_buff
. 244 reference 가 .

```

```

249     dst_clone(n->dst);

```

```

249     include/net/dst.h     dst_clone .
83 extern __inline__
84 struct dst_entry * dst_clone(struct dst_entry * dst)
85 {
86     if (dst)
87         atomic_inc(&dst->use);
88     return dst;
89 }
dst     use     가 . dst_entry
.

```

```

250     n->cloned = 1;
251     n->next = n->prev = NULL;
252     n->list = NULL;
253     n->sk = NULL;
254     n->is_clone = 1;
255     atomic_set(&n->users, 1);
256     n->destructor = NULL;
257     return n;
258 }
259

```

```

skbbuffer가 clone     skbbuffer n     skb     cloned가 1 . 251
252     list     . 254     n
is_clone     1 . 257     skbbuffer n .
clone
.

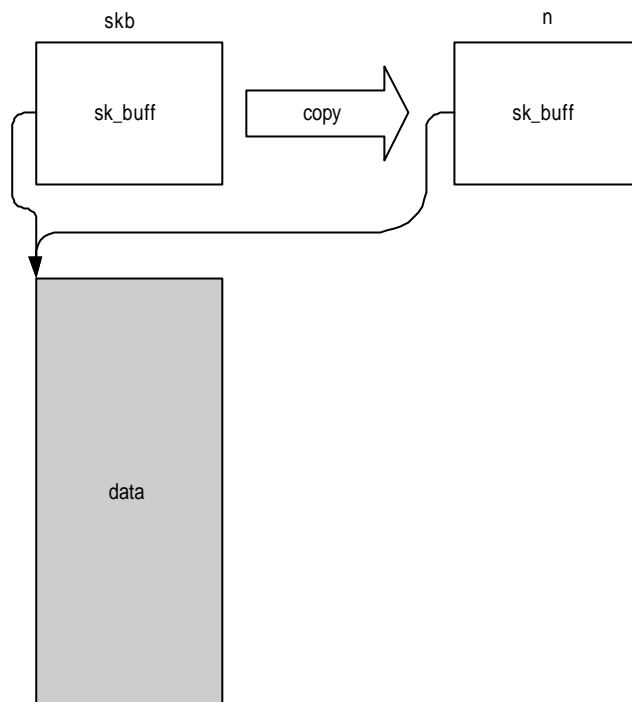
```

2.4 skb_copy

```

skb_copy     skbuff
.

```

3.7 clone

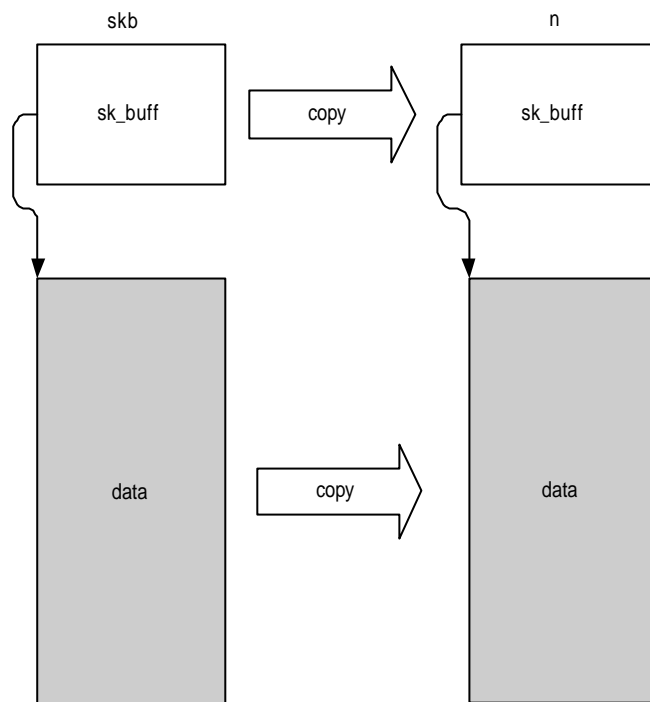
```

264 struct sk_buff *skb_copy(struct sk_buff *skb, int gfp_mask)
265 {
266     struct sk_buff *n;
267     unsigned long offset;
272
273     n=alloc_skb(skb->end - skb->head, gfp_mask);
274     if(n==NULL)
275         return NULL;
276
281     offset=n->head-skb->head;

273     sk_buff      , 281
.

282
283     /* Set the data pointer */
284     skb_reserve(n,skb->data-skb->head);

```



3.8 copy

```

skb_reserve      /include/linux/skbuff.h      data      tail
len
517 extern __inline__ void skb_reserve(struct sk_buff *skb, unsigned int len)
518 {
519     skb->data+=len;
520     skb->tail+=len;
521 }
len      skb      , n

```

```

285     /* Set the tail pointer and length */
286     skb_put(n,skb->len);

```

```

skb_put      skbuff.h
extern __inline__ unsigned char *__skb_put(struct sk_buff *skb, unsigned int len)
{
    unsigned char *tmp=skb->tail;
    skb->tail+=len;
    skb->len+=len;
}

```

```

        return tmp;
    }
    ,      skb      skbuff n      .

287     /* Copy the bytes */
288     memcpy(n->head,skb->head,skb->end-skb->head);

    skb      n      .

289     n->csum = skb->csum;
290     n->list=NULL;
291     n->sk=NULL;
292     n->dev=skb->dev;
293     n->priority=skb->priority;
294     n->protocol=skb->protocol;
295     n->dst=dst_clone(skb->dst);
296     n->h.raw=skb->h.raw+offset;
297     n->nh.raw=skb->nh.raw+offset;
298     n->mac.raw=skb->mac.raw+offset;
299     memcpy(n->cb, skb->cb, sizeof(skb->cb));
300     n->used=skb->used;
301     n->is_clone=0;
302     atomic_set(&n->users, 1);
303     n->pkt_type=skb->pkt_type;
304     n->stamp=skb->stamp;
305     n->destructor = NULL;
306     n->security=skb->security;
307 #ifdef CONFIG_IP_FIREWALL
308     n->fwmark = skb->fwmark;
309 #endif
310     return n;
311 }

    skbbuffer      .      296      298
    ,      ,      가
    ,      ,      281      offset      .

```

2.5 skb_headerinit

skb_headerinit skbuff

```

180 static inline void skb_headerinit(void *p, kmem_cache_t *cache,
181                                     unsigned long flags)
182 {
183     struct sk_buff *skb = p;
184
185     skb->destructor = NULL;
186     skb->pkt_type = PACKET_HOST;      /* Default type */
187     skb->pkt_bridged = 0;             /* Not bridged */
188     skb->prev = skb->next = NULL;
189     skb->list = NULL;
190     skb->sk = NULL;
191     skb->stamp.tv_sec=0;              /* No idea about time */
192     skb->ip_summed = 0;
193     skb->security = 0; /* By default packets are insecure */
194     skb->dst = NULL;
195 #ifdef CONFIG_IP_FIREWALL
196     skb->fwmark = 0;
197 #endif
198     memset(skb->cb, 0, sizeof(skb->cb));
199     skb->priority = 0;
200 }
186     pkt_type          PACKET_HOST          . PACKET_HOST
include/linux/if_packet.h          . packet   type
.
#define PACKET_HOST          0          /* To us          */
#define PACKET_BROADCAST     1          /* To all          */
#define PACKET_MULTICAST     2          /* To group        */
#define PACKET_OTHERHOST     3          /* To someone else */
#define PACKET_OUTGOING      4          /* Outgoing of any type */
/* These ones are invisible by user level */
#define PACKET_LOOPBACK      5          /* MC/BRD frame looped back */
#define PACKET_FASTROUTE     6          /* Fastrouted frame  */

```

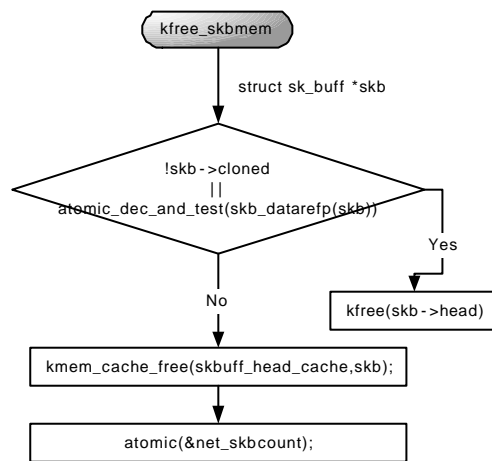
2.6 kfree_skbmem

```

kfree_skbmem      skbbuffer      . kfree_skbmem
                  , __kfree_skbmem      .

202 /*
203  * Free an skbuff by memory without cleaning the state.
204  */
205 void kfree_skbmem(struct sk_buff *skb)
206 {
207     if (!skb->cloned || atomic_dec_and_test(skb_datarefp(skb)))
208         kfree(skb->head);
209
210     kmem_cache_free(skbuff_head_cache, skb);
211     atomic_dec(&net_skbcount);
212 }

```



2.9 kfree_skbmem

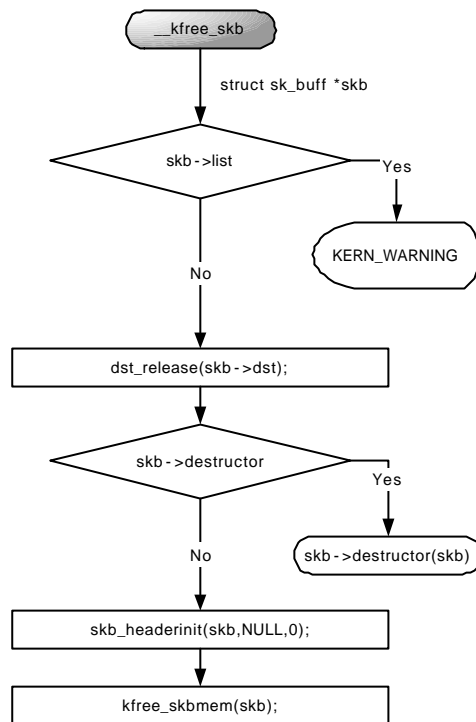
207 , clone , reference 0
, , 가 , kfree()
skbbuffer , net_skbcount .

__kfree_skb

```

214 /*
215  * Free an sk_buff. Release anything attached to the buffer. Clean the state.
216  */
217

```



2.10 __kfree_skb

```

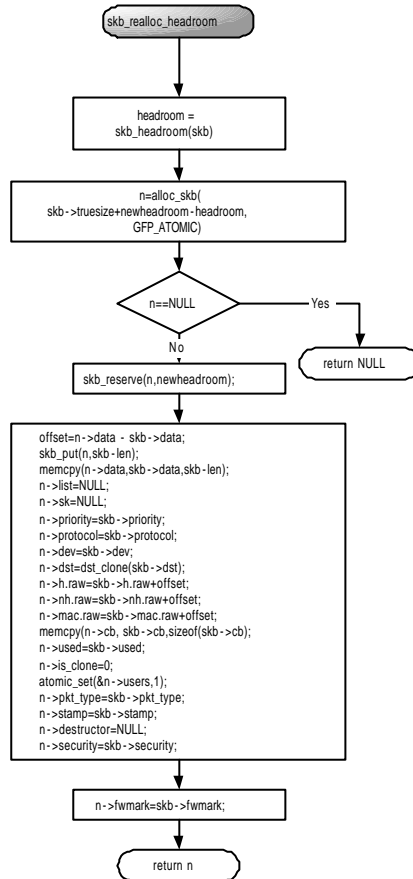
218 void __kfree_skb(struct sk_buff *skb)
219 {
220     if (skb->list)
221         printk(KERN_WARNING "Warning: kfree_skb passed an skb still "
222             "on a list (from %p). \n", __builtin_return_address(0));
223
224     dst_release(skb->dst);
225     if(skb->destructor)
226         skb->destructor(skb);
227     skb_headerinit(skb, NULL, 0); /* clean state */
228     kfree_skbmem(skb);
229 }
  
```

__kfree_skb list가
 가 . 224 dst_release() dst->use
 . destructor가 . skb_headerinit()
 skbuff , kfree_skbmem() .

2.7 skb_realloc_headroom

skb_realloc_headroom

newheadroom



2.11

skb_realloc_headroom

```

313 struct sk_buff *skb_realloc_headroom(struct sk_buff *skb, int newheadroom)
314 {
315     struct sk_buff *n;
316     unsigned long offset;
317     int headroom = skb_headroom(skb);
318

```

317 skb_headroom()

skbuff.h

```

,
extern __inline__ int skb_headroom(struct sk_buff *skb)
{
    return skb->data-skb->head;
}

head    data                headroom    .

319     /*
320     *      Allocate the copy buffer
321     */
322
323     n=alloc_skb(skb->truesize+newheadroom-headroom, GFP_ATOMIC);
324     if(n==NULL)
325         return NULL;

newheadroom    skbuffer    .

326
327     skb_reserve(n,newheadroom);
333     offset=n->data-skb->data;
336     skb_put(n,skb->len);
338     memcpy(n->data,skb->data,skb->len);
339     n->list=NULL;
340     n->sk=NULL;
341     n->priority=skb->priority;
342     n->protocol=skb->protocol;
343     n->dev=skb->dev;
344     n->dst=dst_clone(skb->dst);
345     n->h.raw=skb->h.raw+offset;
346     n->nh.raw=skb->nh.raw+offset;
347     n->mac.raw=skb->mac.raw+offset;
348     memcpy(n->cb, skb->cb, sizeof(skb->cb));
349     n->used=skb->used;
350     n->is_clone=0;
351     atomic_set(&n->users, 1);
352     n->pkt_type=skb->pkt_type;
353     n->stamp=skb->stamp;
354     n->destructor = NULL;
355     n->security=skb->security;

```



```

356 #ifdef CONFIG_IP_FIREWALL
357     n->fwmark = skb->fwmark;
358 #endif
359
360     return n;
361 }

```

skbuffer n

2.8

3 skbuff.c

```

87 void skb_over_panic(struct sk_buff *skb, int sz, void *here)
88 {
89     panic("skput:over: %p:%d put:%d dev:%s",
90         here, skb->len, sz, skb->dev ? skb->dev->name : "<NULL>");
91 }
92
93 void skb_under_panic(struct sk_buff *skb, int sz, void *here)
94 {
95     panic("skput:under: %p:%d put:%d dev:%s",
96         here, skb->len, sz, skb->dev ? skb->dev->name : "<NULL>");
97 }
98
99 void show_net_buffers(void)
100 {
101     printk("Networking buffers in use          : %u \n",
102         atomic_read(&net_skbcount));
103     printk("Total network buffer allocations    : %u \n",
104         atomic_read(&net_allocs));
105     printk("Total failed network buffer allocs : %u \n",
106         atomic_read(&net_fails));
107 #ifdef CONFIG_INET
108     printk("IP fragment buffer size              : %u \n",
109         atomic_read(&ip_frag_mem));
110 #endif
111 }

```

2.9

skbuff.h

```
164 extern __inline__ atomic_t *skb_datarefp(struct sk_buff *skb)
165 {
166     return (atomic_t *) (skb->end);
167 }
```

skb_datarefp()는 skb->end를 atomic_t *로 반환한다.

reference: [skb_datarefp\(\) - skbuff.h](#)

```
169 extern __inline__ int skb_queue_empty(struct sk_buff_head *list)
170 {
171     return (list->next == (struct sk_buff *) list);
172 }
```

skb_queue_empty()는 sk_buff_head list의 next가 list 자체를 가리키는 경우, empty (=1)을 반환한다.

```
174 extern __inline__ void kfree_skb(struct sk_buff *skb)
175 {
176     if (atomic_dec_and_test(&skb->users))
177         __kfree_skb(skb);
178 }
```

skb->users가 0이 되면 __kfree_skb()를 호출하여 skb를 해제한다.

```
180 /* Use this if you didn't touch the skb state [for fast switching] */
181 extern __inline__ void kfree_skb_fast(struct sk_buff *skb)
182 {
183     if (atomic_dec_and_test(&skb->users))
184         kfree_skbmem(skb);
185 }
```

skb->users가 0이 되면, kfree_skbmem()을 호출하여 skb를 해제한다.

```
187 extern __inline__ int skb_cloned(struct sk_buff *skb)
188 {
189     return skb->cloned && atomic_read(skb_datarefp(skb)) != 1;
190 }
```

skb->cloned가 clone

```

192 extern __inline__ int skb_shared(struct sk_buff *skb)
193 {
194     return (atomic_read(&skb->users) != 1);
195 }
skb

204 extern __inline__ struct sk_buff *skb_unshare(struct sk_buff *skb, int pri)
205 {
206     struct sk_buff *nskb;
207     if(!skb_cloned(skb))
208         return skb;
209     nskb=skb_copy(skb, pri);
210     kfree_skb(skb);          /* Free our shared copy */
211     return nskb;
212 }
clone    skb    copy    ,    .    tcpdump
forward

221 extern __inline__ struct sk_buff *skb_peek(struct sk_buff_head *list_)
222 {
223     struct sk_buff *list = ((struct sk_buff *)list_->next;
224     if (list == (struct sk_buff *)list_)
225         list = NULL;
226     return list;
227 }
list    next가 가    sk_buff    .    sk_buff    .

229 extern __inline__ struct sk_buff *skb_peek_tail(struct sk_buff_head *list_)
230 {
231     struct sk_buff *list = ((struct sk_buff *)list_->prev;
232     if (list == (struct sk_buff *)list_)
233         list = NULL;
234     return list;
235 }
list    sk_buff    .

241 extern __inline__ __u32 skb_queue_len(struct sk_buff_head *list_)
242 {
243     return(list_->qlen);

```

```
244 }
```

list .

```
246 extern __inline__ void skb_queue_head_init(struct sk_buff_head *list)
```

```
247 {
```

```
248     list->prev = (struct sk_buff *)list;
```

```
249     list->next = (struct sk_buff *)list;
```

```
250     list->qlen = 0;
```

```
251 }
```

skbuffer list . list가 , prev next list 가

```
260 extern __inline__ void __skb_queue_head(struct sk_buff_head *list, struct sk_buff *newsk)
```

```
261 {
```

```
262     struct sk_buff *prev, *next;
```

```
263
```

```
264     newsk->list = list;
```

```
265     list->qlen++;
```

```
266     prev = (struct sk_buff *)list;
```

```
267     next = prev->next;
```

```
268     newsk->next = next;
```

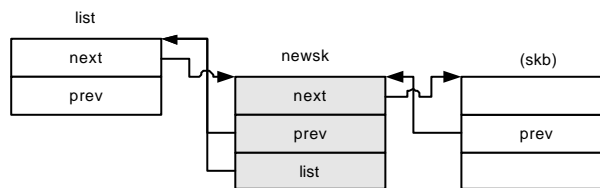
```
269     newsk->prev = prev;
```

```
270     next->prev = newsk;
```

```
271     prev->next = newsk;
```

```
272 }
```

skbuff list skbuff .



```
274 extern spinlock_t skb_queue_lock;
```

```
276 extern __inline__ void skb_queue_head(struct sk_buff_head *list, struct sk_buff *newsk)
```

```
277 {
```

```
278     unsigned long flags;
```

```
279
```

```

280     spin_lock_irqsave(&skb_queue_lock, flags);
281     __skb_queue_head(list, newsk);
282     spin_unlock_irqrestore(&skb_queue_lock, flags);
283 }
skb_queue_head    __skb_queue_head    , lock    atomic

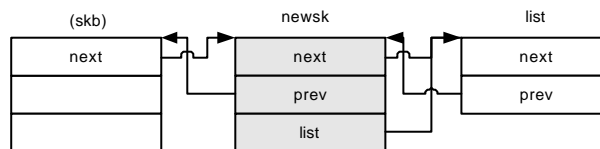
```

```

289 extern __inline__ void __skb_queue_tail(struct sk_buff_head *list, struct sk_buff *newsk)
290 {
291     struct sk_buff *prev, *next;
292
293     newsk->list = list;
294     list->qlen++;
295     next = (struct sk_buff *)list;
296     prev = next->prev;
297     newsk->next = next;
298     newsk->prev = prev;
299     next->prev = newsk;
300     prev->next = newsk;
301 }

```

tail skbuff



```

303 extern __inline__ void skb_queue_tail(struct sk_buff_head *list, struct sk_buff *newsk)
304 {
305     unsigned long flags;
306
307     spin_lock_irqsave(&skb_queue_lock, flags);
308     __skb_queue_tail(list, newsk);
309     spin_unlock_irqrestore(&skb_queue_lock, flags);
310 }
__skb_queue_tail    atomic

```

```

316 extern __inline__ struct sk_buff *__skb_dequeue(struct sk_buff_head *list)
317 {

```

```

318     struct sk_buff *next, *prev, *result;
319
320     prev = (struct sk_buff *) list;
321     next = prev->next;
322     result = NULL;
323     if (next != prev) {
324         result = next;
325         next = next->next;
326         list->qlen--;
327         next->prev = prev;
328         prev->next = next;
329         result->next = NULL;
330         result->prev = NULL;
331         result->list = NULL;
332     }
333     return result;
334 }
list      head      skbuff      .

336 extern __inline__ struct sk_buff *skb_dequeue(struct sk_buff_head *list)
337 {
338     long flags;
339     struct sk_buff *result;
340
341     spin_lock_irqsave(&skb_queue_lock, flags);
342     result = __skb_dequeue(list);
343     spin_unlock_irqrestore(&skb_queue_lock, flags);
344     return result;
345 }
__skb_dequeue()      atomic      .

351 extern __inline__ void __skb_insert(struct sk_buff *newsk,
352     struct sk_buff * prev, struct sk_buff *next,
353     struct sk_buff_head * list)
354 {
355     newsk->next = next;
356     newsk->prev = prev;
357     next->prev = newsk;
358     prev->next = newsk;
359     newsk->list = list;
360     list->qlen++;
361 }

```

prev next newsk .

```
366 extern __inline__ void skb_insert(struct sk_buff *old, struct sk_buff *newsk)
367 {
368     unsigned long flags;
369
370     spin_lock_irqsave(&skb_queue_lock, flags);
371     __skb_insert(newsk, old->prev, old, old->list);
372     spin_unlock_irqrestore(&skb_queue_lock, flags);
373 }
```

old newsk atomic .

```
379 extern __inline__ void __skb_append(struct sk_buff *old, struct sk_buff *newsk)
380 {
381     __skb_insert(newsk, old, old->next, old->list);
382 }
```

old newsk 가 .

```
384 extern __inline__ void skb_append(struct sk_buff *old, struct sk_buff *newsk)
385 {
386     unsigned long flags;
387
388     spin_lock_irqsave(&skb_queue_lock, flags);
389     __skb_append(old, newsk);
390     spin_unlock_irqrestore(&skb_queue_lock, flags);
391 }
```

old newsk 가 atomic .

```
397 extern __inline__ void __skb_unlink(struct sk_buff *skb, struct sk_buff_head *list)
398 {
399     struct sk_buff * next, * prev;
400
401     list->qlen - -;
402     next = skb->next;
403     prev = skb->prev;
404     skb->next = NULL;
405     skb->prev = NULL;
406     skb->list = NULL;
407     next->prev = prev;
408     prev->next = next;
409 }
```

list skb .

```
418 extern __inline__ void skb_unlink(struct sk_buff *skb)
```

```

419 {
420     unsigned long flags;
421
422     spin_lock_irqsave(&skb_queue_lock, flags);
423     if(skb->list)
424         __skb_unlink(skb, skb->list);
425     spin_unlock_irqrestore(&skb_queue_lock, flags);
426 }
__skb_unlink    atomic    .

428 /* XXX: more streamlined implementation */
429 extern __inline__ struct sk_buff *__skb_dequeue_tail(struct sk_buff_head *list)
430 {
431     struct sk_buff *skb = skb_peek_tail(list);
432     if (skb)
433         __skb_unlink(skb, list);
434     return skb;
435 }
list    tail    skb    dequeue    .

437 extern __inline__ struct sk_buff *skb_dequeue_tail(struct sk_buff_head *list)
438 {
439     long flags;
440     struct sk_buff *result;
441
442     spin_lock_irqsave(&skb_queue_lock, flags);
443     result = __skb_dequeue_tail(list);
444     spin_unlock_irqrestore(&skb_queue_lock, flags);
445     return result;
446 }
__skb_dequeue_tail    atomic    .

452 extern __inline__ unsigned char *__skb_put(struct sk_buff *skb, unsigned int len)
453 {
454     unsigned char *tmp=skb->tail;
455     skb->tail+=len;
456     skb->len+=len;
457     return tmp;
458 }

len    .

460 extern __inline__ unsigned char *skb_put(struct sk_buff *skb, unsigned int len)
461 {
462     unsigned char *tmp=skb->tail;

```



```

463     skb->tail+=len;
464     skb->len+=len;
465     if(skb->tail>skb->end)
466     {
467         __label__ here;
468         skb_over_panic(skb, len, &&here);
469 here:
470     }
471     return tmp;
472 }
__skb_put      tail    end      skb_over_panic      .

```

```

474 extern __inline__ unsigned char *__skb_push(struct sk_buff *skb, unsigned int len)
475 {
476     skb->data-=len;
477     skb->len+=len;
478     return skb->data;
479 }
len
.

```

```

481 extern __inline__ unsigned char *__skb_push(struct sk_buff *skb, unsigned int len)
482 {
483     skb->data-=len;
484     skb->len+=len;
485     if(skb->data<skb->head)
486     {
487         __label__ here;
488         skb_under_panic(skb, len, &&here);
489 here:
490     }
491     return skb->data;
492 }
__skb_push      , skb_under_panic      .

```

```

494 extern __inline__ char *__skb_pull(struct sk_buff *skb, unsigned int len)
495 {
496     skb->len-=len;
497     return  skb->data+=len;
498 }
len
.

```

```

500 extern __inline__ unsigned char * skb_pull(struct sk_buff *skb, unsigned int len)
501 {
502     if (len > skb->len)

```

```

503         return NULL;
504     return __skb_pull(skb, len);
505 }
__skb_pull                len

507 extern __inline__ int skb_headroom(struct sk_buff *skb)
508 {
509     return skb->data-skb->head;
510 }

.

512 extern __inline__ int skb_tailroom(struct sk_buff *skb)
513 {
514     return skb->end-skb->tail;
515 }

.

517 extern __inline__ void skb_reserve(struct sk_buff *skb, unsigned int len)
518 {
519     skb->data+=len;
520     skb->tail+=len;
521 }

len

523 extern __inline__ void __skb_trim(struct sk_buff *skb, unsigned int len)
524 {
525     skb->len = len;
526     skb->tail = skb->data+len;
527 }

len

529 extern __inline__ void skb_trim(struct sk_buff *skb, unsigned int len)
530 {
531     if (skb->len > len) {
532         __skb_trim(skb, len);
533     }
534 }
__skb_trim                skb->len    len    .

536 extern __inline__ void skb_orphan(struct sk_buff *skb)
537 {
538     if (skb->destructor)
539         skb->destructor(skb);

```

```

540     skb->destructor = NULL;
541     skb->sk = NULL;
542 }
skb          sk    NULL

544 extern __inline__ void skb_queue_purge(struct sk_buff_head *list)
545 {
546     struct sk_buff *skb;
547     while ((skb=skb_dequeue(list))!=NULL)
548         kfree_skb(skb);
549 }
list

551 extern __inline__ struct sk_buff *dev_alloc_skb(unsigned int length)
552 {
553     struct sk_buff *skb;
554
555     skb = alloc_skb(length+16, GFP_ATOMIC);
556     if (skb)
557         skb_reserve(skb,16);
558     return skb;
559 }
          16          sk_buff          ,          16
          .          16          .

561 extern __inline__ struct sk_buff *
562 skb_cow(struct sk_buff *skb, unsigned int headroom)
563 {
564     headroom = (headroom+15)&~15;
565
566     if ((unsigned)skb_headroom(skb) < headroom || skb_cloned(skb)) {
567         struct sk_buff *skb2 = skb_realloc_headroom(skb, headroom);
568         kfree_skb(skb);
569         skb = skb2;
570     }
571     return skb;
572 }
          16          headroom          .

```

3. Device driver

3.1 device driver

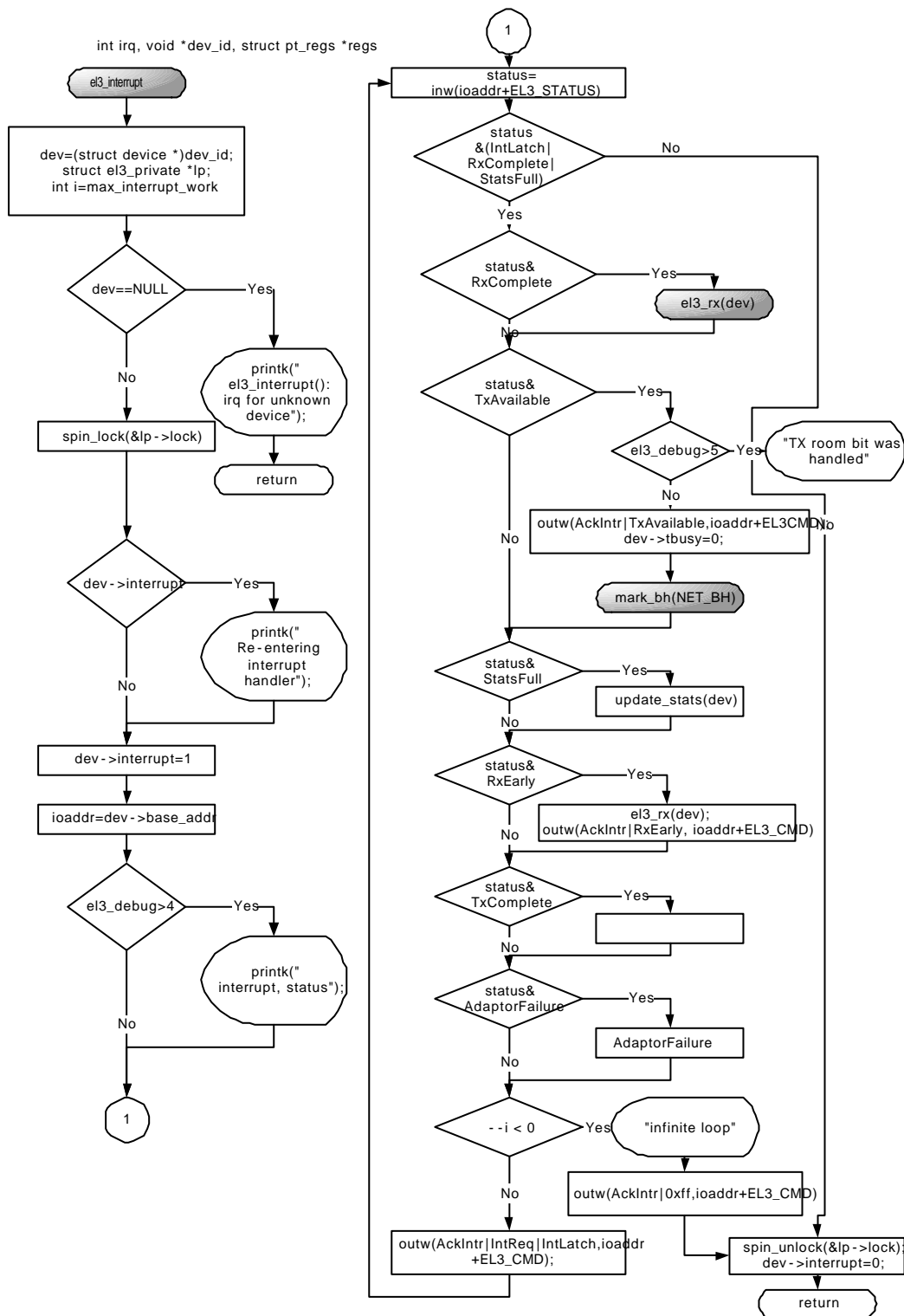
3c509 device driver . driver
가 . 3c509
driver drivers/net/3c509.c .

3.2 el3_interrupt

(NIC) 가 ,
가 , RxComplete
 , TxAvailable .
 .

```
614 /* The EL3 interrupt handler. */
615 static void
616 el3_interrupt(int irq, void *dev_id, struct pt_regs *regs)
617 {
618     struct device *dev = (struct device *)dev_id;
619     struct el3_private *lp;
620     int ioaddr, status;
621     int i = max_interrupt_work;
622
623     if (dev == NULL) {
624         printk ("el3_interrupt(): irq %d for unknown device. \n", irq);
625         return;
626     }
627
628     lp = (struct el3_private *)dev->priv;
629     spin_lock(&lp->lock);
```

device . struct device
include/linux/netdevice.h ,
가 . 621



3.1 el3_interrupt()

```
max_interrupt_work    10                . 628                3c509    가
    가 lp                . 629    spin_lock    include/asm-i386/spinlock.h
```

```
#define spin_lock_init(lock)    do { } while(0)
#define spin_lock(lock)        (void)(lock) /* Not "unused variable". */
#define spin_trylock(lock)    (1)
#define spin_unlock_wait(lock) do { } while(0)
#define spin_unlock(lock)      do { } while(0)
#define spin_lock_irq(lock)    cli()
#define spin_unlock_irq(lock)  sti()
```

```
    CPU    spin_lock
```

```
631    if (dev->interrupt)
632        printk("%s: Re-entering the interrupt handler. \n", dev->name);
633    dev->interrupt = 1;
634
635    ioaddr = dev->base_addr;
636
```

```
dev->interrupt=1                가                . 635
    io base                가                .
```

```
637    if (el3_debug > 4) {
638        status = inw(ioaddr + EL3_STATUS);
639        printk("%s: interrupt, status %4.4x. \n", dev->name, status);
640    }
641
642    while ((status = inw(ioaddr + EL3_STATUS)) &
643        (IntLatch | RxComplete | StatsFull)) {
644
```

```
    while                IntLatch, RxComplete,    StatsFull
```

```
645        if (status & RxComplete)
646            el3_rx(dev);
647
```

RxComplete

el3_rx()

```
648         if (status & TxAvailable) {
649             if (el3_debug > 5)
650                 printk("    TX room bit was handled. \n");
651             /* There's room in the FIFO for a full-sized packet. */
652             outw(AckIntr | TxAvailable, iaddr + EL3_CMD);
653             dev->tbusy = 0;
654             mark_bh(NET_BH);
655         }
```

648 655 가 가 mark_bh(NET_BH) Bottom Half
Handler . mark_bh() Bottom Half Handler

```
656         if (status & (AdapterFailure | RxEarly | StatsFull | TxComplete)) {
657             /* Handle all uncommon interrupts. */
658             if (status & StatsFull)                                     /* Empty
y statistics. */
659                 update_stats(dev);
660             if (status & RxEarly) {                                     /* Rx ea
rly is unused. */
661                 el3_rx(dev);
662                 outw(AckIntr | RxEarly, iaddr + EL3_CMD);
663             }
664             if (status & TxComplete) {                                 /* Really Tx error.
*/
665                 struct el3_private *lp = (struct el3_private *)dev->priv;
666                 short tx_status;
667                 int i = 4;
668
669                 while (--i>0 && (tx_status = inb(iaddr + TX_STATUS)) >
0) {
670                     if (tx_status & 0x38) lp->stats.tx_aborted_errors+
+;
671                     if (tx_status & 0x30) outw(TxReset, iaddr + EL3_
CMD);
672                     if (tx_status & 0x3C) outw(TxEnable, iaddr + EL
```

```

3_CMD);
673                                     outb(0x00, ioaddr + TX_STATUS); /* Pop the stat
us stack. */
674                                     }
675                                     }
676                                     if (status & AdapterFailure) {
677                                         /* Adapter failure requires Rx reset and reinit. */
678                                         outw(RxReset, ioaddr + EL3_CMD);
679                                         /* Set the Rx filter to the current state. */
680                                         outw(SetRxFilter | RxStation | RxBroadcast
681                                             | (dev->flags & IFF_ALLMULTI ? RxMulticast :
0)
682                                             | (dev->flags & IFF_PROMISC ? RxProm : 0),
683                                             ioaddr + EL3_CMD);
684                                         outw(RxEnable, ioaddr + EL3_CMD); /* Re-enable the recei
ver. */
685                                         outw(AckIntr | AdapterFailure, ioaddr + EL3_CMD);
686                                     }
687                                     }
688
656         688         StatsFull, RxEarly,TxComplete,AdaptorFailure
.
.
.
689         if (--i < 0) {
690             printk("%s: Infinite loop in interrupt, status %4.4x. \n",
691                 dev->name, status);
692             /* Clear all interrupts. */
693             outw(AckIntr | 0xFF, ioaddr + EL3_CMD);
694             break;
695         }

        621         i=max_interrupt_work, while i
.         10
.

696         /* Acknowledge the IRQ. */
697         outw(AckIntr | IntReq | IntLatch, ioaddr + EL3_CMD); /* Ack IRQ */
698     }

```


697 가 . 698

```
brace   while   .
```

```

699
700     if (el3_debug > 4) {
701         printk("%s: exiting interrupt, status %4.4x. \n", dev->name,
702               inw(ioaddr + EL3_STATUS));
703     }
704     spin_unlock(&lp->lock);
705     dev->interrupt = 0;
706     return;
707 }

```

704 629 lock .
 interrupt=0 가

3.3 el3_rx

el3_rx() 가 , el3_interrupt()

```

761 static int
762 el3_rx(struct device *dev)
763 {
764     struct el3_private *lp = (struct el3_private *)dev->priv;
765     int ioaddr = dev->base_addr;
766     short rx_status;
767
768     if (el3_debug > 5)
769         printk("    In rx_packet(), status %4.4x, rx_status %4.4x. \n",
770               inw(ioaddr+EL3_STATUS), inw(ioaddr+RX_STATUS));
771     while ((rx_status = inw(ioaddr + RX_STATUS)) > 0) {

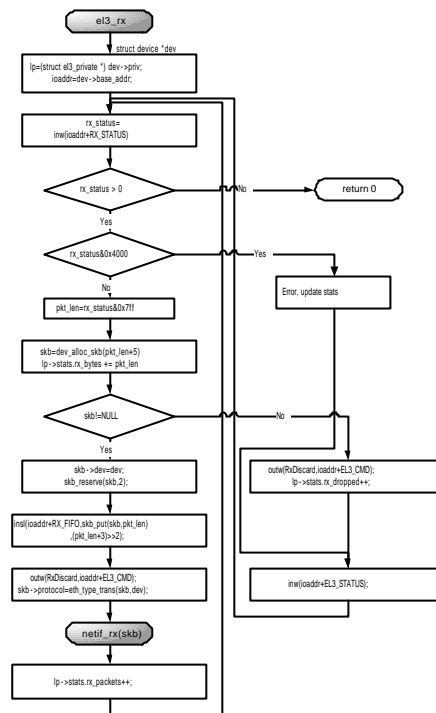
```

RX

```

772         if (rx_status & 0x4000) { /* Error, update stats. */
773             short error = rx_status & 0x3800;
774

```



3.2 el3_rx

```

775 outw(RxDiscard, iaddr + EL3_CMD);
776 lp->stats.rx_errors++;
777 switch (error) {
778     case 0x0000: lp->stats.rx_over_errors++; break;
779     case 0x0800: lp->stats.rx_length_errors++; break;
780     case 0x1000: lp->stats.rx_frame_errors++; break;
781     case 0x1800: lp->stats.rx_length_errors++; break;
782     case 0x2000: lp->stats.rx_frame_errors++; break;
783     case 0x2800: lp->stats.rx_crc_errors++; break;
784 }

```

772 784 .

```

785     } else {
786         short pkt_len = rx_status & 0x7ff;
787         struct sk_buff *skb;
788
789         skb = dev_alloc_skb(pkt_len+5);
790         lp->stats.rx_bytes += pkt_len;

```

```

786 rx_status 11
    pkt_len 789 skbbuffer
    5 790 가
.

791 if (el3_debug > 4)
792     printk("Receiving packet size %d status %4.4x. \n",
793         pkt_len, rx_status);
794 if (skb != NULL) {
795     skb->dev = dev;
796     skb_reserve(skb, 2); /* Align IP on 16 byte */
797
798     /* 'skb->data' points to the start of sk_buff data area. */
799 #ifdef __powerpc__
800     insl_unswapped(ioaddr+RX_FIFO, skb_put(skb,pkt_len),
801         (pkt_len + 3) >> 2);
802 #else
803     insl(ioaddr + RX_FIFO, skb_put(skb,pkt_len),
804         (pkt_len + 3) >> 2);
805 #endif
806
807     outw(RxDiscard, ioaddr + EL3_CMD); /* Pop top Rx packe
t. */
808     skb->protocol = eth_type_trans(skb,dev);
809     netif_rx(skb);
810     lp->stats.rx_packets++;
811     continue;
812 }

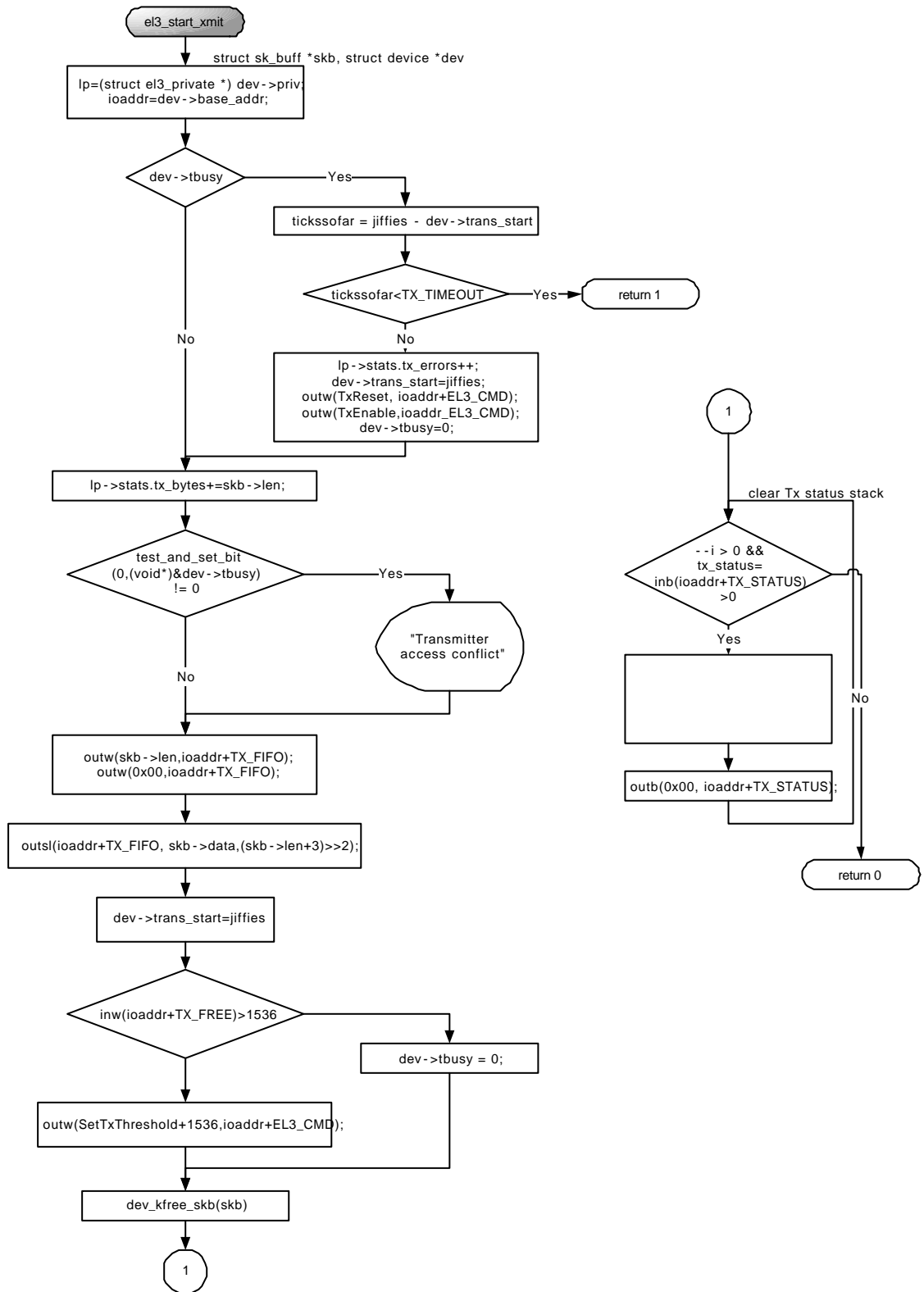
795 skb dev 796
    2 가 789 dev_alloc_skb 796
skb_reserve 16+2, 18 803
    skb insl
    in() string long(4 )
    >>2 가 4
808 eth_type_trans() net/ethernet/eth.c

```


3.4 el3_start_xmit

el3_start_xmit()

.



```

507 static int
508 el3_start_xmit(struct sk_buff *skb, struct device *dev)
509 {
510     struct el3_private *lp = (struct el3_private *)dev->priv;
511     int ioaddr = dev->base_addr;
512
513     /* Transmitter timeout, serious problems. */
514     if (dev->tbusy) {
515         int ticksssofar = jiffies - dev->trans_start;
516         if (ticksssofar < TX_TIMEOUT)
517             return 1;
518         printk("%s: transmit timed out, Tx_status %2.2x status %4.4x "
519             "Tx FIFO room %d. \n",
520             dev->name, inb(ioaddr + TX_STATUS), inw(ioaddr + EL3_STATUS),
521             inw(ioaddr + TX_FREE));
522         lp->stats.tx_errors++;
523         dev->trans_start = jiffies;
524         /* Issue TX_RESET and TX_START commands. */
525         outw(TxReset, ioaddr + EL3_CMD);
526         outw(TxEnable, ioaddr + EL3_CMD);
527         dev->tbusy = 0;
528     }

```

tbusy가
tbusy=0
. jiffies

include/linux/sched.h

```

529
530     lp->stats.tx_bytes += skb->len;
531
532     if (el3_debug > 4) {
533         printk("%s: el3_start_xmit(length = %u) called, status %4.4x. \n",
534             dev->name, skb->len, inw(ioaddr + EL3_STATUS));
535     }

```

530

```

555     /* Avoid timer-based retransmission conflicts. */

```

```

556 if (test_and_set_bit(0, (void*)&dev->tbusy) != 0)
557     printk("%s: Transmitter access conflict. \n", dev->name);

556     dev->tbusy      1      .      test_and_set_bit
include/asm-i386/bitops.h      ,      1      ,
.

558 else {
559     /* Put out the doubleword header... */
560     outw(skb->len, ioaddr + TX_FIFO);
561     outw(0x00, ioaddr + TX_FIFO);
562     outsl(ioaddr + TX_FIFO, skb->data, (skb->len + 3) >> 2);

        skb->data      .      4
(skb->len+3)>>2      .

563     dev->trans_start = jiffies;
564     if (inw(ioaddr + TX_FREE) > 1536) {
565         dev->tbusy = 0;
566     } else
567         /* Interrupt us when the FIFO has room for max-sized packet. */
568         outw(SetTxThreshold + 1536, ioaddr + EL3_CMD);
569 }

        dev->tbusy      0      ,      ,      가
.

570 dev_kfree_skb (skb);

skbuffer      .      dev_kfree_skb()      kfree_skb()      skbuff.h
.

571 /* Clear the Tx status stack. */
572 {
573     short tx_status;
574     int i = 4;
575
576     while (--i > 0      &&      (tx_status = inb(ioaddr + TX_STATUS)) > 0) {
577         if (tx_status & 0x38) lp->stats.tx_aborted_errors++;
578         if (tx_status & 0x30) outw(TxReset, ioaddr + EL3_CMD);

```



```

607             if (tx_status & 0x3C) outw(TxEnable, ioaddr + EL3_CMD);
608             outb(0x00, ioaddr + TX_STATUS); /* Pop the status stack. */
609         }
610     }
611     return 0;
612 }

```

```

599         610             Tx status stack
        . 611             가

```

4. netdevice

4.1 netdevice

device
net/core/dev.c , device include/linux/netdevice.h

4.2 device

device 가
include/linux/netdevice.h

char	*name		
unsigned long	rmem_end		
unsigned long	rmem_start		
unsigned long	mem_end		
unsigned long	mem_start		
unsigned long	base_addr		IO
unsigned int	irq		IRQ
volatile unsigned char	start		
unsigned long	interrupt		
unsigned long	tbusy		
struct device	*next		가
int	(*init)(struct device *dev)		
void	(*destructor)(struct device *dev)		
int	ifindex		index
int	iflink		link
unsigned char	if_port		
unsigned char	dma		DMA
struct net_device_stats*	(*get_stats)(struct device *dev)		
struct iw_statistics*	(*get_wireless_stats)(struct device *dev)		wireless
unsigned long	trans_start		
			(jiffies)
unsigned long	last_rx		
unsigned short	flags		flag
unsigned short	gflags		gflag
unsigned	mtu		MTU
unsigned short	type		
			hardware type
unsigned short	hard_header_len		hardware
void	*priv		
			가 .
unsigned char	broadcast[MAX_ADDR_LEN]		broadcast
unsigned char	pad	dev_addr	8
			pad

unsigned char	dev_addr[MAX_ADDR_LEN]	
unsigned char	addr_len	
struct dev_mc_list	*mc_list	Multicast MAC
int	mc_count	Multicast
int	promiscuity	
int	allmulti	
unsigned long	pkt_queue	packet
struct device	*slave	slave device
void	*atalk_ptr	AppleTalk link
void	*ip_ptr	IPv4
void	*dn_ptr	Decnet
struct Qdisc	*qdisc	Q discipline
struct Qdisc	*qdisc_sleeping	Q discipline sleeping
struct Qdisc	*qdisc_list	Q discipline list
unsigned long	tx_queue_len	
int	bridge_port_id	
int	(*open)(struct device *dev)	open
int	(*stop)(struct device *dev)	stop
int	(*hard_start_xmit)(struct sk_buff *skb, struct device *dev)	
int	(*hard_header)(struct sk_buff *skb, struct device *dev, unsigned short type, void *daddr, void *saddr, unsigned len)	
int	(*rebuild_header)(struct sk_buff *skb)	
void	(*set_multicast_list)(struct device *dev)	Multicast
int	(*set_mac_address)(struct device *dev, void *daddr)	MAC
int	(*do_ioctl)(struct device *dev, struct ifreq *ifr, int cmd)	io
int	(*set_config)(struct device *dev, struct ifmap *map)	
int	(*hard_header_cache)(struct neighbour *neigh, struct hh_cache *hh)	
void	(*header_cache_update)(struct hh_cache *hh, struct device *dev, unsigned char *haddr)	
int	(*change_mtu)(struct device *dev, int new_mtu)	MTU
int	(*hard_header_parse)(struct sk_buff *skb, unsigned char *haddr)	
int	(*neigh_setup)(struct device *dev, struct neighbour *parms)	
int	(*accept_fastpath)(struct device *, struct dst_entry *)	

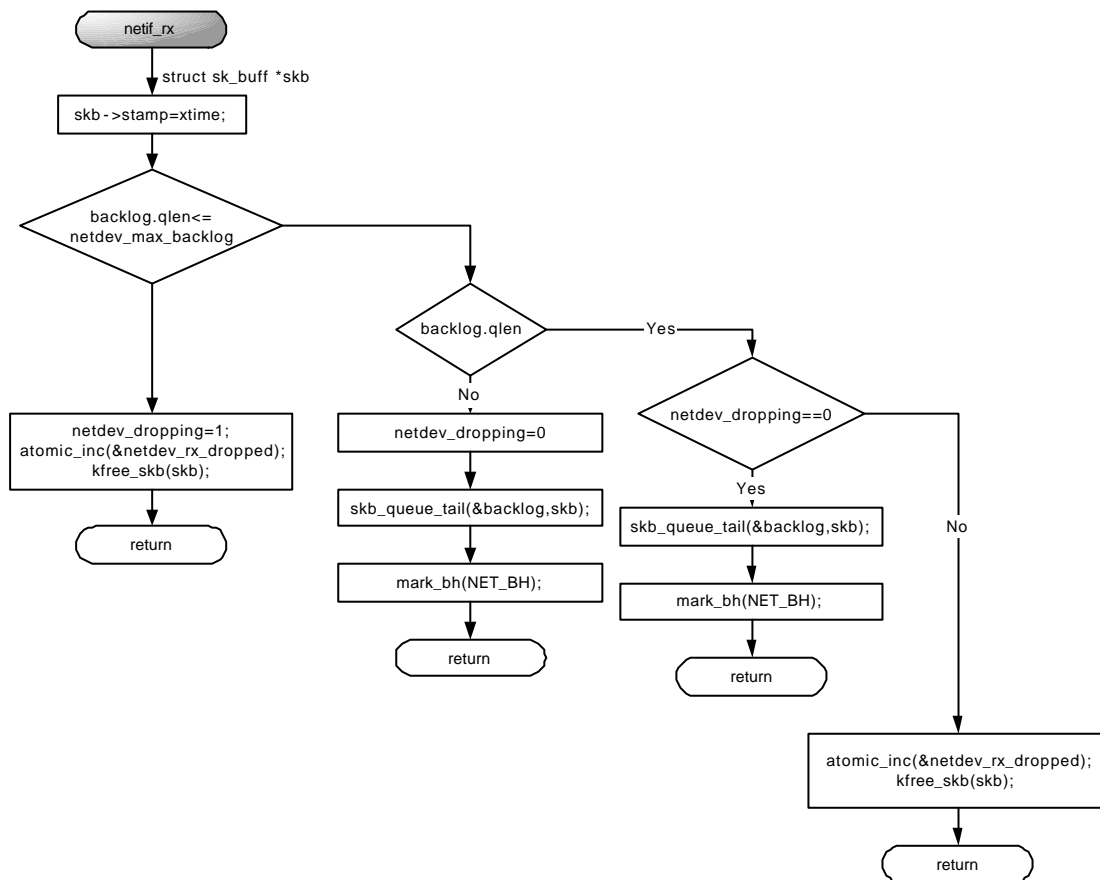
4.3 netif_rx

가

netif_rx

netif_rx()

net/core/dev.c



4.1

```

757 void netif_rx(struct sk_buff *skb)
758 {
763     skb->stamp = xtime;
770     if (backlog.qlen <= netdev_max_backlog) {
771         if (backlog.qlen) {

```

```

772             if (netdev_dropping == 0) {
773                 skb_queue_tail(&backlog,skb);
774                 mark_bh(NET_BH);
775                 return;
776             }
777             atomic_inc(&netdev_rx_dropped);
778             kfree_skb(skb);
779             return;
780         }

backlog             tail   skb   enqueue   mark_bh(NET_BH)
. 777             779             netdev_dropping   0
skb               drop   . netdev_dropping             791
    backlog가      1             . backlog             dev.c             156             struct
sk_buff_head      . backlog             .

785             netdev_dropping = 0;
787             skb_queue_tail(&backlog,skb);
788             mark_bh(NET_BH);
789             return;
790     }

backlog.qlen      , 785             789             ,   가   backlog
    skb           mark_bh(NET_BH)             ,   .

791     netdev_dropping = 1;
792     atomic_inc(&netdev_rx_dropped);
793     kfree_skb(skb);
794 }

backlog   가   netdev_dropping   1   ,   skb   .

```

4.4 mark_bh()

```

mark_bh()           Bottom Half Handler           .
include/asm-i386/softirq.h           NET_BH
include/linux/interrupt.h   9           . mark_bh()
    .           bh_active           9           BH handler가

```

net_bh()

```
100 extern inline void mark_bh(int nr)
101 {
102     set_bit(nr, &bh_active);
103 }
```

bh_active kernel/softirq.c

```
23 atomic_t bh_mask_count[32];
24 unsigned long bh_active = 0;
25 unsigned long bh_mask = 0;
26 void (*bh_base[32])(void);
```

```
26         active
kernel/softirq.c         __initfunc()         mark_bh(NET_BH)
, net_bh()가         _initfunc()
```

```
1893 __initfunc(int net_dev_init(void))
1894 {
1895     struct device *dev, **dp;
1905     skb_queue_head_init(&backlog);
2027     init_bh(NET_BH, net_bh);
2031     dev_mcast_init();
2044     return 0;
2045 }
```

2027 NET_BH net_bh init_bh()
. init_bh() kernel/softirq.c .

```
76 extern inline void init_bh(int nr, void (*routine)(void))
77 {
78     unsigned long flags;
79
80     bh_base[nr] = routine;
81     atomic_set(&bh_mask_count[nr], 0);
82 }
```

```

83     spin_lock_irqsave(&i386_bh_lock, flags);
84     bh_mask |= 1 << nr;
85     spin_unlock_irqrestore(&i386_bh_lock, flags);
86 }
80     NET_BH                                bh_base    net_bh    .

```

4.5 net_bh()

```

835 void net_bh(void)
836 {
837     struct packet_type *ptype;
838     struct packet_type *pt_prev;
839     unsigned short type;
840     unsigned long start_time = jiffies;
849
850     NET_PROFILE_ENTER(net_bh);
857
850     NET_PROFILE_ENTER                                include/net/profile.h
profile
.

858     if (qdisc_head.forw != &qdisc_head)
859         qdisc_run_queues();
858     ,                                discipline                                qdisc_head    가
. (net/sched/sch_generic.c).                                qdisc_head.forw = qdisc_head
.                                discipline    가
discipline    , qdisc_run_queues()    .

.

875     while (!skb_queue_empty(&backlog))
876     {
877         struct sk_buff * skb;
878
880         if (jiffies - start_time > 1)
881             goto net_bh_break;
882
886         skb = skb_dequeue(&backlog);
887
backlog q                                가                                886                                backlog
sk
.

916         skb->h.raw = skb->nh.raw = skb->data;

```



```

918         if (skb->mac.raw < skb->head || skb->mac.raw > skb->data) {
919             printk(KERN_CRIT "%s: wrong mac.raw ptr, proto=%04x\n", skb->dev
v->name, skb->protocol);
920             kfree_skb(skb);
921             continue;
922         }
928         type = skb->protocol;

```

Transport Header Network Header

```

. 918-922 . 928 type .
skb->protocol .

```

```

945         pt_prev = NULL;
946         for (ptype = ptype_all; ptype!=NULL; ptype=ptype->next)
947         {
948             if (!ptype->dev || ptype->dev == skb->dev) {
949                 if(pt_prev)
950                 {
951                     struct sk_buff *skb2=skb_clone(skb, GFP_ATOMIC
C);
952                     if(skb2)
953                         pt_prev->func(skb2,skb->dev, pt_prev);
954                 }
955                 pt_prev=ptype;
956             }
957         }

```

```

ptype_all    ETH_P_ALL    type    가
. struct packet_type    include/net/netdevice.h .
946-957    pt_prev    type    type    . 948    dev
가 NULL    wildcard    , .
             ptype    2 ,
clone .

```

```

959         for (ptype = ptype_base[ntohs(type)&15]; ptype != NULL; ptype = ptype->next)
960         {
961             if (ptype->type == type && (!ptype->dev || ptype->dev==skb->dev))
962             {
963                 if(pt_prev)
964                 {
965                     struct sk_buff *skb2;
966                     skb2=skb_clone(skb, GFP_ATOMIC);
967                     if(skb2)
968                         pt_prev->func(skb2, skb->dev, pt_pre

```

```

v);
980                                     }
982                                     pt_prev=ptype;
983                                     }
984                                     } /* End of protocol list loop */

                                     ptype_base[]                                     type                                     type
                                     clone                                     .                                     type   ETH_P_ALL   가
                                     .

990                                     if(pt_prev)
991                                     pt_prev->func(skb, skb->dev, pt_prev);
996                                     else {
997                                     kfree_skb(skb);
998                                     }
999   }

                                     type                                     .

1009   if (qdisc_head.forw != &qdisc_head)
1010                                     qdisc_run_queues();

                                     .

1025   NET_PROFILE_LEAVE(net_bh);
1026   return;
1027
1032 }

```

PROFILE

5. Linux Traffic Control

1. Linux

Linux (2.2.X) Traffic Control Code

4

(Queue Discipline)

(Class)

(Filter)

(Policing)

가

가

. 가

FIFO

가

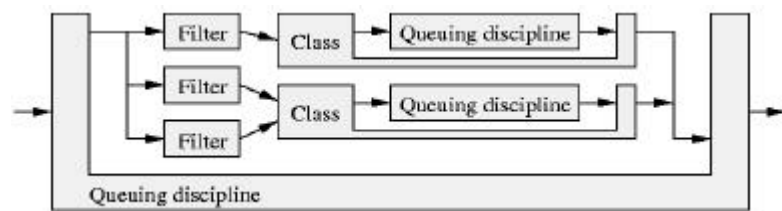
가

가

가

가

5.1



5.1

enqueue 가 ,

가

3. Linux (Queueing Disciplines)

. 5.1

5.1

	tc	Kernel	prefix
			tc
Queueing discipline	qdisc	sch_	q_
Class	class	(sch_)	(q_)
Filter	filter	cls_	f_

```
include/net/pkt_sched.h struct Qdisc_ops (
    enqueue . 가 가
    , enqueue ,
enqueue .
    dequeue .
    가 NULL .
    requeue . dequeue ,
    .
    drop .
    init .
    reset . ,
    가 reset
    .
    destroy . ,
```

```

        .
        dump
        .

        struct Qdisc
        .
        . (net/core/dev.c dev_queue_xmit
        ), (include/linux/netdevice.c struct device
        qdisc ) enqueue 가 dev_queue_xmit
        include/net/pkt_sched.h qdisc_wakeup ,
        . qdisc_wakeup net/sched/sch_generic.c
        qdisc_restart . qdisc_restart ,
        .
        . ,
        hard_start_xmit . ,
        requeue . (
        ) qdisc_wakeup . TBF(Token
        Bucket Filter)가 , net/core/dev.c net_bh run_queues
        qdisc_restart . net_bh
        'bottom-half" , 5.2
        .

```


4. (Class)

```

class ID
internal ID . internal ID
. internal ID 0 “
” “hot found” . class ID u32
internal ID unsigned long 가 , internal ID
. get change class ID
. Class ID .
Class ID ID , major number 가
ID가 , minor number가
.
include/net/pkt_sched.h struct Qdisc_class_ops .
FIFO
.
graft
.
get class ID , internal ID .
가 usage count , 가 .
put get 가 , 가
. 가 usage count .
count가 0 put .
change ,
. change
init , include/

```



```

linux/pkt_sched.h
    delete
    walk
    tcf_chain
    bind_tcf
    unbind_tcf
    put
    dump_class
    enqueue
    tc_classify
    struct tcf_result
    ID(class)
    (policing decision)
    TC_POLICE_OK
    skb->priority
    sk->priority
    sk_buff

```

가

callback

가

가

sch_cbq

unbind_tcf

dump

include/net/pkt_cls.h

tc_classify

include/net/pkt_cls.h

tc_result class ID(classid), internal

tc_classify -1

include/linux/pkt_cls.h

TC_POLICE_OK 0

ID

skb->priority

include/linux/skbuff.h

struct

sk->priority

```

include/net/sock.h struct sock . sk->priority
SO_PRIORITY (net/core/sock.c) .
,
. skb->priority IPv4 TOS ,
Diffserv .
enqueue 가 , enqueue 가
. 가
. FIFO ,
WF2Q+ , .

```

5. (Filter)

```

가
. 가 enqueue .
.
. 가 .
include/linux/if_ether.h , skb->protocol
.
가 .
가 . 32
가 . classID , major number
minor number . 0
. 가 internal ID 가 , get
. 가 가 .
. include/net/pkt_cls.h struct
tcf_proto_ops .
classify ( ) ,
TC_POLICE_.. . 가

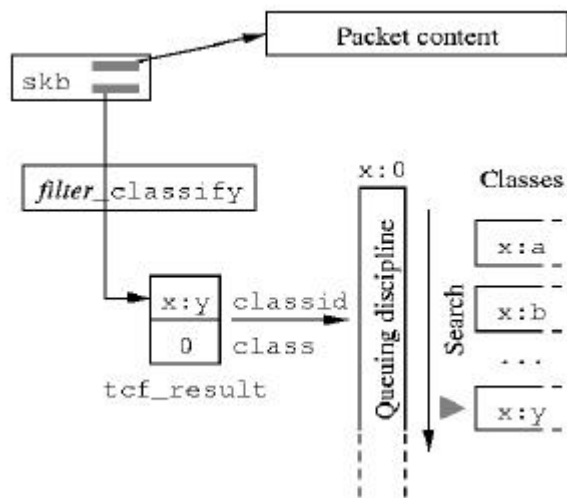
```

```

TC_POLICE_UNSPEC, struct tcf_result res, class
ID 가 internal class ID가 internal class ID
res->class 0
init
destory, sch_cbq
destroy 가
get ID
put get 가
change,
가
delete
destroy
walk callback
dump

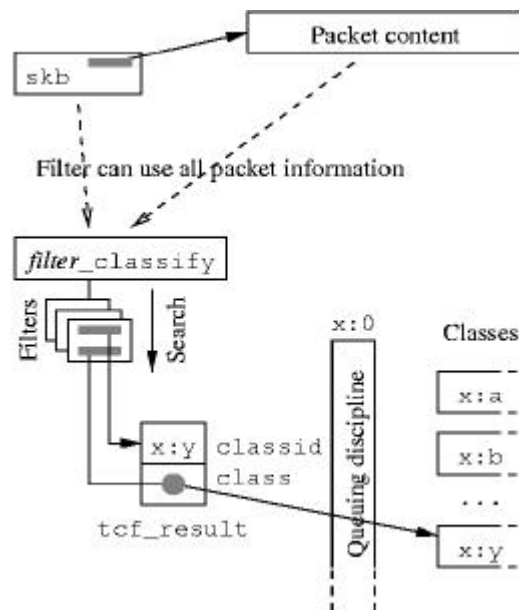
RSVP, cls_rsvp.h
, cls_rsvp.c cls_rsvp6.c include,
가
가, cls_fw cls_route,
class ID
, cls_fw firewall marking
(generic) ( 5.3)

```



5.3 Generic Filter

(cls_rsvp cls_u32) 가
(specific) .
internal filter ID internal ID .
internal filter ID "filter ID "
,
(specific filter)가
, internal ID ,
5.4
가 internal ID .



5.4 Specific Filter