# Sheet 4

# sk_buff structure

1    **struct sk_buff**

2

3    /*
4     *    Definitions for the 'struct sk_buff' memory handlers.
5     *
6     *    Authors:
7     *        Alan Cox, <gw4pts@gw4pts.ampr.org>
8     *        Florian La Roche, <rzsfl@rz.uni-sb.de>
9     *
10    *    This program is free software; you can redistribute it and/or
11    *    modify it under the terms of the GNU General Public License
12    *    as published by the Free Software Foundation; either version
13    *    2 of the License, or (at your option) any later version.
14    */

15

16   #define SKB_DATA_ALIGN(X)   (((X) + (SMP_CACHE_BYTES-1)) & ~(SMP_CACHE_BYTES-1))
17   #define SKB_MAX_ORDER(X,ORDER)(((PAGE_SIZE<<(ORDER)) - (X) - sizeof(struct
18   skb_shared_info))&~(SMP_CACHE_BYTES-1))
19   #define SKB_MAX_HEAD(X)         (SKB_MAX_ORDER((X),0))
20   #define SKB_MAX_ALLOC       (SKB_MAX_ORDER(0,2))

21

22   /* A. Checksumming of received packets by device.
23    *
24    *    NONE: device failed to checksum this packet.
25    *        skb->csum is undefined.

```
26   *
27   *   UNNECESSARY: device parsed packet and wouldbe verified checksum.
28   *        skb->csum is undefined.
29   *        It is bad option, but, unfortunately, many of vendors do this.
30   *        Apparently with secret goal to sell you new device, when you
31   *        will add new protocol to your host. F.e. IPv6. 8)
32   *
33   *   HW: the most generic way. Device supplied checksum of _all_
34   *        the packet as seen by netif_rx in skb->csum.
35   *        NOTE: Even if device supports only some protocols, but
36   *        is able to produce some skb->csum, it MUST use HW,
37   *        not UNNECESSARY.
38   *
39   * B. Checksumming on output.
40   *
41   *   NONE: skb is checksummed by protocol or csum is not required.
42   *
43   *   HW: device is required to csum packet as seen by hard_start_xmit
44   *   from skb->h.raw to the end and to record the checksum
45   *   at skb->h.raw+skb->csum.
46   *
47   *   Device must show its capabilities in dev->features, set
48   *   at device setup time.
49   *   NETIF_F_HW_CSUM  - it is clever device, it is able to checksum
50   *              everything.
51   *   NETIF_F_NO_CSUM - loopback or reliable single hop media.
```

```
52    *   NETIF_F_IP_CSUM - device is dumb. It is able to csum only
53    *              TCP/UDP over IPv4. Sigh. Vendors like this
54    *              way by an unknown reason. Though, see comment above
55    *              about CHECKSUM_UNNECESSARY. 8)
56    *
57    *   Any questions? No questions, good.            --ANK
58    */
59
60    #ifdef __i386__
61    #define NET_CALLER(arg) (*(((void**)&arg)-1))
62    #else
63    #define NET_CALLER(arg) __builtin_return_address(0)
64    #endif
65
66    struct sk_buff_head {
67        /* These two members must be first. */
68        struct sk_buff          *next;
69        struct sk_buff          *prev;
70        __u32                   qlen;
71        spinlock_t              lock;
72    };
73
74    struct sk_buff;
75
76    #define MAX_SKB_FRAGS 6
77
```

```
78   typedef struct skb_frag_struct skb_frag_t;
79
80   struct skb_frag_struct {
81       struct page        *page;
82       __u16              page_offset;
83       __u16              size;
84   };
85
86   /* This data is invariant across clones and lives at
87    * the end of the header data, ie. at skb->end.
88    */
89   struct skb_shared_info {
90       atomic_t           dataref;
91       unsigned int       nr_frags;
92       struct sk_buff     *frag_list;
93       skb_frag_t         frags[MAX_SKB_FRAGS];
94   };
95
96
```

This is a massively important structure. It is the way of representing packets within the kernel.  I have deleted some stuff for the purposes of clarity.

```
100  struct sk_buff {
```
Linking these buffers together. The reason this must be first is that we can cast the packet to sk_buff_head, defined above.
```
102      /* These two members must be first. */
```

```
103          struct sk_buff          *next;          /* Next buffer in list              */
104          struct sk_buff          *prev;          /* Previous buffer in list          */
105          struct sk_buff_head     *list;          /* List we are on                   */
106
```

Back pointer to the sock structure we belong to

```
108          struct sock             *sk;            /* Socket we are owned by    */
109
```

The stamp is the time that the last protocol touched this buffer. Actually, this is a bit more involved than I'm making out –
useful for scheduling.

```
112          struct timeval          stamp;          /* Time we arrived                  */
113
```

In the administration of network buffers the identity of the device used for sending or receiving the packet must be known.

```
115          struct net_device       *dev;           /* Device we arrived on/are leaving by */
116
```

Just what you'd expect from a transport layer header, but note the overlay. You'll find the definitions in
include/linux/tcp.h::tcphdr, include/linux/udp.h::udphdr, include/linux/icmp.h::icmphdr, etc. So, for example, a udp header
is given by:

```
        struct udphdr {
                __u16 source;
                __u16 dest;
                __u16 len;
                __u16 check;
        };
```

```
126        /* Transport layer header */
127        union
128        {
129            struct tcphdr      *th;
130            struct udphdr      *uh;
131            struct icmphdr     *icmph;
132            struct igmphdr     *igmph;
133            struct iphdr       *ipiph;
134            struct spxhdr      *spxh;
135            unsigned char      *raw;
136        } h;
137
```

Again, no surprises here. E.g. from include/linux/ip.h::iphdr we see:

```
struct iphdr {
#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;
```

```
150              __u16 tot_len;
151              __u16 id;
152              __u16 frag_off;
153              __u8  ttl;
154              __u8  protocol;
155              __u16 check;
156              __u32 saddr;
157              __u32 daddr;
158              /*The options start here. */
159         };
160         /* Network layer header */
161         union
162         {
163              struct iphdr      *iph;
164              struct ipv6hdr    *ipv6h;
165              struct arphdr     *arph;
166              struct ipxhdr     *ipxh;
167              unsigned char     *raw;
168         } nh;
169
```

170    Still nothing unusual. So e.g. [include/linux/if_ether.h::ethhdr](include/linux/if_ether.h::ethhdr)

```
171         struct ethhdr
172         {
173              unsigned char    h_dest[ETH_ALEN];      /* destination eth addr    */
```

```
174            unsigned char     h_source[ETH_ALEN];    /* source ether addr      */
175            unsigned short    h_proto;               /* packet type ID field   */
176        };
177        /* Link layer header */
178        union
179        {
180            struct ethhdr     *ethernet;
181            unsigned char     *raw;
182        } mac;
183
```

184 This related to destination cache information.

```
185        struct  dst_entry *dst;
186
```

187 Private data for each layer. E.g. the ip layer keeps include/net/ip.h::inet_skb_parm (basically IP options) in there, whereas
188 TCP keeps include/net/tcp.h::tcp_skb_cb (sequence numbers, flags, etc.) in there.

```
189        /*
190         * This is the control buffer. It is free to use for every
191         * layer. Please put your private variables there. If you
192         * want to keep them across layers you have to do a skb_clone()
193         * first. This is owned by whoever has the skb queued at the moment.
194         */
195        char              cb[48];
196
```

197 Comment notwithstanding, len holds the length of the packet (including headers), and data_len the length of the data part.
198 csum holds the checksum if it has been calculated. See comment at head of file re checksumming.

```
199        unsigned int      len;              // Length of actual data
200        unsigned int      data_len;
201        unsigned int      csum;             // Checksum
```
202 **This is the length of this buffer, including the length of this struct, used for memory management purposes.**
```
203        unsigned int truesize;             // Buffer size
204
```
205 **Management parameters.**
```
206        unsigned char     cloned,           // head may be cloned (check refcnt to be sure).
207                          pkt_type,         // Packet class
208                          ip_summed;        // Driver fed us an IP checksum
209        __u32             priority;         // Packet queueing priority
210        unsigned short    protocol;         // Packet protocol from driver.
211        unsigned short    security;         // Security level of packet
```
212 **Actually, see include/linux/skbuff.h::skb_get – this is a reference count to this sk_buff**
```
213        atomic_t          users;            // User count - see datagram.c,tcp.c
214
```
215 **This is a really important bit – it's where the data resides. The head pointer points to the first part of the buffer (i.e. the bit**
216 **containing the header), the data pointer points to the part of the buffer containing the data and the tail pointer to whatever**
217 **follows the data. End, naturally points to the end. There are a lot of helper functions both in this file and in net/core/skbuff.c**
218 **to allow manipulation of these pointers, the addition of extra space and so forth. See below.**
```
219        unsigned char     *head;           /* Head of buffer            */
220        unsigned char     *data;           /* Data head pointer         */
221        unsigned char     *tail;           /* Tail pointer              */
222        unsigned char     *end;            /* End pointer               */
```

```
223
224       void (*destructor)(struct sk_buff *);  /* Destruct function      */
225
226   };
227
```

227 **Sending UDP packets – the code**
228
229 OK, so lets take a quick look at what happens to the sk_buff when we send a UDP packet (net/ipv4/udp.c). This what gets
230 passed to net/ipv4/udp.c::udp_sendmsg:

231 ```
int udp_sendmsg(struct sock *sk, struct msghdr *msg, int len)
```
232
233 The msghdr here is defined in include/linux/socket.h as:
234 ```
struct msghdr {
235     void            *msg_name;          /* Socket name              */
236     int              msg_namelen;       /* Length of name           */
237     struct iovec    *msg_iov;       /* Data blocks              */
238     __kernel_size_t  msg_iovlen;        /* Number of blocks         */
239     void            *msg_control; /* Per protocol magic     */
240                                             /* (eg BSD file descriptor passing) */
241     __kernel_size_t  msg_controllen;  /* Length of cmsg list    */
242     unsigned         msg_flags;
243 };
```
244
245 The data blocks are in an array of iovecs (defined in include/linux/uio.h), each of which is a structure with two fields
246 of interest:
247 ```
struct iovec {
248     void            *iov_base;          /* BSD uses caddr_t, 1003.1g void *) */
249     __kernel_size_t  iov_len;       /* Must be size_t (1003.1g)    */
250 };
```

251

252  So, we're being passed an array of pointers to odd bits of data of interest rather than a contiguous area of memory. This is
253  pretty standard within unix. Now, let's go back to the code of net/ipv4/udp.c::udp_sendmsg. The next thing of interest is the
254  declaration:

255  `struct udpfakehdr ufh;`

256

257       For this, we need to look earlier in the file – we see that the first part of this is reserved for a real udp header followed
258  by some other info.

```
259      struct udpfakehdr {
260          struct udphdr uh;
261          u32             saddr;
262          u32             daddr;
263          struct iovec *iov;
264          u32             wcheck;
265      };
```

266

267  The fields of this header are filled in (with the exception of the checksum, which is set to zero) and iov is made to point to
268  the iov we were passed. We then call net/ipv4/ip_output.c::ip_build_xmit thus:

```
269  err = ip_build_xmit( sk,
270                  (sk->no_check == UDP_CSUM_NOXMIT ? udp_getfrag_nosum : udp_getfrag),
271                  &ufh, ulen, &ipc, rt, msg->msg_flags);
```

272

273  The definition of this routine is below and we care about the first four fields in this context.

274  `int ip_build_xmit( struct sock      *sk,`

```
275                     int       getfrag(const void *, char *, unsigned int, unsigned int),
276                     const void          *frag,
277                     unsigned            length,
278                     struct ipcm_cookie *ipc,
279                     struct rtable      *rt,
280                     int                 flags)
281
```

282  **Within this, we have the definition:**

```
283  struct sk_buff *skb;
```

284

285  **This next call itself calls** net/core/skbuff.c::alloc_skb. **This allocates a sk_buff from a central store, and initialises it with**
286  **head=tail=data all pointing to the same allocated block of memory. hh_len essentially respresents the MAC header length,**
287  **rounded up to the next multiple of 16bytes.**

```
288  int hh_len = (rt->u.dst.dev->hard_header_len + 15)&~15;
289  skb = sock_alloc_send_skb(sk, length+hh_len+15, flags&MSG_DONTWAIT, &err);
```

290

291  **The first call moves the data and tail pointers forward by hh_len, to give us some header room and the second moves the**
292  **tail pointer forward to give us more data room. When we're done, iph points to the data part of the structure.**

```
293  skb_reserve(skb, hh_len);
294  iph = (struct iphdr *)skb_put(skb, length);
```

295

296  **We use our callback to get the data out of the fake header we were passed and to stick it in the data part, possibly after an**
297  **ip header. This is either** net/ipv4/udp.c::udp_getfrag **or** net/ipv4/udp.c::udp_getfrag_nosum, **depending on whether we**
298  **need to do checksumming or not.**

```
299  if(!sk->protinfo.af_inet.hdrincl) {
300      <fill in IP header details>
301      err = getfrag(frag, ((char *)iph)+iph->ihl*4, 0, length-iph->ihl*4);
302  }
303  else
304      err = getfrag(frag, (void *)iph, 0, length);
305
```

306   **Ok, so let's take net/ipv4/udp.c::udp_getfrag_nosum. The code we care about is:**

```
307  static intudp_getfrag_nosum(const void *p, char *to, unsigned int offset,
308                                    unsigned int fraglen)
309  {
310      struct udpfakehdr *ufh = (struct udpfakehdr *)p;
311
312      Copy the header part of the fake header
313      memcpy(to, ufh, sizeof(struct udphdr));
314      Now copy the data from the iovec into our buffer. See net/core/iovec.c::memcpy_fromiovecend
315      return memcpy_fromiovecend(to+sizeof(struct udphdr), ufh->iov, offset,
316                                      fraglen-sizeof(struct udphdr));
317  }
318
```

319   **Jumping back to net/ipv4/ip_output.c::ip_build_xmit, we see the following. This does network filtering, then jumps to the**
320   **routine named as the last parameter.**

```
321  err = NF_HOOK(PF_INET, NF_IP_LOCAL_OUT, skb, NULL, rt->u.dst.dev,
322  output_maybe_reroute);
```

323

**Then we come here to send the packet.**

```
output_maybe_reroute(struct sk_buff *skb)
{
    return skb->dst->output(skb);
}
```

329