# Network buffers
## The BSD, Unix SVR4 and Linux approaches
## (BSD4.4,SVR4.2,Linux2.6.2)

*Roberto Innocente, Olumide S.Adewale*

SISSA - Trieste

## 1. Introduction

On many operating systems, network buffers are a special category of buffers to be used for network operations. This is because they are used in a peculiar way and there is a need to make the most common of those operations extremely efficient. Most of the network operations result in :

- addition or stripping of headers or trailers

- transfer of data to/from devices

- joining or splitting buffers

According to Gilder's law, an empirical law recently being cited very often, the total communication bandwidth triples every twelve months. In contrast, according to Moore's law, cpu performance only doubles every 18 months. This will create an increasing gap between network and cpu performance.

We analyze how the traditional in-kernel networking is performed on 3 major network stacks : Bsd Unix, Unix SystemV and Linux, looking for possible improvements.

This area has recently seen the emerging of new software paradigms (User Level Networking), to keep up with the performance improvement of communications.
User-level networking gives a boost of performance pinning down some user memory where network cards perform directly their I/O operations and memory connected network cards, provide the protection mechanism by the standard page protection mechanism that most processors have.

## 2. I/O architecture

Some computers have a uniform view of memory and I/O devices, everything is mapped in a single address space (Vax, Alpha). On these architectures CSRs (Control and Status Registers) and memory on I/O devices are mapped in memory, and read and write memory cycles on the bus may result in reads and writes to/from the I/O devices depending on the address. Other architectures feel I/O devices are peculiar and they have a separate address space for them (x86). On these architectures a signal on the bus denotes an I/O operation instead of a standard memory operation, and special input and output instructions are used (I/O ports and I/O memory 64 KB).

## 3. History of network code

The roots of the BSD network code are back to the 4.2 BSD TCP/IP implementation in 1983. The 4.4 BSD Lite implementation, which we discuss appeared in 1994. The first STREAMS implementation appeared on SVR3 in 1986, then there was a complete new implementation which we discuss, on SVR4 in 1988 and later on SVR4.2MP in late 1993. The Linux TCP/IP network code started with the Ross Biro NET-1 implementation in release 0.98 in 1992 with device drivers written by Donald Becker. After R.Biro quit, then Fred van Kempen worked on the NET-2 project rewriting major parts of the initial release. NET-2

appeared publicly after Alan Cox debugging of the code  as NET-2D in kernel 0.99.10 in 1993.  A. Kut-
nezsov and A. Kleen and D.Miller.  NET-3 was incorporated in 1.2.x and the current release of the network
code in 2.2 kernels and on is NET-4.

## 4.  The Berkeley approach

The organization of the network buffers in BSD derives from the observation that packets on networks are
for the most part either very small or near the maximum size.  This bimodal distribution of packets on net-
works leads to the choice of a very small fixed size combined descriptor and buffer called an mbuf (mem-
ory buffer), with the possibility to keep there a pointer to an external page of data in the case of a large
packet.  Therefore the provision of a small structure of 128 bytes, that can keep the header and the data if
the packet is small.

_____ *sys/sys/mbuf.h*

```
67     struct m_hdr {
68          struct    mbuf *mh_next;      /* next buffer in chain */
69          struct    mbuf *mh_nextpkt;   /* next chain in queue/record */
70          caddr_t   mh_data;       /* location of data */
71          int  mh_len;                 /* amount of data in this mbuf */
72          short mh_type;        /* type of data in this mbuf */
73          short mh_flags;       /* flags; see below */
74     };
```
_____ *sys/sys/mbuf.h*

The two possibilities are never used at the same time so that if the data is stored in the internal area, the
mbuf can't use an external page and viceversa. m_len is the size of the data in this mbuf, while m_data
is a pointer to the beginning of the data. m_next is a pointer that can link multiple mbufs in a chain
(singly linked list). m_pktnext is a pointer that links multiple packets (mbuf chains) on the device or
socket queue.  There are 108 bytes available for data in an mbuf, but the code tries to be smart, keeping
some space in front and after the data for the addition of headers.
Operations on networks buffers frequently result to addition or stripping of data in the beginning or the end
of a packet.  Mbufs in BSD can be joined in a single linked list through their m_next pointer. This is usu-
ally called an mbuff chain, and operations on packets are mapped on BSD to insertion or deletion on the
head of a chain.  In the first mbuf some bytes have to be used for storing the total lenght of a packet and
eventually the receiving interface, these are kept in the pkthdr structure.

_____ *sys/sys/mbuf.h*

```
77     struct     pkthdr {
78          struct     ifnet *rcvif;       /* rcv interface */
79          int  len;            /* total packet length */
80     };
```
_____ *sys/sys/mbuf.h*

The total length of a packet split in multiple mbufs is stored in the m_pkthdr.len field of the first mbuf
of the chain, while in this case m_len is only the amount of data in each mbuf.  In the first mbuf of a
chain, 8 bytes of these 108 available are used to specify the length and interface for the packet and so only
100 bytes of data are available in it.  If the data is from 0 to 84 bytes, then a single mbuf is allocated and the
data is placed in it leaving 16 bytes free in front of it for an eventual header.  If the data is from 85 to 100
bytes then the data is placed at the beginning of the mbuf.  If the data is from 101 to 207 bytes[*] a chain of 2
mbufs is allocated and the data is placed starting at the beginning of the first mbuf.  If the data is  more than
207 then an mbuf with an external data area is allocated. In this case the mbuf has a pointer to the external

_____
[*] in fact there would be space for up to 208 bytes, but because of a little mistake in the code that uses <
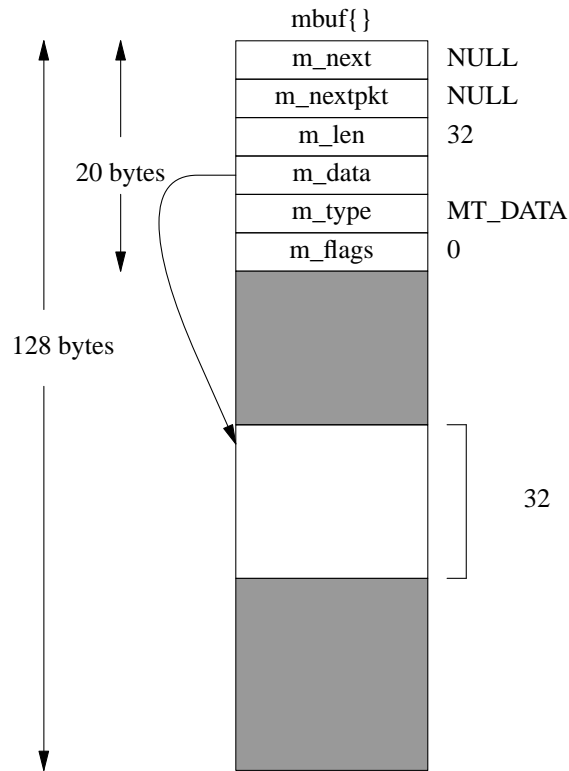instead of <=, this is what happens

*Figure 1 An mbuf with some data*

data area (also called a page cluster : because it is the same page used for virtual memory) and the data is stored from the beginning of the pagecluster.

─────────────────────────────────────────────────────────────────────────────── *sys/sys/mbuf.h*

```
83    struct m_ext {
84        caddr_t   ext_buf;        /* start of buffer */
85        void (*ext_free)();       /* free routine if not the usual */
86        u_int ext_size;       /* size of buffer, for ext_free */
87    };
```

─────────────────────────────────────────────────────────────────────────────── *sys/sys/mbuf.h*

| simple mbuf | packet header mbuf | pagecluster mbuf | packet header pagecluster | |
|---|---|---|---|---|
| | | | | m_next |
| | | | | m_nextpkt |
| 0 - 108 | 0 - 100 | 208-2048 | 208-2048 | m_len |
| | | | | m_data |
| MT_xxx | MT_xxx | MT_xxx | MT_xxx | m_type |
| 0 | M_PKTHDR | M_EXT | M_PKTHDR \| M_EXT | m_flag |
| | | | | m_pkthdr.len |
| | | | | m_pkthdr.rcvif |
| | | | | m_ext.ext_buf |
| | | | | m_ext.ext_free |
| | | 2048 | 2048 | m_ext.ext_size |

m_hdr{} (20 bytes)

pkthdr{} (8 bytes)

m_ext{} (12 bytes)

108-byte (=MLEN) buffer: m_dat

100-byte (=MHLEN) buffer: m_pktdat

2048-byte cluster (external buffer)
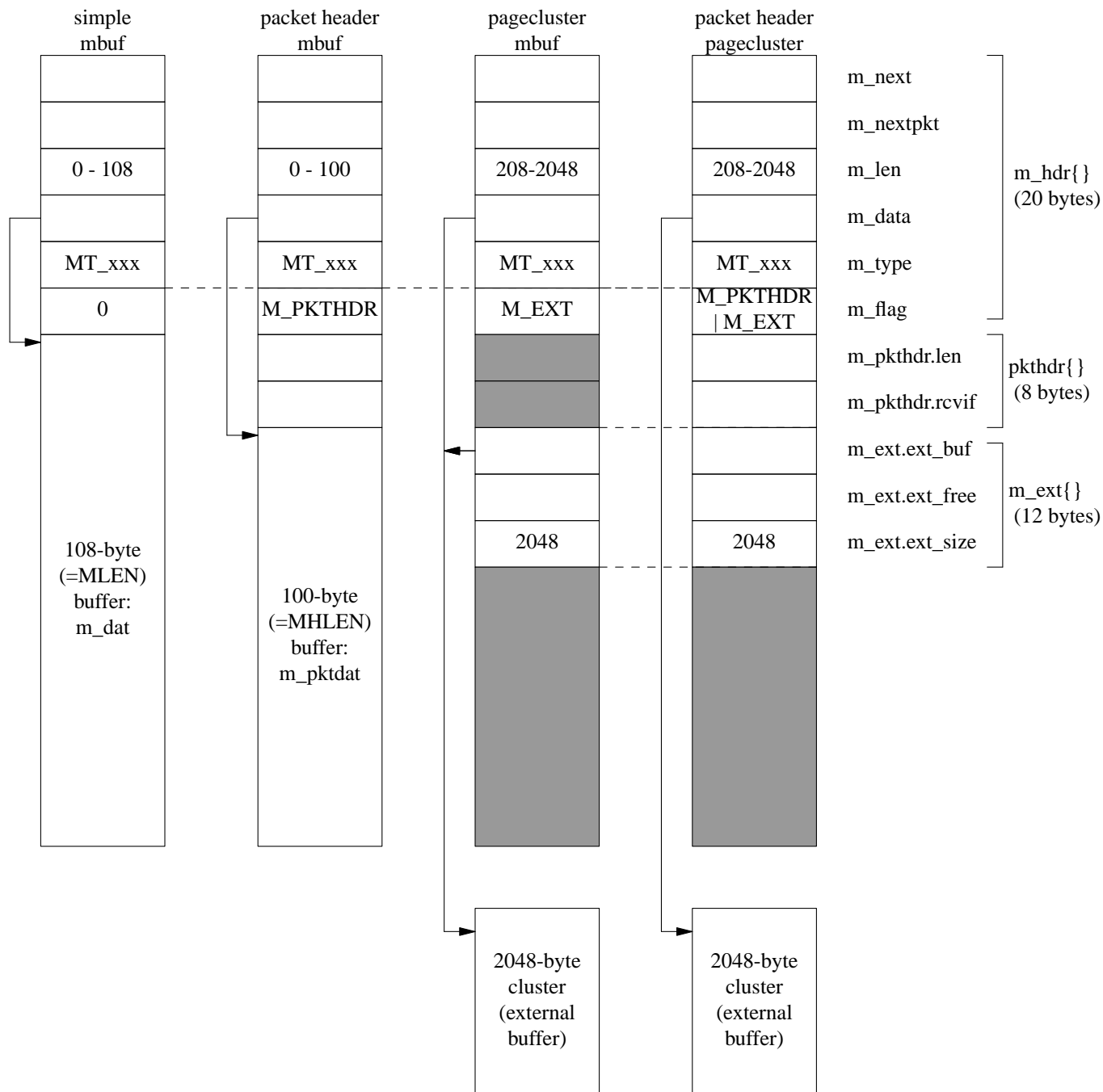
2048-byte cluster (external buffer)

*Figure 2 Various mbuf types*

Instead of using a general structure with all the fields necessary, BSD chooses to save space and uses a sovrapposition of them specified by a C union. The different layout used is chosen by the `M_EXT`, `M_PKTHDR` flags :

```
89      struct mbuf {
90          struct    m_hdr m_hdr;
```

```
 91          union {
 92                  struct {
 93                          struct    pkthdr MH_pkthdr;    /* M_PKTHDR set */
 94                          union {
 95                                  struct    m_ext MH_ext;   /* M_EXT set */
 96                                  char MH_databuf[MHLEN];
 97                          } MH_dat;
 98                  } MH;
 99                  char M_databuf[MLEN];              /* !M_PKTHDR, !M_EXT */
100          } M_dat;
101     };
102     #define   m_next       m_hdr.mh_next
103     #define   m_len        m_hdr.mh_len
104     #define   m_data       m_hdr.mh_data
105     #define   m_type       m_hdr.mh_type
106     #define   m_flags      m_hdr.mh_flags
107     #define   m_nextpkt m_hdr.mh_nextpkt
108     #define   m_act        m_nextpkt
109     #define   m_pkthdr  M_dat.MH.MH_pkthdr
110     #define   m_ext      M_dat.MH.MH_dat.MH_ext
111     #define   m_pktdat  M_dat.MH.MH_dat.MH_databuf
112     #define   m_dat      M_dat.M_databuf
```

—————————————————————————————————————————————————————————————— *sys/sys/mbuf.h*

## 4.1. Memory issues

Mbuf clusters are allocated with the standard kernel memory allocator malloc and mapped in a dedicated area of the kernel virtual space. An initial allocation is done at system startup and it can be expanded up to a configurable maximum (sysctl). Memory allocated to mbuf clusters is never given back. An array of reference counts mclrefcnt is statically allocated at system startup, its size is NMBCLUSTERS that was 512 for gateways and 256 for other systems. But on the x86 there is no need for page clustering like on the Vax and so CLSIZE=1 and NBPG=4096, and so CLBYTES=CLSIZE*NBPG=4096 and MCLBYTES=1024. Then the size of the kernel window into which to map mbuf clusters is 2 MB in pages NKMEMCLUSTERS=(2048*1024/CLBYTES)=512. The initial memory mapped for mbuf clusters is VM_MBUF_SIZE=NMBCLUSTERS*MCLSIZE=512*1024 enough to map the max number of clusters(kmem_suballoc). Then every time this is needed a page is allocated with kmem_alloc and assigned to the mb_map (the window in the kernel virtual space dedicated to mbuf clusters) if there are still free slots.

## 4.2. mbuf macros and functions

There are many macros and functions defined to help with the use of mbufs. mbufs are allocated at the device level for an incoming packet with m_devget:

—————————————————————————————————————————————————————— *usr/src/sys/kern/uipc_mbuf.c*

```
591     struct mbuf *
592     m_devget(buf, totlen, off0, ifp, copy)
593          char *buf;
594          int totlen, off0;
595          struct ifnet *ifp;
596          void (*copy)();
```

_____ *usr/src/sys/kern/uipc_mbuf.c*

This function according to the size of the packet will allocate 1 or 2 normal mbufs or an mbuf with an external cluster, and will copy the data from the linear device buffer to the mbuf chain. In other places mbufs are allocated with the m_get function or MGET macro :

_____ *usr/src/sys/kern/uipc_mbuf.c*

```
157     struct mbuf *
158     m_get(nowait, type)
159          int nowait, type;
160     {
161          register struct mbuf *m;
162
163          MGET(m, nowait, type);
164          return (m);
165     }
```

_____ *usr/src/sys/kern/uipc_mbuf.c*

_____ *usr/src/sys/sys/mbuf.h*

```
170     #define   MGET(m, how, type) { \
171          MALLOC((m), struct mbuf *, MSIZE, mbtypes[type], (how)); \
172          if (m) { \
173               (m)->m_type = (type); \
174               MBUFLOCK(mbstat.m_mtypes[type]++;) \
175               (m)->m_next = (struct mbuf *)NULL; \
176               (m)->m_nextpkt = (struct mbuf *)NULL; \
177               (m)->m_data = (m)->m_dat; \
178               (m)->m_flags = 0; \
179          } else \
180               (m) = m_retry((how), (type)); \
181     }
```

_____ *usr/src/sys/sys/mbuf.h*

mbufs that have to be packet headers are allocated instead with m_gethdr and MGETHDR function and macro respectively, that set the M_PKTHDR flag in the m_flags field of the mbuf header and the m_data pointer is initialized to point after the pkthdr structure. External pages are allocated for mbufs using the MCLGET macro :

_____ *usr/src/sys/sys/mbuf.h*

```
224     #define   MCLGET(m, how) \
225          { MCLALLOC((m)->m_ext.ext_buf, (how)); \
226           if ((m)->m_ext.ext_buf != NULL) { \
227               (m)->m_data = (m)->m_ext.ext_buf; \
228               (m)->m_flags |= M_EXT; \
229               (m)->m_ext.ext_size = MCLBYTES;  \
230           } \
231          }
```

_____ *usr/src/sys/sys/mbuf.h*

This function allocates an external page and sets the m_data pointer in the existing mbuf to point to the beginning of the allocated page. mbufs are freed using the macro MFREE or the functions m_free or m_freem :

_____ *usr/src/sys/sys/mbuf.h*

```
261    #define   MFREE(m, nn) \
262          { MBUFLOCK(mbstat.m_mtypes[(m)->m_type]--;) \
263            if ((m)->m_flags & M_EXT) { \
264                MCLFREE((m)->m_ext.ext_buf); \
265            } \
266            (nn) = (m)->m_next; \
267            FREE((m), mbtypes[(m)->m_type]); \
268          }
```
———————————————————————————————————————————————————————— *usr/src/sys/sys/mbuf.h*

m_free frees just one mbuf, while m_freem frees all the mbufs in the chain. It frequently happens you want to put some data at the end of a buffer to leave all the possible room for headers. The M_ALIGN and MH_ALIGN macros conveniently set the m_data pointer of an mbuf to place an object of size len at the end of the mbuf :

———————————————————————————————————————————————————————— *usr/src/sys/sys/mbuf.h*

```
285    #define   M_ALIGN(m, len) \
286          { (m)->m_data += (MLEN - (len)) &~ (sizeof(long) - 1); }
287    /*
288     * As above, for mbufs allocated with m_gethdr/MGETHDR
289     * or initialized by M_COPY_PKTHDR.
290     */
291    #define   MH_ALIGN(m, len) \
292          { (m)->m_data += (MHLEN - (len)) &~ (sizeof(long) - 1); }
293
```
———————————————————————————————————————————————————————— *sys/sys/mbuf.h*

The M_PREPEND macro instead is used to prepend len bytes at the head of the current data in the mbuf. If there is not enough space, a new mbuf is linked in front of the mbuf chain :

———————————————————————————————————————————————————————— *usr/src/sys/sys/mbuf.h*

```
318    #define   M_PREPEND(m, plen, how) { \
319          if (M_LEADINGSPACE(m) >= (plen)) { \
320                (m)->m_data -= (plen); \
321                (m)->m_len += (plen); \
322          } else \
323                (m) = m_prepend((m), (plen), (how)); \
324          if ((m) && (m)->m_flags & M_PKTHDR) \
325                (m)->m_pkthdr.len += (plen); \
326    }
```
———————————————————————————————————————————————————————— *sys/sys/mbuf.h*

m_adj is a function that is used to trim len bytes from the head or the tail of the data :

———————————————————————————————————————————————————————— *sys/kern/uipc_mbuf.c*

```
381    void
382    m_adj(mp, req_len)
383          struct mbuf *mp;
384          int req_len;
```
———————————————————————————————————————————————————————— *sys/kern/uipc_mbuf.c*

m_cat is used instead to concatenate two mbufs together :

———————————————————————————————————————————————————————— *sys/kern/uipc_mbuf.c*

```
360     void
361     m_cat(m, n)
362          register struct mbuf *m, *n;
363     {
364          while (m->m_next)
365                m = m->m_next;
366          while (n) {
367                if (m->m_flags & M_EXT ||
368                     m->m_data + m->m_len + n->m_len >= &m->m_dat[MLEN]) {
369                      /* just join the two chains */
370                      m->m_next = n;
371                      return;
372                }
373                /* splat the data from one into the other */
374                bcopy(mtod(n, caddr_t), mtod(m, caddr_t) + m->m_len,
375                     (u_int)n->m_len);
376                m->m_len += n->m_len;
377                n = m_free(n);
378          }
379     }
```
———————————————————————————————————————————————— *sys/kern/uipc_mbuf.c*

There are different functions to copy mbufs m_copy, m_copydata, m_copyback, m_copym.
m_copym creates a new mbuf chain and copies len bytes starting at offset off0 from the old mbuf chain,
m_copy is the same as the m_copym except that it calls m_copym with a forth argument of nowait. If the
len argument is M_COPYALL, then all the remaining data in the mbuff after the offset is copied.

———————————————————————————————————————————————— *usr/src/sys/sys/mbuf.h*

```
337     /* compatiblity with 4.3 */
338     #define  m_copy(m, o, l)  m_copym((m), (o), (l), M_DONTWAIT)
```

———————————————————————————————————————————————— *[usr/src/sys/kern/uipc_mbuf.c]*

```
253     struct mbuf *
254     m_copym(m, off0, len, wait)
255          register struct mbuf *m;
256          int off0, wait;
257          register int len;
258     {
259          register struct mbuf *n, **np;
260          register int off = off0;
261          struct mbuf *top;
262          int copyhdr = 0;
263
264          if (off < 0 || len < 0)
265                panic("m_copym");
266          if (off == 0 && m->m_flags & M_PKTHDR)
267                copyhdr = 1;
268          while (off > 0) {
269                if (m == 0)
270                      panic("m_copym");
271                if (off < m->m_len)
272                      break;
273                off -= m->m_len;
274                m = m->m_next;
```

```
275               }
276           np = &top;
277           top = 0;
278           while (len > 0) {
279                 if (m == 0) {
280                       if (len != M_COPYALL)
281                             panic("m_copym");
282                       break;
283                 }
284                 MGET(n, wait, m->m_type);
285                 *np = n;
286                 if (n == 0)
287                       goto nospace;
288                 if (copyhdr) {
289                       M_COPY_PKTHDR(n, m);
290                       if (len == M_COPYALL)
291                             n->m_pkthdr.len -= off0;
292                       else
293                             n->m_pkthdr.len = len;
294                       copyhdr = 0;
295                 }
296                 n->m_len = min(len, m->m_len - off);
297                 if (m->m_flags & M_EXT) {
298                       n->m_data = m->m_data + off;
299                       mclrefcnt[mtocl(m->m_ext.ext_buf)]++;
300                       n->m_ext = m->m_ext;
301                       n->m_flags |= M_EXT;
302                 } else
303                       bcopy(mtod(m, caddr_t)+off, mtod(n, caddr_t),
304                             (unsigned)n->m_len);
305                 if (len != M_COPYALL)
306                       len -= n->m_len;
307                 off = 0;
308                 m = m->m_next;
309                 np = &n->m_next;
310           }
311           if (top == 0)
312                 MCFail++;
313           return (top);
314    nospace:
315           m_freem(top);
316           MCFail++;
317           return (0);
318    }
```
——————————————————————————————————————————————————— *[usr/src/sys/kern/uipc_mbuf.c]*

264-265  If the offset off is less than zero or the len requested is negative, then the kernel panics.

266-267 If a packet header mbuf is being copied since its beginning (from offset 0), then  a flag is set to force the copy of the packet header information (copyhdr).

268-275 The code runs through the mbuf chain skipping mbufs until offset bytes has being skipped. If this is not possible because the end of the chain is reached then it panics with the m_copym message.

278-310 This while loop tries to copy len bytes from the current position in the old mbuf chain to a chain of newly allocated mbufs.  If the argument len is M_COPYALL, then it copies up to the end of the old mbuf

chain. The pointer to the first mbuf allocated is kept in top. For the first mbuf allocated if the copyhdr is set, then the packet header info are copied (the len is decreased by the offset that is skipped). The data inside the mbuf is then copied, but if the mbuf had an external cluster, then the data is not copied, only the pointers are copied and the refcount for the pagecluster is incremented.

311- The new mbuf chain is then returned.

`m_copydata` copies data from an mbuf chain to a linear buffer starting at offset `off` bytes and m_copyback copies from a linear buffer to an mbuf chain starting `off` bytes from the beginning :

——————————————————————————————————— *[usr/src/sys/kern/uipc_mbuf.c]*

```
324     void
325     m_copydata(m, off, len, cp)
326           register struct mbuf *m;
327           register int off;
328           register int len;
329           caddr_t cp;
330     {
331           register unsigned count;
332
333           if (off < 0 || len < 0)
334                 panic("m_copydata");
335           while (off > 0) {
336                 if (m == 0)
337                       panic("m_copydata");
338                 if (off < m->m_len)
339                       break;
340                 off -= m->m_len;
341                 m = m->m_next;
342           }
343           while (len > 0) {
344                 if (m == 0)
345                       panic("m_copydata");
346                 count = min(m->m_len - off, len);
347                 bcopy(mtod(m, caddr_t) + off, cp, count);
348                 len -= count;
349                 cp += count;
350                 off = 0;
351                 m = m->m_next;
352           }
353     }
```

——————————————————————————————————— *[usr/src/sys/kern/uipc_mbuf.c]*

333-334 If offset or len is negative something is wrong and the kernel will panic.  335-342 All the mbufs in the chain prior to `off` offset bytes are skipped, if it happens that there are not enough bytes in the chain then the kernel will panic

342-352 The code runs through the mbuf chain and copies up to `len` bytes from the mbuf chain to the destionation buffer pointed to by `cp`, if the data finish before all the requested bytes have been copied then the kernel will panic.

——————————————————————————————————— *[usr/src/sys/net/rtsock.c]*

```
377     void
378     m_copyback(m0, off, len, cp)
379           struct     mbuf *m0;
380           register int off;
```

```
381            register int len;
382            caddr_t cp;
383     {
384            register int mlen;
385            register struct mbuf *m = m0, *n;
386            int totlen = 0;
387
388            if (m0 == 0)
389                    return;
390            while (off > (mlen = m->m_len)) {
391                    off -= mlen;
392                    totlen += mlen;
393                    if (m->m_next == 0) {
394                            n = m_getclr(M_DONTWAIT, m->m_type);
395                            if (n == 0)
396                                    goto out;
397                            n->m_len = min(MLEN, len + off);
398                            m->m_next = n;
399                    }
400                    m = m->m_next;
401            }
402            while (len > 0) {
403                    mlen = min (m->m_len - off, len);
404                    bcopy(cp, off + mtod(m, caddr_t), (unsigned)mlen);
405                    cp += mlen;
406                    len -= mlen;
407                    mlen += off;
408                    off = 0;
409                    totlen += mlen;
410                    if (len == 0)
411                            break;
412                    if (m->m_next == 0) {
413                            n = m_get(M_DONTWAIT, m->m_type);
414                            if (n == 0)
415                                    break;
416                            n->m_len = min(MLEN, len);
417                            m->m_next = n;
418                    }
419                    m = m->m_next;
420            }
421     out: if (((m = m0)->m_flags & M_PKTHDR) && (m->m_pkthdr.len < totlen))
422                    m->m_pkthdr.len = totlen;
423     }
```
_____ *[usr/src/sys/net/rtsock.c]*

388-389 If the mbuf chain is empty then the function returns

390-401 It skips off bytes along the mbuf chain, eventually allocating new zeroed mbufs that are appended to the end of the chain

402-420 len bytes are now copied from the linear buffer pointed by cp to the mbuf chain, allocating all eventually needed mbufs. 421-422 If the original first mbuf was a packet header, then the total packet length is adjusted in the packet header according to the new size.

The M_COPY_PKTHDR macro copies the information used by packet header mbufs (pkthdr.len and pkthdr.rcvif), sets the flag M_COPYFLAGS and the m_data pointer to point just immediately after the

packet header info.

*——————————————————————————————— [usr/src/sys/sys/uipc_mbuf.h]*

```
275     #define   M_COPY_PKTHDR(to, from) { \
276           (to)->m_pkthdr = (from)->m_pkthdr; \
277           (to)->m_flags = (from)->m_flags & M_COPYFLAGS; \
278           (to)->m_data = (to)->m_pktdat; \
279     }
```

*——————————————————————————————— [usr/src/sys/sys/uipc_mbuf.h]*

m_pullup rearranges an mbuf chain so that the first len bytes are stored contiguously inside the first mbuf of the chain, it is used to keep the headers contiguous, to allow the parsing of them with common pointer arithmetic.

*——————————————————————————————— usr/src/sys/kern/uipc_mbuf.c*

```
465     struct mbuf *
466     m_pullup(n, len)
467           register struct mbuf *n;
468           int len;
469     {
470           register struct mbuf *m;
471           register int count;
472           int space;
473
474           /*
475            * If first mbuf has no cluster, and has room for len bytes
476            * without shifting current data, pullup into it,
477            * otherwise allocate a new mbuf to prepend to the chain.
478            */
479           if ((n->m_flags & M_EXT) == 0 &&
480               n->m_data + len < &n->m_dat[MLEN] && n->m_next) {
481               if (n->m_len >= len)
482                   return (n);
483               m = n;
484               n = n->m_next;
485               len -= m->m_len;
486           } else {
487               if (len > MHLEN)
488                   goto bad;
489               MGET(m, M_DONTWAIT, n->m_type);
490               if (m == 0)
491                   goto bad;
492               m->m_len = 0;
493               if (n->m_flags & M_PKTHDR) {
494                   M_COPY_PKTHDR(m, n);
495                   n->m_flags &= ~M_PKTHDR;
496               }
497           }
498           space = &m->m_dat[MLEN] - (m->m_data + m->m_len);
499           do {
500               count = min(min(max(len, max_protohdr), space), n->m_len);
501               bcopy(mtod(n, caddr_t), mtod(m, caddr_t) + m->m_len,
502                   (unsigned)count);
503               len -= count;
```

```
504              m->m_len += count;
505              n->m_len -= count;
506              space -= count;
507              if (n->m_len)
508                   n->m_data += count;
509              else
510                   n = m_free(n);
511         } while (len > 0 && n);
512         if (len > 0) {
513              (void) m_free(m);
514              goto bad;
515         }
516         m->m_next = n;
517         return (m);
518   bad:
519         m_freem(n);
520         MPFail++;
521         return (0);
522    }
```
———————————————————————————————————————————————————— *usr/src/sys/kern/uipc_mbuf.c*

479-485 If the mbuf is not using a pagecluster, and there is enough room to keep all the `len` bytes inside it, then if the data in the first mbuf is already more than the data requested the mbuf pointer is returned without any further action. Otherwise `m->m_len` bytes are already in the first mbuf, and only `len-m->m_len` bytes need to be copied at the tail of the current data in the mbuf.

486-488 If the request was for more bytes than a packet header mbuf can keep (MHLEN = 100 ) then it tries to free the mbuf and returns an error.

489-491 It allocates a new mbuf of the same type, if the allocation is not successful then it returns an error.

492-497 The size of the data in the newly allocated mbuf is set to 0. If it was a packet header then the header data is copied from the old mbuf to this new one (pkthdr.len and pkthdr.rcvif), and then sets the m_data pointer to point just after the packet header information. Finally the M_PKTDHR flag is reset in the old mbuf, that will not be anymore a packet header.

The situation at this point is that we have allocated a new mbuf to prepend to the chain or not, in any case m points to what should be the first mbuf in the chain and n points to the second.

498    The space available at the end of the data in the `mbuf` is the difference between the two pointers, the one that points to the last byte of the mbuf ( `&m->m_dat[MLEN]` ) and the one that points to the last byte of the stored data ( `m->m_data+m->m_len` ).

499- It runs through the mbuf chain, copying until enough bytes have been copied, and eventually freeing the mbufs that have been completely exhausted. (m_free returns a pointer to the second mbuf in the chain if it exists)

512- If it was not possible to copy the requested number of bytes then the first mbuf is freed, and an error is returned. 516- This is the usual exit of the function. It links the first mbuf with the chain of other mbufs, and the complete chain is then returned. The dtom and mtod macros are frequently used :

———————————————————————————————————————————————————— *usr/src/sys/sys/mbuf.h*
```
61    #define   mtod(m,t)  ((t)((m)->m_data))
62    #define   dtom(x)        ((struct mbuf *)((long)(x) & ~(MSIZE-1)))
63    #define   mtocl(x)   (((u_long)(x) - (u_long)mbutl) >> MCLSHIFT)
64    #define   cltom(x)   ((caddr_t)((u_long)mbutl + ((u_long)(x) << MCLSHIFT)))
```
———————————————————————————————————————————————————— *usr/src/sys/sys/mbuf.h*

the mtod macro casts the m_data pointer to the correct type, the dtom macro gets the mbuf pointer from the

the m_data pointer truncating it to the previous 128-bytes(MSIZE) aligned address.

### 4.3. BSD Sharing of buffers

BSD allows only the sharing of external pages(page clusters). This is done keeping a reference count for each page cluster, the `mclrefcnt` array. The m_copy and m_free functions use this counter if the mbuf has an associated page cluster. m_copy instead of copying the data just increments the mclrefcnt counter and m_free puts the page cluster back in the free list only if the reference count becomes zero.

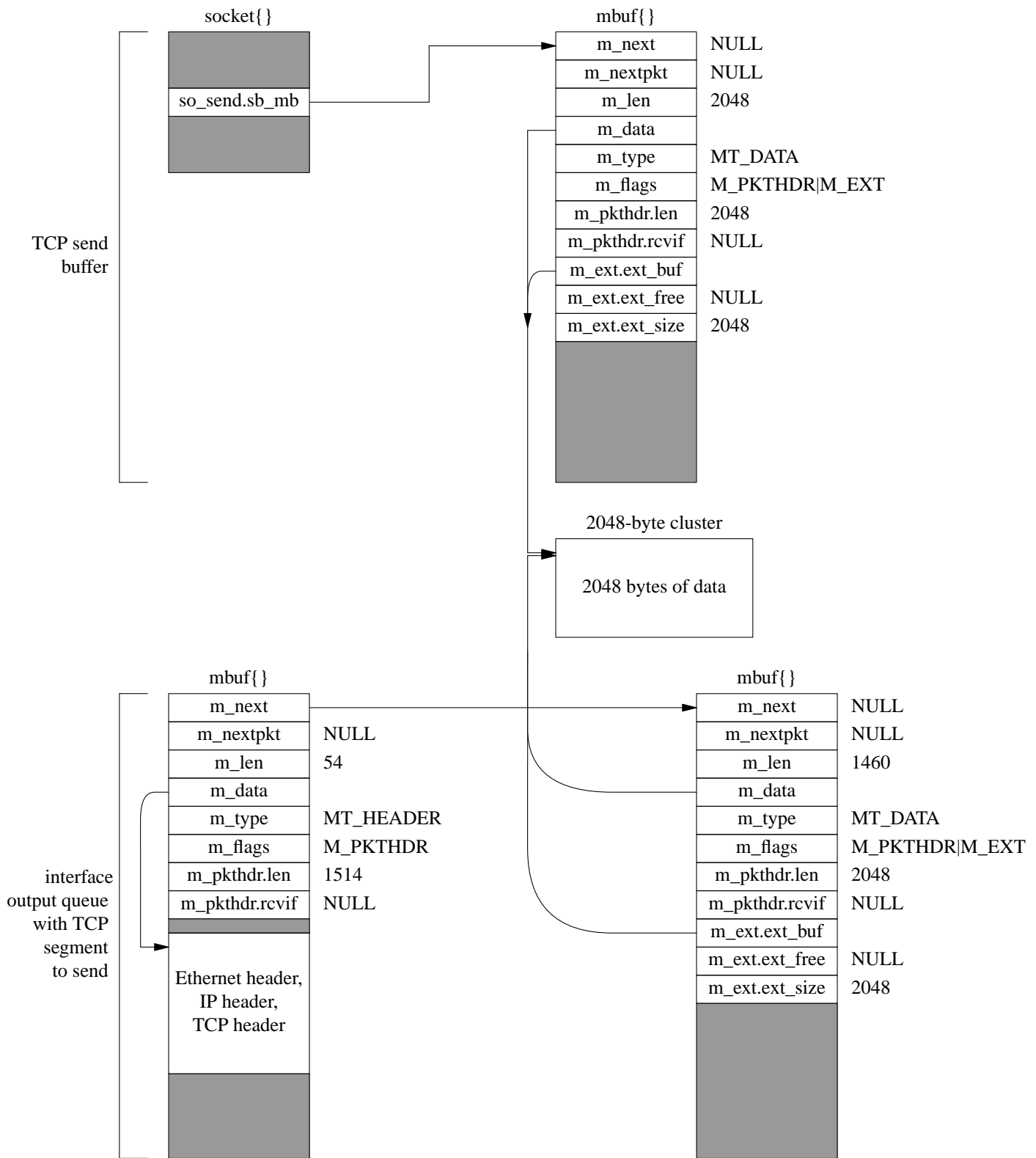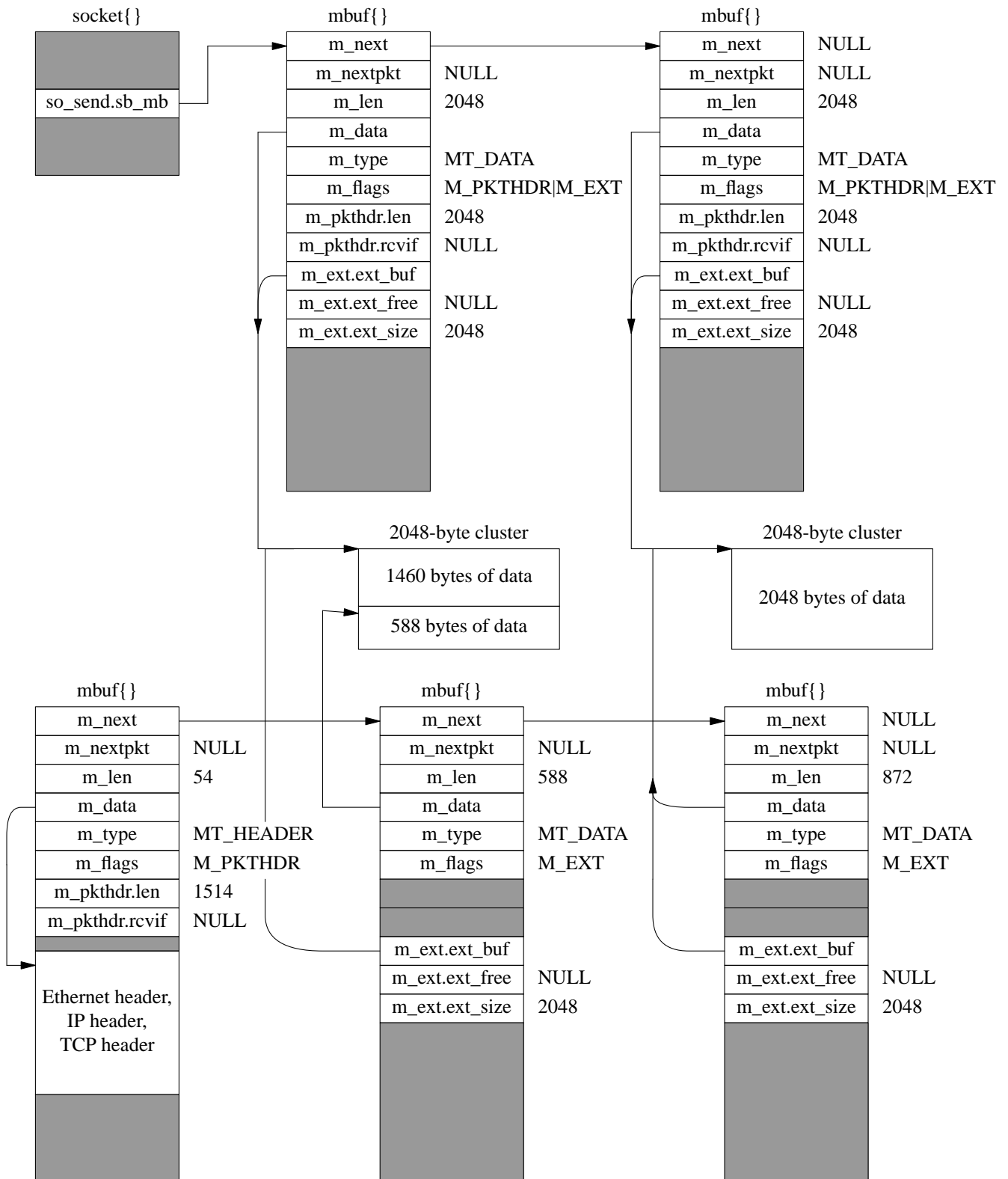socket{ }

mbuf{ }

| m_next | NULL |
| m_nextpkt | NULL |
| m_len | 2048 |
| m_data | |
| m_type | MT_DATA |
| m_flags | M_PKTHDR|M_EXT |
| m_pkthdr.len | 2048 |
| m_pkthdr.rcvif | NULL |
| m_ext.ext_buf | |
| m_ext.ext_free | NULL |
| m_ext.ext_size | 2048 |

so_send.sb_mb

TCP send buffer

2048-byte cluster

2048 bytes of data

mbuf{ }

mbuf{ }

| m_next | |
| m_nextpkt | NULL |
| m_len | 54 |
| m_data | |
| m_type | MT_HEADER |
| m_flags | M_PKTHDR |
| m_pkthdr.len | 1514 |
| m_pkthdr.rcvif | NULL |

Ethernet header, IP header, TCP header

interface output queue with TCP segment to send

| m_next | NULL |
| m_nextpkt | NULL |
| m_len | 1460 |
| m_data | |
| m_type | MT_DATA |
| m_flags | M_PKTHDR|M_EXT |
| m_pkthdr.len | 2048 |
| m_pkthdr.rcvif | NULL |
| m_ext.ext_buf | |
| m_ext.ext_free | NULL |
| m_ext.ext_size | 2048 |

*Figure 3 Sharing of a pagecluster*

| socket{ } | | mbuf{ } | | | mbuf{ } | |
|---|---|---|---|---|---|---|
| | | m_next | | | m_next | NULL |
| | | m_nextpkt | NULL | | m_nextpkt | NULL |
| so_send.sb_mb | | m_len | 2048 | | m_len | 2048 |
| | | m_data | | | m_data | |
| | | m_type | MT_DATA | | m_type | MT_DATA |
| | | m_flags | M_PKTHDR\|M_EXT | | m_flags | M_PKTHDR\|M_EXT |
| | | m_pkthdr.len | 2048 | | m_pkthdr.len | 2048 |
| | | m_pkthdr.rcvif | NULL | | m_pkthdr.rcvif | NULL |
| | | m_ext.ext_buf | | | m_ext.ext_buf | |
| | | m_ext.ext_free | NULL | | m_ext.ext_free | NULL |
| | | m_ext.ext_size | 2048 | | m_ext.ext_size | 2048 |

| 2048-byte cluster | | 2048-byte cluster |
|---|---|---|
| 1460 bytes of data | | 2048 bytes of data |
| 588 bytes of data | | |

| mbuf{ } | | | mbuf{ } | | | mbuf{ } | |
|---|---|---|---|---|---|---|---|
| m_next | | | m_next | | | m_next | NULL |
| m_nextpkt | NULL | | m_nextpkt | NULL | | m_nextpkt | NULL |
| m_len | 54 | | m_len | 588 | | m_len | 872 |
| m_data | | | m_data | | | m_data | |
| m_type | MT_HEADER | | m_type | MT_DATA | | m_type | MT_DATA |
| m_flags | M_PKTHDR | | m_flags | M_EXT | | m_flags | M_EXT |
| m_pkthdr.len | 1514 | | | | | | |
| m_pkthdr.rcvif | NULL | | | | | | |
| | | | m_ext.ext_buf | | | m_ext.ext_buf | |
| Ethernet header, IP header, TCP header | | | m_ext.ext_free | NULL | | m_ext.ext_free | NULL |
| | | | m_ext.ext_size | 2048 | | m_ext.ext_size | 2048 |

*Figure 4 Partial sharing of a pageclusters*

### 4.4.  Queues of packets and mbuf chains

As we have seen a packet on BSD is stored as a single mbuf or a chain of mbufs linked through the `m_next` pointer. Packets are linked on queues ( for instance the device or socket queue) through the `m_nextpkt` pointer of the only or the first mbuf of the chain.

| head of queue |
|---|
| tail of queue |

| mbuf{ } | | | |
|---|---|
| m_next | *next mbuf in chain* |
| m_nextpkt | |
| m_len | 42 |
| m_data | |
| m_type | MT_DATA |
| m_flags | M_PKTHDR |
| m_pkthdr.len | 192 |
| m_pkthdr_rcvif | NULL |

mbuf packet header

Ethernet header, IP header, UDP header

| mbuf{ } | | | |
|---|---|
| m_next | *next mbuf in chain* |
| m_nextpkt | NULL |
| m_len | 100 |
| m_data | |
| m_type | MT_DATA |
| m_flags | 0 |

100 bytes of data

| mbuf{ } | | |
|---|---|
| m_next | NULL |
| m_nextpkt | NULL |
| m_len | 50 |
| m_data | |
| m_type | MT_DATA |
| m_flags | 0 |

50 bytes of data

| mbuf{ } | | | |
|---|---|
| m_next | *next mbuf in chain* |
| m_nextpkt | |
| m_len | 54 |
| m_data | |
| m_type | MT_DATA |
| m_flags | M_PKTHDR |
| m_pkthdr.len | 1514 |
| m_pkthdr_rcvif | NULL |

mbuf packet header

Ethernet header, IP header, TCP header

| mbuf{ } | | |
|---|---|
| m_next | NULL |
| m_nextpkt | NULL |
| m_len | 1460 |
| m_data | |
| m_type | MT_DATA |
| m_flags | M_EXT |
| m_ext.ext_buf | |
| m_ext.ext_free | |
| m_ext.ext_size | 2048 |

2048-byte cluster

1460 bytes of data

*Figure 5 A UDP and a TCP packets enqueued, possibly at the device queue*

### 4.5. TCP packet creation

The code allocates a header mbuf for the ip and tcp header using the `MGETHDR` macro. The header is then created leaving in front of it a space sufficient for a maximum link header (this parameter is sysctl configurable with a lower bound of 16, sufficient for storing a 14 byte Ethernet header and keeping the IP header 16 byte aligned for efficiency reason).



*Figure 6 A TCP/IP packet with with headers in a prepended mbuf*

If the data in the mbuf chain can fit inside it, then it is copied inside this header mbuf. Otherwise the header mbuf is prepended to the mbuf chain created by m_copy (we know that this function doesnt copy external pageclusters, just copies the pointer to them).

—————————————————————————————————————————————————————————————— *[sys/netinet/tcp_output.c]*

```
360             MGETHDR(m, M_DONTWAIT, MT_HEADER);
361             if (m == NULL) {
362                     error = ENOBUFS;
363                     goto out;
364             }
365             m->m_data += max_linkhdr;
366             m->m_len = hdrlen;
367             if (len <= MHLEN - hdrlen - max_linkhdr) {
368                     m_copydata(so->so_snd.sb_mb, off, (int) len,
369                         mtod(m, caddr_t) + hdrlen);
370                     m->m_len += len;
371             } else {
372                     m->m_next = m_copy(so->so_snd.sb_mb, off, (int) len);
373                     if (m->m_next == 0) {
374                             (void) m_free(m);
375                             error = ENOBUFS;
376                             goto out;
377                     }
378             }
```

### 4.6.  UDP packet creation

The addition of IP and UDP headers is done prepending an mbuf to the mbuf chain where the udp and ip headers are placed at the end of its data region to allow for the addition of whatever header is eventually needed.



*Figure 7 A UDP/IP packet with headers*

```
411          /*
412           * Calculate data length and get a mbuf
413           * for UDP and IP headers.
414           */
415          M_PREPEND(m, sizeof(struct udpiphdr), M_DONTWAIT);
416          if (m == 0) {
417                  error = ENOBUFS;
418                  goto release;
419          }
420
421          /*
422           * Fill in mbuf with extended UDP header
423           * and addresses and length put into network format.
424           */
425          ui = mtod(m, struct udpiphdr *);
426          ui->ui_next = ui->ui_prev = 0;
427          ui->ui_x1 = 0;
428          ui->ui_pr = IPPROTO_UDP;
429          ui->ui_len = htons((u_short)len + sizeof (struct udphdr));
430          ui->ui_src = inp->inp_laddr;
431          ui->ui_dst = inp->inp_faddr;
432          ui->ui_sport = inp->inp_lport;
```

```
433        ui->ui_dport = inp->inp_fport;
434        ui->ui_ulen = ui->ui_len;
435
436        /*
437         * Stuff checksum and output datagram.
438         */
439        ui->ui_sum = 0;
440        if (udpcksum) {
441            if ((ui->ui_sum = in_cksum(m, sizeof (struct udpiphdr) + len)) == 0)
442              ui->ui_sum = 0xffff;
443        }
444        ((struct ip *)ui)->ip_len = sizeof (struct udpiphdr) + len;
445        ((struct ip *)ui)->ip_ttl = inp->inp_ip.ip_ttl;     /* XXX */
446        ((struct ip *)ui)->ip_tos = inp->inp_ip.ip_tos;     /* XXX */
447        udpstat.udps_opackets++;
448        error = ip_output(m, inp->inp_options, &inp->inp_route,
449            inp->inp_socket->so_options & (SO_DONTROUTE | SO_BROADCAST),
450            inp->inp_moptions);
```
———————————————————————————————————————————————— *[sys/netinet/udp_usrreq.c]*


### 4.7.  Ethernet header addition

The addition of the ethernet header happens in the `ether_output` function, where it uses the
M_PREPEND macro to verify that there are 14 bytes  available in front of the packet.

———————————————————————————————————————————————— *[sys/net/if_ethersubr.c]*
```
264        M_PREPEND(m, sizeof (struct ether_header), M_DONTWAIT);
265        if (m == 0)
266            senderr(ENOBUFS);
267        eh = mtod(m, struct ether_header *);
268        type = htons((u_short)type);
269        bcopy((caddr_t)&type,(caddr_t)&eh->ether_type,
270            sizeof(eh->ether_type));
271        bcopy((caddr_t)edst, (caddr_t)eh->ether_dhost, sizeof (edst));
272        bcopy((caddr_t)ac->ac_enaddr, (caddr_t)eh->ether_shost,
273            sizeof(eh->ether_shost));
```
———————————————————————————————————————————————— *[sys/net/if_ethersubr.c]*

This funtion is called inside [sys/netinet/ip_output.c] ip_output() as the if_output method for an ethernet
network device.

———————————————————————————————————————————————— *[sys/netinet/ip_output.c]*
```
279  sendit:
280        /*
281         * If small enough for interface, can just send directly.
282         */
283        if ((u_short)ip->ip_len <= ifp->if_mtu) {
284            ip->ip_len = htons((u_short)ip->ip_len);
285            ip->ip_off = htons((u_short)ip->ip_off);
286            ip->ip_sum = 0;
287            ip->ip_sum = in_cksum(m, hlen);
288            error = (*ifp->if_output)(ifp, m,
289                    (struct sockaddr *)dst, ro->ro_rt);
```

```
290                goto done;
291            }
```
——————————————————————————————————————————— *[sys/netinet/ip_output.c]*


The linearization of the mbuf chains happens in BSD for example in the `leput` (Lance ethernet put ) function for an ethernet interface.  This function  in the BSD code copies all the data from the mbuf chain contiguously to the hardware buffer.

——————————————————————————————————————————— *[sys/hp300/dev/if_le.c]*

```
717    /*
718     * Routine to copy from mbuf chain to transmit
719     * buffer in board local memory.
720     */
721    leput(lebuf, m)
722        register char *lebuf;
723        register struct mbuf *m;
724    {
725        register struct mbuf *mp;
726        register int len, tlen = 0;
727
728        for (mp = m; mp; mp = mp->m_next) {
729            len = mp->m_len;
730            if (len == 0)
731                continue;
732            tlen += len;
733            bcopy(mtod(mp, char *), lebuf, len);
734            lebuf += len;
735        }
736        m_freem(m);
737        if (tlen < LEMINSIZE) {
738            bzero(lebuf, LEMINSIZE - tlen);
739            tlen = LEMINSIZE;
740        }
741        return(tlen);
742    }
```
——————————————————————————————————————————— *[sys/hp300/dev/if_le.c]*


### 4.8.  IP fragmentation/defragmentation in BSD

IP fragments are recognized by having the MF bit set to 1 or the IP offset field different from 0 in the IP header.  BSD stores such datagrams in a queue specific for the IP source and destination addresses. The IP source and destination addresses are stored in the head of the queue.

——————————————————————————————————————————— *[sys/netinet/ip_var.h]*

```
54    struct ipq {
55        struct     ipq *next,*prev;      /* to other reass headers */
56        u_char     ipq_ttl;         /* time for reass q to live */
57        u_char     ipq_p;                 /* protocol of this fragment */
58        u_short    ipq_id;                /* sequence id for reassembly */
59        struct     ipasfrag *ipq_next,*ipq_prev;
60                        /* to ip headers of fragments */
61        struct     in_addr ipq_src,ipq_dst;
```

```
  62    };
```

The IP address fields inside the IP packets are re-used to store backward and forward pointers between IP fragments to keep them linked in a doubly linked list. The mbuf structure can easily be retrieved from the m->data pointer by simply truncating it to a multiple of the mbuf size (MSIZE=128), as it is done by the dtom macro :

```
  62    #define  dtom(x)          ((struct mbuf *)((long)(x) & ~(MSIZE-1)))
```

There is a problem when the data is stored in an external page as in an mbuf cluster. In that case knowing the address of the IP header would not be sufficient to retrieve the mbuf header using the dtom function (this function because of this is deprecated).

mbuf{ }

| | |
|---|---|
| m_next | NULL |
| m_nextpkt | NULL |
| m_len | 296 |
| m_data | |
| m_type | MT_DATA |
| m_flags | M_PKTHDR\|M_EXT |
| m_pkthdr.len | 296 |
| m_pkthdr.rcvif | ptr |
| m_ext.ext_buf | |
| m_ext.ext_free | NULL |
| m_ext.ext_size | 2048 |

doubly linked list
of IP fragments

doubly linked lis
of IP fragments

IP hdr

276 bytes of data

2048-byte cluster

*Figure 8 mbuf_five.pic*

For this reason, in such a case the `m_pullup` function is called to eventually copy the IP header in a newly allocated mbuf that will be prepended to the one being processed.

| mbuf{ }            |          |     | mbuf{ }            |          |
|--------------------|----------|-----|--------------------|----------|
| m_next             |          |     | m_next             | NULL     |
| m_nextpkt          | NULL     |     | m_nextpkt          | NULL     |
| m_len              | 40       |     | m_len              | 256      |
| m_data             |          |     | m_data             |          |
| m_type             | MT_DATA  |     | m_type             | MT_DATA  |
| m_flags            | M_PKTHDR |     | m_flags            | M_EXT    |
| m_pkthdr.len       | 296      |     |                    |          |
| m_pkthdr.rcvif     | ptr      |     |                    |          |
|                    |          |     | m_ext.ext_buf      |          |
| IP header          |          |     | m_ext.ext_free     | NULL     |
|                    |          |     | m_ext.ext_size     | 2048     |
| next 20 bytes      |          |     |                    |          |
| of datagram        |          |     |                    |          |

doubly linked list of IP datagram

doubly linked list of IP datagram

next 256 bytes of datagram

2048-byte cluster

*Figure 9 mbuf_six.pic*

It is the ip_reass function in the ip_input.c file that takes care of managing IP fragments. An mbuf (or mbuf chain) with an IP fragment is stored on an ip fragment queue after having decreased its size (m_len) by the amount of the IP header length and adjusted its `m_data` pointer to point after the IP header. When ip_reass detects the queue of fragments is complete, then it goes through the links stored in the ip headers of the packets, concatenates the mbufs throught their m_next pointer(this reduces an eventual list of mbuf chains to a simple single mbuf chain), restores the IP source and destination in the IP header of the first mbuf, adjusts the m_data of the first mbuf to expose again the IP header and its m_len to comprise the IP header length. Then the packet len is computed as the sum of the lengths of the fragments and stored in the m_pkthdr.len of the first mbuf. Then it returns the mbuf chain.

──────────────────────────────────────────────────────────────────────── *[sys/netinet/ip_input.c]*

```
 140     void
```

```
141    ipintr()
142    {
143          register struct ip *ip;
144          register struct mbuf *m;
145          register struct ipq *fp;
146          register struct in_ifaddr *ia;
147          int hlen, s;
148
149    next:
150          /*
151           * Get next datagram off input queue and get IP header
152           * in first mbuf.
153           */
154          s = splimp();
155          IF_DEQUEUE(&ipintrq, m);
156          splx(s);
157          if (m == 0)
158                return;
159    #ifdef    DIAGNOSTIC
160          if ((m->m_flags & M_PKTHDR) == 0)
161                panic("ipintr no HDR");
162    #endif
163          /*
164           * If no IP addresses have been set yet but the interfaces
165           * are receiving, can't do anything with incoming packets yet.
166           */
167          if (in_ifaddr == NULL)
168                goto bad;
169          ipstat.ips_total++;
170          if (m->m_len < sizeof (struct ip) &&
171              (m = m_pullup(m, sizeof (struct ip))) == 0) {
172                ipstat.ips_toosmall++;
173                goto next;
174          }
175          ip = mtod(m, struct ip *);
176          if (ip->ip_v != IPVERSION) {
177                ipstat.ips_badvers++;
178                goto bad;
179          }
180          hlen = ip->ip_hl << 2;
181          if (hlen < sizeof(struct ip)) {/* minimum header length */
182                ipstat.ips_badhlen++;
183                goto bad;
184          }
185          if (hlen > m->m_len) {
186                if ((m = m_pullup(m, hlen)) == 0) {
187                      ipstat.ips_badhlen++;
188                      goto next;
189                }
190                ip = mtod(m, struct ip *);
191          }
192          if (ip->ip_sum = in_cksum(m, hlen)) {
193                ipstat.ips_badsum++;
194                goto bad;
```

```
195                }
196
197                /*
198                 * Convert fields to host representation.
199                 */
200                NTOHS(ip->ip_len);
201                if (ip->ip_len < hlen) {
202                        ipstat.ips_badlen++;
203                        goto bad;
204                }
205                NTOHS(ip->ip_id);
206                NTOHS(ip->ip_off);
207
208                /*
209                 * Check that the amount of data in the buffers
210                 * is as at least much as the IP header would have us expect.
211                 * Trim mbufs if longer than we expect.
212                 * Drop packet if shorter than we expect.
213                 */
214                if (m->m_pkthdr.len < ip->ip_len) {
215                        ipstat.ips_tooshort++;
216                        goto bad;
217                }
218                if (m->m_pkthdr.len > ip->ip_len) {
219                        if (m->m_len == m->m_pkthdr.len) {
220                                m->m_len = ip->ip_len;
221                                m->m_pkthdr.len = ip->ip_len;
222                        } else
223                                m_adj(m, ip->ip_len - m->m_pkthdr.len);
224                }
225
226                /*
227                 * Process options and, if not destined for us,
228                 * ship it on.  ip_dooptions returns 1 when an
229                 * error was detected (causing an icmp message
230                 * to be sent and the original packet to be freed).
231                 */
232                ip_nhops = 0;           /* for source routed packets */
233                if (hlen > sizeof (struct ip) && ip_dooptions(m))
234                        goto next;
235
236                /*
237                 * Check our list of addresses, to see if the packet is for us.
238                 */
239                for (ia = in_ifaddr; ia; ia = ia->ia_next) {
240     #define    satosin(sa)     ((struct sockaddr_in *)(sa))
241
242                        if (IA_SIN(ia)->sin_addr.s_addr == ip->ip_dst.s_addr)
243                                goto ours;
244                        if (
245     #ifdef     DIRECTED_BROADCAST
246                        ia->ia_ifp == m->m_pkthdr.rcvif &&
247     #endif
248                        (ia->ia_ifp->if_flags & IFF_BROADCAST)) {
```

```
249                      u_long t;
250
251                      if (satosin(&ia->ia_broadaddr)->sin_addr.s_addr ==
252                          ip->ip_dst.s_addr)
253                            goto ours;
254                      if (ip->ip_dst.s_addr == ia->ia_netbroadcast.s_addr)
255                            goto ours;
256                      /*
257                       * Look for all-0's host part (old broadcast addr),
258                       * either for subnet or net.
259                       */
260                      t = ntohl(ip->ip_dst.s_addr);
261                      if (t == ia->ia_subnet)
262                            goto ours;
263                      if (t == ia->ia_net)
264                            goto ours;
265                  }
266          }
267      if (IN_MULTICAST(ntohl(ip->ip_dst.s_addr))) {
268              struct in_multi *inm;
269  #ifdef MROUTING
270              extern struct socket *ip_mrouter;
271
272              if (ip_mrouter) {
273                      /*
274                       * If we are acting as a multicast router, all
275                       * incoming multicast packets are passed to the
276                       * kernel-level multicast forwarding function.
277                       * The packet is returned (relatively) intact; if
278                       * ip_mforward() returns a non-zero value, the packet
279                       * must be discarded, else it may be accepted below.
280                       *
281                       * (The IP ident field is put in the same byte order
282                       * as expected when ip_mforward() is called from
283                       * ip_output().)
284                       */
285                      ip->ip_id = htons(ip->ip_id);
286                      if (ip_mforward(m, m->m_pkthdr.rcvif) != 0) {
287                          ipstat.ips_cantforward++;
288                          m_freem(m);
289                          goto next;
290                      }
291                      ip->ip_id = ntohs(ip->ip_id);
292
293                      /*
294                       * The process-level routing demon needs to receive
295                       * all multicast IGMP packets, whether or not this
296                       * host belongs to their destination groups.
297                       */
298                      if (ip->ip_p == IPPROTO_IGMP)
299                          goto ours;
300                      ipstat.ips_forward++;
301              }
302  #endif
```

```
303                /*
304                 * See if we belong to the destination multicast group on the
305                 * arrival interface.
306                 */
307                IN_LOOKUP_MULTI(ip->ip_dst, m->m_pkthdr.rcvif, inm);
308                if (inm == NULL) {
309                    ipstat.ips_cantforward++;
310                    m_freem(m);
311                    goto next;
312                }
313            goto ours;
314        }
315        if (ip->ip_dst.s_addr == (u_long)INADDR_BROADCAST)
316            goto ours;
317        if (ip->ip_dst.s_addr == INADDR_ANY)
318            goto ours;
319
320        /*
321         * Not for us; forward if possible and desirable.
322         */
323        if (ipforwarding == 0) {
324            ipstat.ips_cantforward++;
325            m_freem(m);
326        } else
327            ip_forward(m, 0);
328        goto next;
329
330    ours:
331        /*
332         * If offset or IP_MF are set, must reassemble.
333         * Otherwise, nothing need be done.
334         * (We could look in the reassembly queue to see
335         * if the packet was previously fragmented,
336         * but it's not worth the time; just let them time out.)
337         */
338        if (ip->ip_off &~ IP_DF) {
339            if (m->m_flags & M_EXT) {        /* XXX */
340                if ((m = m_pullup(m, sizeof (struct ip))) == 0) {
341                    ipstat.ips_toosmall++;
342                    goto next;
343                }
344                ip = mtod(m, struct ip *);
345            }
346            /*
347             * Look for queue of fragments
348             * of this datagram.
349             */
350            for (fp = ipq.next; fp != &ipq; fp = fp->next)
351                if (ip->ip_id == fp->ipq_id &&
352                    ip->ip_src.s_addr == fp->ipq_src.s_addr &&
353                    ip->ip_dst.s_addr == fp->ipq_dst.s_addr &&
354                    ip->ip_p == fp->ipq_p)
355                    goto found;
356            fp = 0;
```

```
357    found:
358
359               /*
360                * Adjust ip_len to not reflect header,
361                * set ip_mff if more fragments are expected,
362                * convert offset of this to bytes.
363                */
364               ip->ip_len -= hlen;
365               ((struct ipasfrag *)ip)->ipf_mff &= ~1;
366               if (ip->ip_off & IP_MF)
367                     ((struct ipasfrag *)ip)->ipf_mff |= 1;
368               ip->ip_off <<= 3;
369
370               /*
371                * If datagram marked as having more fragments
372                * or if this is not the first fragment,
373                * attempt reassembly; if it succeeds, proceed.
374                */
375               if (((struct ipasfrag *)ip)->ipf_mff & 1 || ip->ip_off) {
376                     ipstat.ips_fragments++;
377                     ip = ip_reass((struct ipasfrag *)ip, fp);
378                     if (ip == 0)
379                           goto next;
380                     ipstat.ips_reassembled++;
381                     m = dtom(ip);
382               } else
383                     if (fp)
384                           ip_freef(fp);
385          } else
386               ip->ip_len -= hlen;
387
388          /*
389           * Switch out to protocol's input routine.
390           */
391          ipstat.ips_delivered++;
392          (*inetsw[ip_protox[ip->ip_p]].pr_input)(m, hlen);
393          goto next;
394     bad:
395          m_freem(m);
396          goto next;
397     }

405     struct ip *
406     ip_reass(ip, fp)
407          register struct ipasfrag *ip;
408          register struct ipq *fp;
409     {
410          register struct mbuf *m = dtom(ip);
411          register struct ipasfrag *q;
412          struct mbuf *t;
413          int hlen = ip->ip_hl << 2;
414          int i, next;
415
416          /*
```

```
417              * Presence of header sizes in mbufs
418              * would confuse code below.
419              */
420             m->m_data += hlen;
421             m->m_len -= hlen;
422
423             /*
424              * If first fragment to arrive, create a reassembly queue.
425              */
426             if (fp == 0) {
427                     if ((t = m_get(M_DONTWAIT, MT_FTABLE)) == NULL)
428                             goto dropfrag;
429                     fp = mtod(t, struct ipq *);
430                     insque(fp, &ipq);
431                     fp->ipq_ttl = IPFRAGTTL;
432                     fp->ipq_p = ip->ip_p;
433                     fp->ipq_id = ip->ip_id;
434                     fp->ipq_next = fp->ipq_prev = (struct ipasfrag *)fp;
435                     fp->ipq_src = ((struct ip *)ip)->ip_src;
436                     fp->ipq_dst = ((struct ip *)ip)->ip_dst;
437                     q = (struct ipasfrag *)fp;
438                     goto insert;
439             }
440
441             /*
442              * Find a segment which begins after this one does.
443              */
444             for (q = fp->ipq_next; q != (struct ipasfrag *)fp; q = q->ipf_next)
445                     if (q->ip_off > ip->ip_off)
446                             break;
447
448             /*
449              * If there is a preceding segment, it may provide some of
450              * our data already.  If so, drop the data from the incoming
451              * segment.  If it provides all of our data, drop us.
452              */
453             if (q->ipf_prev != (struct ipasfrag *)fp) {
454                 i = q->ipf_prev->ip_off + q->ipf_prev->ip_len - ip->ip_off;
455                 if (i > 0) {
456                         if (i >= ip->ip_len)
457                                 goto dropfrag;
458                         m_adj(dtom(ip), i);
459                         ip->ip_off += i;
460                         ip->ip_len -= i;
461                 }
462             }
463
464             /*
465              * While we overlap succeeding segments trim them or,
466              * if they are completely covered, dequeue them.
467              */
468             while (q != (struct ipasfrag *)fp && ip->ip_off + ip->ip_len > q->ip_off)
{
469                     i = (ip->ip_off + ip->ip_len) - q->ip_off;
```

```
470              if (i < q->ip_len) {
471                      q->ip_len -= i;
472                      q->ip_off += i;
473                      m_adj(dtom(q), i);
474                      break;
475              }
476              q = q->ipf_next;
477              m_freem(dtom(q->ipf_prev));
478              ip_deq(q->ipf_prev);
479         }
480
481   insert:
482         /*
483          * Stick new segment in its place;
484          * check for complete reassembly.
485          */
486         ip_enq(ip, q->ipf_prev);
487         next = 0;
488         for (q = fp->ipq_next; q != (struct ipasfrag *)fp; q = q->ipf_next) {
489              if (q->ip_off != next)
490                      return (0);
491              next += q->ip_len;
492         }
493         if (q->ipf_prev->ipf_mff & 1)
494              return (0);
495
496         /*
497          * Reassembly is complete; concatenate fragments.
498          */
499         q = fp->ipq_next;
500         m = dtom(q);
501         t = m->m_next;
502         m->m_next = 0;
503         m_cat(m, t);
504         q = q->ipf_next;
505         while (q != (struct ipasfrag *)fp) {
506              t = dtom(q);
507              q = q->ipf_next;
508              m_cat(m, t);
509         }
510
511         /*
512          * Create header for new ip packet by
513          * modifying header of first packet;
514          * dequeue and discard fragment reassembly header.
515          * Make header visible.
516          */
517         ip = fp->ipq_next;
518         ip->ip_len = next;
519         ip->ipf_mff &= ~1;
520         ((struct ip *)ip)->ip_src = fp->ipq_src;
521         ((struct ip *)ip)->ip_dst = fp->ipq_dst;
522         remque(fp);
523         (void) m_free(dtom(fp));
```

```
524          m = dtom(ip);
525          m->m_len += (ip->ip_hl << 2);
526          m->m_data -= (ip->ip_hl << 2);
527          /* some debugging cruft by sklower, below, will go away soon */
528          if (m->m_flags & M_PKTHDR) { /* XXX this should be done elsewhere */
529              register int plen = 0;
530              for (t = m; m; m = m->m_next)
531                  plen += m->m_len;
532              t->m_pkthdr.len = plen;
533          }
534          return ((struct ip *)ip);
535
536  dropfrag:
537          ipstat.ips_fragdropped++;
538          m_freem(m);
539          return (0);
540  }
```
*—————————————————————————————————————————— [sys/netinet/ip_input.c]*


## 5. Unix SVR4 STREAMS

STREAMS is a general programming model to develop communication services on Unix. It defines the programming interface between the various components of the programming model: a stream head, zero or more modules and a device driver. In this model a user process interacts with a STREAMS queue head, this one can interact directly or through a stack of modules with a device driver.

*Figure 10 stream.pic*

All the communications between STREAMS components are done through `messages`. A `message` in STREAMS is described by one or more message blocks (msgb) linked through the `b_cont` pointer. Messages are kept on queues in a doubly linked list through the `b_next`, `b_prev` pointers of the first message block(msgb) of the message.

*Figure 11 stream_two.pic*

```
struct msgb {
    struct msgb     *b_next;         /*next msg on queue*/
    struct msgb     *b_prev;         /*previous msg on queue*/
    struct msgb     *b_cont;         /*next msg block of message*/
    unsigned char   *b_rptr;         /*1st unread byte in bufr*/
    unsigned char   *b_wptr;         /*1st unwritten byte in bufr*/
    struct datab    *b_datap;        /*data block*/
    unsigned char   b_band;          /*message priority*/
    unsigned short  b_flag;          /*see below - Message flags*/
};
typedef struct msgb mblk_t;
```

The data, kept in a `data buffer` is described by the data block( `datab` ) that is accessed through the `b_datap` pointer in the message block.

```
struct datab {
    unsigned char   *db_base;        /* first byte of buffer */
    unsigned char   *db_lim;         /* last byte+1 of buffer */
    unsigned char   db_ref;          /* msg count ptg to this blk */
    unsigned char   db_type;         /* msg type */
};
typedef struct datab dblk_t;
```

*Figure 12 stream_one.pic*

A mechanism to avoid unnecessary copies is provided through the sharing of the data block, where a refer-enc counter `db_ref` is kept and incremented for each new entity that references the same data block. This is called duplication on STREAMS documentation, and is performed through the use of the `dupmsg` or `dupb` function.

*Figure 13 A shared data block*

With Unix STREAMS a module can create a new message block that shares the data block and data buffer of an existing one. The read and write pointers (b_rptr, b_wptr) are kept inside the message block and so two entities can share the data block and the associated data buffer but have a different viewport on the data. For instance a tcp module and an ip module can share a data block and buffer even if the tcp module will not consider the IP header.

### 5.1.  STREAMS functions

Message blocks (mblk) are allocated through the  :

```
mblk_t *allocb(size_t size,uint_t pri);
```

where size is the number of data bytes available in the message block, and pri is the priority of the request. This call allocates all the necessary structures : the msgb, the datab and a conveniently sized data buffer and it appropriately initializes all the fields. This example sends an error code of 1 byte on a stream :

```
1  send_error(q,err)
2    queue_t *q;
3    unsigned char err;
4  {
5    mblk_t *bp;
6
7    if ((bp = allocb(1, BPRI_HI)) == NULL) /* allocate msg. block */
8        return(0);
9
10   bp->b_datap->db_type = M_ERROR;    /* set msg type to M_ERROR */
11   *bp->b_wptr++ = err;               /* increment write pointer */
12
13   if (!(q->q_flag & QREADR))     /* if not read queue     */
14       q = RD(q);                 /*    get read queue     */
15   putnext(q,bp);                 /* send message upstream */
16   return(1);
17 }
```

A message block can be allocated for an existing data buffer. In this case the esballoc (extended streams buffer allocation) function is used :

```
mblk_t *esballoc(uchar *base, size_t size, uint_t pri,frtn_t *fr_rtnp);
```

that allocates only the msgb and datab, and sets up the fields to point to the user supplied base buffer (it is used for instance with devices having dual-ported memory to avoid copying the data to a kernel buffer). The caller can in this case specify a message specific free function that will be called to deallocate the data buffer when it will be released. The allocation of memory for a message block in STREAMS can be cumbersome. Three different areas need to be allocated : a msgb, a datab and a data buffer. In SVR4 memory allocation is done in a smart way through the use of 128 bytes mdbblocks.

```
struct mdbblock {
    mblk_t m_mblock;
    void   (*m_func)(void);
    dblk_t d_dblock;
    char   databuf[FASTBUF];
}
```

These areas of memory can keep a msgb and a datab, and eventually the remaining space can be used for a small data buffer. Therefore the allocation of a msgb for a small buffer size reduces to only one memory allocation.

*Figure 14 An mddblock structure just after allocation*

In case the size requested is more than the free space available in an mdbblock then an external buffer of proper size is allocate with the standard kmem_alloc() function. LiS (Linux STREAMS) is an implementation of STREAMS for Linux that follows the SVR4 memory allocation design. In LiS the size of the mdbblock internal data buffer(FASTBUF) on Intel x86 is 78 bytes.

Message blocks can be linked through their b_cont field to form a chain using the function :

```
void linkb(mblk_t* mp,mblk_t* bp);
```

This function concatenates the 2 messages putting the bp message at the tail of the mp message.
Message blocks can be copied with copyb :

```
struct msgb *bp = copyb(struct msgb *mp);
```

this will allocate a new msgb and datab and data buffer and will copy over the data between the rptr and wptr in the old msgb. If it is successfull it returns a pointer to the new msgb otherwise it returns NULL. This function will not follow the b_cont pointer, and will copy a single msgb. The copymsg function instead copies all the msgb in a message following the b_cont links :

```
mblk_t* copymsg(mblk_t* bp);
```

Example 1: : Using copyb

For each message in the list, test to see if the  downstream
queue is full with the canputnext(9F) function (line 21). If
it is not full, use copyb to copy a  header  message  block,
and dupmsg(9F) to duplicate the data to be retransmitted. If
either operation fails, reschedule a  timeout  at  the  next
valid interval.

Update the new header block  with  the  correct  destination
address  (line  34),  link  the  message to it (line 35), and
send it downstream (line  36).  At  the  end  of  the  list,
reschedule this routine.

```
 1  struct retrans {
 2        mblk_t                *r_mp;
 3        int                   r_address;
 4        queue_t               *r_outq;
 5        struct retrans     *r_next;
 6  };
 7
 8  struct protoheader {
      ...
 9    int                     h_address;
      ...
10  };
11
12  mblk_t *header;
13
14  void
15  retransmit(struct retrans *ret)
16  {
17        mblk_t *bp, *mp;
18        struct protoheader *php;
19
20        while (ret) {
21          if (!canputnext(ret->r_outq)) {    /* no room */
22                ret = ret->r_next;
23                continue;
24          }
25          bp = copyb(header);                /* copy header msg. block */
26          if (bp == NULL)
27                break;
28          mp = dupmsg(ret->r_mp);            /* duplicate data */
29          if (mp == NULL) {                  /* if unsuccessful */
30              freeb(bp);                     /* free the block */
31              break;
32          }
33          php = (struct protoheader *)bp->b_rptr;
34          php->h_address = ret->r_address;   /* new header */
35          bp->bp_cont = mp;                  /* link the message */
36          putnext(ret->r_outq, bp);          /* send downstream */
37          ret = ret->r_next;
38        }
```

```
39        /* reschedule */
40        (void) timeout(retransmit, (caddr_t)ret, RETRANS_TIME);
41  }
```

If its not necessary to modify the data in the buffer, the dupb function can be called to allocate a new msgb that points to the same datab, and so shares the data buffer :

```
mblk_t* dupb(mblk_t bp);
mblk_t *dupmsg(mblk_t *mp);
```

The dupb function doesn't follow the b_cont link and does its job only on one msgb. The dupmsg follows the b_cont link and so duplicates all the msgb of a chain.  Message blocks can be freed with freeb :

```
void  freeb(mblk_t *bp);
void freemsg(mblk_t* bp);
```

The freeb function decrements the reference counter of the datab, if it becomes zero it really gives back the memory for the data buffer and the datab structure, then in any case deallocates the msgb area.  The freeb function doesn't follow the b_cont link, while the freemsg function can be used to deallocate all the parts of a message and this is done following the b_cont link.
A message can be trimmed at the head or the tail using the adjmsg() function :

```
int adjmsg(mblk_t* mp, int len)
```

where if len is positive len bytes are trimmed from the beginning and if negative len bytes are trimmed from the tail (this can be used by the protocol layers to strip headers or trailers).  The general functions to get or put a message on a queue are :

```
int putq(queue_t q,mblk_t* bp);
mblk_t*  getq(queue_t q);
int insq(queue_t* q,mblk_t* emp,mblk_t* mp);
```

putq puts a message on the specified queue based on its priority (db_type/b_band), getq gets the next available message and insq inserts the mp message immediately before emp in the queue q.
The memory allocation can't sleep waiting for resources, so it can return with a NULL, meaning that there is no memory available.  In this case it is possible to register a callback function that will be invoked as soon as a buffer of the requested size will be available.

```
int bufcall(uint size,int pri, void (*func)(), long arg);
```

The STREAMS framework will invoke (*func)() at the time it decides memory is available again.
Each STREAMS stream head, module or driver has to define a put procedure (pointed to by the  qi_putp pointer in the qinit structure):

```
int put(queue_t* q, mblk_t* mp);
```

for each of the 2 queues that process incoming messages and eventually forward them.  All stream heads and eventually modules or drivers can also define a service procedure (pointed to by the qi_srvp pointer in the qinit structure) :

```
int service(queue* q);
```

to process all deferred messages, when the STREAMS scheduler will detect that a blocking condition is resolved and will invoke it.
In the STREAMS framework TCP/IP has been provided through a message based interface called TPI(Transport Provider Interface) and through the procedural TLI (Transport Layer Interface) APIs.

*Figure 15 stream_seven.pic*

As the socket interface came long before TPI/TLI and was widely diffused it has been provided on top of the STREAMS TCP, UDP and IP modules through a STREAMS sockmod module and a library.

STREAMS entities that can queue messages have a service procedure defined. There are two functions that can check if there is room in a queue. If the queue doesn't have a service function then they will search for one along the stream in the same direction :

```
int bcanput(queue_t* q,unsigned char pri);
int canput(queue_t* q);
```

bcanput checks if there is room for one message at priority pri, canput is equivalent to bcanput(q,0).
An operation that is frequently performed is that of sending a message to the next module (or driver or stream head) in a stream. This is done invoking the putp function registered by that module in the qinit structure. This operation can be conveniently boiled down in a  putnext macro :

```
#define putnext(q, mp) ((*(q)->q_next->q_qinfo->qi_putp)((q)->next,(mp)))
```

Many operations on headers and trailers are done through some pointer arithmetic (for instance the source address field is at a certain offset from the network header). These operations on pointers can be done only if that area of memory is contiguous. Therefore a special utility function is used to copy in a contiguous data buffer a specified number of bytes from a chain of message blocks.

```
int pullupmsg(struct msgb* mp,int len);
```

This function returns a chain of msgb in which the first msgb contains len bytes of data in its contiguous data buffer. If len is -1 then the utility concatenates all the blocks of the same type at the beginning of the message. In the LiS implementation the first msgb is reused while a new datab and data buffer are allocated for the first block.
STREAMS queues are always allocated in pairs, one is the read side and the other the write side, the first

passes messages downstream (towards the driver), the second passes messages upstream (toward the stream head).  Inside the queue structure there is a pointer q_other that connects each other, and makes it simple to get the read or write side or the other side of the queue pair.

```
taken from the LiS implementation for Linux :
typedef
struct queue
{
  EXPORT
        struct  qinit   *q_qinfo;        /* procs and limits for queue [I]*/
        struct  msgb    *q_first;        /* first data block [Z]*/
        struct  msgb    *q_last;         /* last data block [Z]*/
        struct  queue   *q_next;         /* next Q of stream [Z]*/
        struct  queue   *q_link;         /* to next Q for the scan list [Z]*/
        void            *q_ptr;          /* module private data for free */
        ulong           q_count;         /* number of bytes on Q [Z]*/
        volatile ulong  q_flag;          /* queue state [Z]*/
  SHARE
        long            q_minpsz;        /* min packet size accepted [I]*/
        long            q_maxpsz;        /* max packet size accepted [I]*/
        ulong           q_hiwat;         /* queue high water mark [I]*/
        ulong           q_lowat;         /* queue low water mark [I]*/
  PRIVATE
        struct qband    *q_bandp;        /* separate flow information */
        struct queue    *q_other;        /* for RD()/WR()/OTHER() */
        void            *q_str;          /* pointer to stream's stdata */
        struct  queue   *q_scnxt;        /* next q in the scan list */
        ulong            q_magic;        /* magic number */
        lis_semaphore_t  q_lock;         /* for put, srv, open, close */
        lis_atomic_t     q_lock_nest;    /* for nested entries */
        void            *q_taskp;        /* owner of the q_lock */
        lis_spin_lock_t  q_isr_lock;     /* for ISR protection */
        lis_semaphore_t *q_wakeup_sem ;  /* helps sync closes */
} queue_t;
```

Special calls are defined to perform these operations :

```
        queue_t RD(queue_t* q);
        queue_t WR(queue_t* q);
        queue_t OTHERQ(queue_t* q);
```

RD returns a pointer to the read side, WR returns a pointer to the write side and OTHERQ returns a pointer to the other side of the queue pair whatever it is. Read or  write side is easily recognized because of the QREADR flag set or unset in the q_flag field of the queue structure respectively.  The qinit structure pointed to by the q_qinfo pointer in the queue structure, is a table of pointers to functions defined to be used with the queue and a structure with info about the module and another with some statistical data.  The only required procedure is the put procedure, if the module implements queueing, then at least the service procedure should also be defined.

```
typedef
struct  qinit {
#if defined(USE_VOID_PUT_PROC)
        void    (*qi_putp)(queue_t*, mblk_t*);  /* put procedure */
        void    (*qi_srvp)(queue_t*);           /* service procedure */
#else
        int     (*qi_putp)(queue_t*, mblk_t*);  /* put procedure */
        int     (*qi_srvp)(queue_t*);           /* service procedure */
#endif
        int     (*qi_qopen)(queue_t *, dev_t *, int, int, cred_t *);   /* open procedure */
        int     (*qi_qclose)(queue_t *, int,cred_t *);   /* close procedure */
        int     (*qi_qadmin)(void);         /* debugging */
        struct lis_module_info *qi_minfo;   /* module information structure */
        struct module_stat *qi_mstat;   /* module statistics structure */
} qinit_t;
```

Each module should define 2 such qinit structures, one for the read side and one for the write side, pointers to them are kept inside a table called streamtab :

```
typedef struct streamtab {
  SHARE
        struct qinit *st_rdinit; /* read queue */
        struct qinit *st_wrinit; /* write queue */
        struct qinit *st_muxrinit; /* mux read queue */
        struct qinit *st_muxwinit; /* mux write queue */
} streamtab_t;
```

The user defines the read and write qinit structures, and creates the streamtab structure that links the 2 together.

*Figure 16 streamtab structure linking the read and write side of a queue*

A module is pushed on the top of a stream with the I_PUSH ioctl call or calling stropen. The STREAMS then calls qattach with the streamtab as an argument. qattach allocates a pair of queues with allocq, links the queues onto the stream and calls the module's qopen routine.

The standard read and write calls can't differentiate between message types and priorities, they simply read and write unqualified streams of bytes. Therefore the putmsg,getmsg,putpmsg,getpmsg functions are added (the putpmsg and getpmsg functions simply have one additional argument that specifies a priority for the message). The put functions accept a control and a data buffer, and will build a message with a msgb of type M_PROTO or M_DATA or one with a msgb of type M_PROTO and one of type M_DATA in a chain. The get functions will perform the reverse.

```
int putmsg(int fd, void *ctlptr, void *dataptr, int flags)
int getmsg(int fd, void *ctlptr, void *dataptr, int *flagsp)
int putpmsg(int fd, void *ctlptr, void *dataptr, int band, int flags) ;
int getpmsg(int fd, void *ctlptr, void *dataptr, int *bandp, int *flagsp) ;
```

**5.2. Avoiding unecessary copies sharing or cloning network buffers**

Sometimes it is necessary that a protocol layer or a module, a device driver needs to keep a copy of a network buffer. For instance a reliable protocol can need to keep a copy for retransmission (TCP). Or a packet capture driver can need it to pass it to another reader. In many cases who requires the copy doesn't need to alter it. A mechanism of lazy copy or COW or sharing of network buffers therefore could provide great benefits. Linux has two different mechanisms to avoid unnecessary copies of network buffers : one in which everything is shared (the header and the data, this is obtained calling skb_get for instance, and denoted by the skb->users field in the header being different from 1), and the other in which only the data is shared (denoted by incrementing the shared_info field dataref), also called cloning, where two different headers, point to the same data part. It is important to stress that in the Linux cloning the pointers inside the headers can be modified without the necessity to copy the data (adding/stripping headers).

## 6. Linux skbuffs

The skbuff system is a memory management facility explicitly thought for the network code. It is a small software layer that makes use of the general memory allocation facilities offered by the Linux kernel. Starting with version 2.2 the Linux kernel introduced the slab system for allocation of small memory areas and eventually the caching of information between allocations for specific structures. The slab allocator keeps continguous physical memory for each slab.

### 6.1. Linux memory management

There are different levels of memory management on Linux. There is a base level that allocates contiguous pages of memory, this is the physical page allocator, and an upper allocator for smaller memory areas that is organized to take advantage of object reuse.



*Figure 17 skb_mem.pic*

### 6.1.1. Physical page allocation

The physical page allocator in Linux is patterned after the binary buddy system [1] This allocator keeps pages in a power of 2 list, this power of 2 is called order and in the current implementation the maximum of it (MAX_ORDER) is 10 and so there are lists for chunks of physical pages from $2^0$=1(4 KB) to $2^9$=512(2

MB). If a block of the desired size is not available a higher order block is split and the 2 halves become buddies. When eventually running through the list you find a buddy free, you look if the other is free too and eventually you join them and link the resulting higher order block to the higher order list. This system has a coarse granularity and it is not suitable for allocation of small memory areas, also because this kind of allocation can be expensive in terms of the required time. At the very end all memory allocation is anyway based on this method.

### 6.1.2. Small areas allocation

This allocator is essentially based on the slab allocator proposed by Bonwick in 1994 [Bonwick 94], and first implemented on Sun Solaris. It is based on the concept of object based allocation. This system relies on the physical page allocator. Named slabs are created and given a set of contiguous pages for the allocation of specific structures, so that the time required for their complete initialization is avoided. Anyway this same system is used also for the allocation of generic small memory areas, this is done providing generic slabs in power of 2 sizes.

### 6.1.2.1. Named slabs

These are used for the object based allocation. At their creation a constructor and destructor function are specified. For instance tcp_open_request, tcp_bind_bucket, tcp_tw_bucket, sock, skbuff_head_cache are named slabs.
Named slabs are created with the function:

```
kmem_cache_t* kmem_cache_create(const char* name,size_t size,size_t offset,
          unsigned long flags,
          void (*ctor) (void*,kmem_cache_t*,unsigned long),
          void (*dtor) (void*,kmem_cache_t*,unsigned long))
```

Then the allocation of an object from a slab is done with :

```
void* kmem_cache_alloc(kmem_cache_t* cachep, int flags )
```

and its deallocation with :

```
void kmem_cache_free(kmem_cache_t* cachep, void* objp)
```

### 6.1.2.2. Size-n generic slabs

Two power of 2 lists of generic slabs ranging in size from $2^5$=32 bytes to $2^{17}$=128KB of memory chunks are kept. One is suitable for DMA (in x86 architecture this memory should be below 16 MB) and one is general. They are named size-32(DMA), ..., size-131072(DMA) and size-32, ...., size-131072 respectively.
Allocation from these generic slabs happens when the generic kernel allocator

```
void * kmalloc(size_t size,int flags)
```

function is invoked.
Info about slabs can be obtained from /proc/slabinfo where the columns are : slab name, active objects, number of objects, object size.

### 6.1.3. Skbuff use of Linux memory management

The skbuff system uses in two ways the slab allocator.[2] The skbuff heads are allocated from a named slab called skbuff_head_cache, created when the skb_init function is called at network initialization time.

_____ _[net/core/skbuff.c]_

```
1107    void __init skb_init(void)
1108    {
1109          skbuff_head_cache = kmem_cache_create("skbuff_head_cache",
1110                                      sizeof(struct sk_buff),
1111                                      0,
1112                                      SLAB_HWCACHE_ALIGN,
1113                                      NULL, NULL);
1114          if (!skbuff_head_cache)
1115              panic("cannot create skbuff cache");
1116    }
```
——————————————————————————————————————————————————— *[net/core/skbuff.c]*

While the data areas of skbuffs, not being able of taking any advantage from the previous allocations and varying in size are allocated from size-N generic slabs using the kmalloc() function.

——————————————————————————————————————————————————— *[net/core/skbuff.c]*

```
 125    struct sk_buff *alloc_skb(unsigned int size, int gfp_mask)
 126    {
 127          struct sk_buff *skb;
 128          u8 *data;
 129
 130          /* Get the HEAD */
 131          skb = kmem_cache_alloc(skbuff_head_cache,
 132                          gfp_mask & ~__GFP_DMA);
 133          if (!skb)
 134              goto out;
 135
 136          /* Get the DATA. Size must match skb_add_mtu(). */
 137          size = SKB_DATA_ALIGN(size);
 138          data = kmalloc(size + sizeof(struct skb_shared_info), gfp_mask);
 139          if (!data)
 140              goto nodata;
 141
 142          memset(skb, 0, offsetof(struct sk_buff, truesize));
 143          skb->truesize = size + sizeof(struct sk_buff);
 144          atomic_set(&skb->users, 1);
 145          skb->head = data;
 146          skb->data = data;
 147          skb->tail = data;
 148          skb->end  = data + size;
 149
 150          atomic_set(&(skb_shinfo(skb)->dataref), 1);
 151          skb_shinfo(skb)->nr_frags  = 0;
 152          skb_shinfo(skb)->tso_size = 0;
 153          skb_shinfo(skb)->tso_segs = 0;
 154          skb_shinfo(skb)->frag_list = NULL;
 155    out:
 156          return skb;
 157    nodata:
 158          kmem_cache_free(skbuff_head_cache, skb);
 159          skb = NULL;
 160          goto out;
 161    }
```

——————————————————————————————————————————— *[net/core/skbuff.c]*

### 7. Fundamental data structures

*File :* `[include/linux/skbuff.h]`

#### 7.1. sk_buff

The most important data structure is the `sk_buff`. This is the skbuff header where all the status information for a linear skbuff are kept. Every skbuff has an sk_buff structure that holds all the pointers to the data areas. The skbuff header is allocated from the relative memory slab. Its size is 160 bytes. The skbuffs are moved from queues at socks to/from queues at devices.

This is done through the use of the `next` and `prev` pointers that can link skbuffs in a doubly linked list. The head of the list where they are linked to is pointed to by the `list` pointer. This head can be the send/receive queue at the sock/device.
The sock, if any, associated with the skbuff is pointed to by the `sk` pointer.
And the device from where the data arrived or is leaving by is pointed to by the `dev` and `real_dev` pointers. The real_dev is used for example in bonding and VLAN drivers. In bonding, when a packet is received, if the device on which it is received has a master, then the real_dev is set to the dev contents and the dev field is set to the master device    :

——————————————————————————————————————————— *[net/core/dev.c]*

```
1531    static __inline__ void skb_bond(struct sk_buff *skb)
1532    {
1533        struct net_device *dev = skb->dev;
1534
1535        if (dev->master) {
1536            skb->real_dev = skb->dev;
1537            skb->dev = dev->master;
1538        }
1539    }
```

——————————————————————————————————————————— *[net/core/dev.c]*

Pointers to the transport header(tcp/udp) `h`, network layer header (ip) `nh`, link layer header `mac` are filled as soon as known.
A pointer to a destination cache entry is kept in `dst`. Security information (keys and so on) for IPSec are pointed to by the `sp` pointer.
A free area of 48 bytes called control buffer ( `cb` ) is left for specific protocol layers necessities (that area can be used to pass info between protocol layers).
The `len` field keeps the length in bytes of the data area, this is the total data area, encompassing also the eventual pages of data of a fragmented skbuff.
The `data_len` field is the length in bytes of the data area, not in the linear part of the skbuff. If this field is different from zero, then the skbuff is fragmented. The difference `len - data_len` is the amount of data in the linear part of the skbuff, and is also called `headlen` (not to be confused with the size of `headroom`). The `csum` field keeps the eventual checksum of the data. The `local_df` field is used to signal if the real path mtu discovery was requested or not. local_df == 1 means the IP_PMTUDISC_DO was not requested, local_df == 0 means the IP_PMTUDISC_DO was requested and so an icmp error should be generated if we receive fragments. The `cloned` field signals the skbuff has been cloned and so if a user wants to write on it, the skbuff should be copied. The `pkt_type` field describes the destination of the packet (for us, for someone else, broadcast, multicast .. ) according to the following definitions :

——————————————————————————————————————————— *[include/linux/if_packet.h]*

```
24      #define PACKET_HOST       0              /* To us       */
25      #define PACKET_BROADCAST  1              /* To all      */
26      #define PACKET_MULTICAST  2              /* To group       */
27      #define PACKET_OTHERHOST  3              /* To someone else  */
28      #define PACKET_OUTGOING       4              /* Outgoing of any type */
29      /* These ones are invisible by user level */
30      #define PACKET_LOOPBACK       5              /* MC/BRD frame looped back */
31      #define PACKET_FASTROUTE  6              /* Fastrouted frame  */
```
———————————————————————————————————————————————— *[include/linux/if_packet.h]*

The `ip_summed` field tells if the driver supplied an ip checksum. It can be NONE, HW or UNNECES-
SARY :

———————————————————————————————————————————————— *[include/linux/skbuff.h]*
```
35      #define CHECKSUM_NONE 0
36      #define CHECKSUM_HW 1
37      #define CHECKSUM_UNNECESSARY 2
```
———————————————————————————————————————————————— *[include/linux/skbuff.h]*

On input, CHECKSUM_NONE means the device failed to checksum the packet and so csum is undefined,
CHECKSUM_UNNECESSARY means that the checksum has already been verified, but the problem is that it
is not known in which way (for example as an ipv6 or an ipv4 packet ..), so it is an unrecommended flag.
CHECKSUM_HW means the device provides the checksum in the csum field. On output, CHECKSUM_NONE
means checksum provided by protocol or not required, CHECKSUM_HW means the device is required to
checksum the packet (from the header h.raw to the end of the data and put the checksum in the csum field).
The `priority` field keeps the priority level according to :

———————————————————————————————————————————————— *[include/linux/pkt_sched.h]*
```
1      #ifndef __LINUX_PKT_SCHED_H
2      #define __LINUX_PKT_SCHED_H
3
4      /* Logical priority bands not depending on specific packet scheduler.
5         Every scheduler will map them to real traffic classes, if it has
6         no more precise mechanism to classify packets.
7
8         These numbers have no special meaning, though their coincidence
9         with obsolete IPv6 values is not occasional :-). New IPv6 drafts
10        preferred full anarchy inspired by diffserv group.
11
12        Note: TC_PRIO_BESTEFFORT does not mean that it is the most unhappy
13        class, actually, as rule it will be handled with more care than
14        filler or even bulk.
15     */
16
17     #define TC_PRIO_BESTEFFORT     0
18     #define TC_PRIO_FILLER             1
19     #define TC_PRIO_BULK           2
20     #define TC_PRIO_INTERACTIVE_BULK   4
21     #define TC_PRIO_INTERACTIVE        6
22     #define TC_PRIO_CONTROL            7
23
24     #define TC_PRIO_MAX            15
25
```

———————————————————————————————————————————————— *[include/linux/pkt_sched.h]*

they are used by traffic control mechanisms.

The `security` field keeps the level of security.
The `truesize` field keeps the real size occupied by the skbuff, that is it adds the size of the header to the size of the data when the skbuff is allocate in `alloc_skb()`:

———————————————————————————————————————————————— *[net/core/skbuff.c]*

```
140                     goto nodata;
141
142             memset(skb, 0, offsetof(struct sk_buff, truesize));
143             skb->truesize = size + sizeof(struct sk_buff);
```
———————————————————————————————————————————————— *[net/core/skbuff.c]*

When a copy is made, the skbuff header is copied up to the truesize field, because the remaining fields are pointers to the data areas and so need to be replaced.
The `head`, end pointers, are pointers to the boundaries of the available space.
The `data`, tail pointers are pointers to the beginning and end of the already used data area.

———————————————————————————————————————————————— *[include/linux/skbuff.h]*

```
185      *   @nf_bridge: Saved data about a bridged frame - see br_netfilter.c
186      *      @private: Data which is private to the HIPPI implementation
187      *   @tc_index: Traffic control index
188      */
189
190     struct sk_buff {
191             /* These two members must be first. */
192             struct sk_buff        *next;
193             struct sk_buff        *prev;
194
195             struct sk_buff_head   *list;
196             struct sock           *sk;
197             struct timeval        stamp;
198             struct net_device     *dev;
199             struct net_device     *real_dev;
200
201             union {
202                     struct tcphdr   *th;
203                     struct udphdr   *uh;
204                     struct icmphdr  *icmph;
205                     struct igmphdr  *igmph;
206                     struct iphdr    *ipiph;
207                     unsigned char   *raw;
208             } h;
209
210             union {
211                     struct iphdr    *iph;
212                     struct ipv6hdr  *ipv6h;
213                     struct arphdr   *arph;
214                     unsigned char   *raw;
215             } nh;
```

```
216
217         union {
218               struct ethhdr   *ethernet;
219               unsigned char   *raw;
220         } mac;
221
222         struct  dst_entry    *dst;
223         struct    sec_path *sp;
224
225         /*
226          * This is the control buffer. It is free to use for every
227          * layer. Please put your private variables there. If you
228          * want to keep them across layers you have to do a skb_clone()
229          * first. This is owned by whoever has the skb queued ATM.
230          */
231         char           cb[48];
232
233         unsigned int        len,
234                        data_len,
235                        csum;
236         unsigned char       local_df,
237                        cloned,
238                        pkt_type,
239                        ip_summed;
240         __u32          priority;
241         unsigned short      protocol,
242                        security;
243
244         void           (*destructor)(struct sk_buff *skb);
245    #ifdef CONFIG_NETFILTER
246         unsigned long        nfmark;
247         __u32          nfcache;
248         struct nf_ct_info    *nfct;
249    #ifdef CONFIG_NETFILTER_DEBUG
250         unsigned int      nf_debug;
251    #endif
252    #ifdef CONFIG_BRIDGE_NETFILTER
253         struct nf_bridge_info     *nf_bridge;
254    #endif
255    #endif /* CONFIG_NETFILTER */
256    #if defined(CONFIG_HIPPI)
257         union {
258               __u32     ifield;
259         } private;
260    #endif
261    #ifdef CONFIG_NET_SCHED
262         __u32           tc_index;                /* traffic control index */
263    #endif
264
265         /* These elements must be at the end, see alloc_skb() for details.  */
266         unsigned int        truesize;
267         atomic_t        users;
268         unsigned char       *head,
```

*[include/linux/skbuff.h]*

## 7.2. skb_shared_info

The `skb_shared_info` structure is used by the fragmented skbuffs. It has a meaning when the `data_len` field in the skbuff header is different from zero. This field counts the data not in the linear part of the skbuff.

The `dataref` field counts the number of references to the fragmented part of the skbuff, so that a writer knows if it is necessary to copy it.

The `nr_frags` field keeps the number of pages in which this skbuff is fragmented. This kind of fragmentation is done for interfaces supporting scatter and gather. This feature is described in the netdevice structure by the NETIF_F_SG flag. (3com 3c59x , 3com typhoon, Intel e100, ... ) When an skbuff is to be allocated, if the mss is larger than a page then if the interface supports scatter and gather a linear skbuff of a single page is allocated with alloc_skb and then the other pages are allocated and added to the frags array.

The `tso_size,tso_segs` fields were added to support cards able to perform by themselves the tcp segmentation (they are described by the NETIF_F_TSO TCP Segmentation Offload). The tso_size comes from the mss, and is the max size that should be used by the card for segments. (3Com Typhoon family 3c990, 3cr990 supports it if the array of pages is <= 32)

The `frag_list` pointer is used when the skbuff is fragmented in a list. The frag_list pointer is only used to connect the 1st fragment to the second, the other skbuff are linked through the standard next pointer. Skbuffs fragmented in a list are used for the reassembly of ip fragments. This is eventually done when the interface supports the NETIF_F_FRAG_LIST feature. There are no devices in the standard linux kernel tree that support this feature at the moment (except the trivial loopback).

The `frags` array keeps the pointers to the page structures in which the skbuff has been fragmented. The last used page pointer is `nr_frags-1` and there is space for up to MAX_SKB_FRAGS. This was only 6 in previous versions, now it is sufficient to accomodate a maximum length tcp segment (64 KB).

*[include/linux/skbuff.h]*

```
124
125     struct sk_buff;
```

*[include/linux/skbuff.h]*

*[include/linux/skbuff.h]*

```
138     /* This data is invariant across clones and lives at
139      * the end of the header data, ie. at skb->end.
140      */
141     struct skb_shared_info {
142         atomic_t   dataref;
143         unsigned int   nr_frags;
144         unsigned short  tso_size;
145         unsigned short  tso_segs;
```

*[include/linux/skbuff.h]*

The `shared_info` structure that we said is placed at the `skb->end` is usually accessed using the macro `skb_shinfo`

*[include/linux/skbuff.h]*

```
306     #define skb_shinfo(SKB)        ((struct skb_shared_info *)((SKB)->end))
```

*[include/linux/skbuff.h]*

for instance :
`skb_shinfo(skb)->frag_list = 0;`

```
skb_shinfo(skb)->dataref
```

## 8.  Skbuff organization

The network buffers in Linux are formed by 2 different entities : a fixed size header, and a variable size data area.  Until recently the data area of the skbuff was unique and physically contiguous (aka linear). And the Linux kernel was cited because of the efficiency it could obtain with dumb interfaces against other popular OSs like BSD, in which frequently because of the small size of the network buffers, you could have a list of them for a single network packet.  In version 2.4 of the stable branch of the kernel fragmented skbuffs were introduced in two forms and for two completely different reasons. To keep the information necessary to describe this more complicate skbuff organizations a small area placed immediately after the skbuff data area (skb->end) was allocated and reserved.  Placing this info there could allow the sharing of the information between multiple clones and so the structure that keeps the information is called `skb_shared_info`.  The possibility that an skbuff can have data in an array of associated unmapped pages was added to allow the efficient use of interfaces able to perform gather and scatter in hardware.  To efficiently manage ip fragments, a chain of skbuffs that could be passed along the network stack like a single skbuff was introduced.

### 8.1.  Linux sharing and cloning

Linux has the most flexible way of providing  mechanisms to avoid unnecessary copies.  With *sharing* the complete skbuff (headers and data area) is shared between two entities, with *cloning* two different headers point to the same data area.
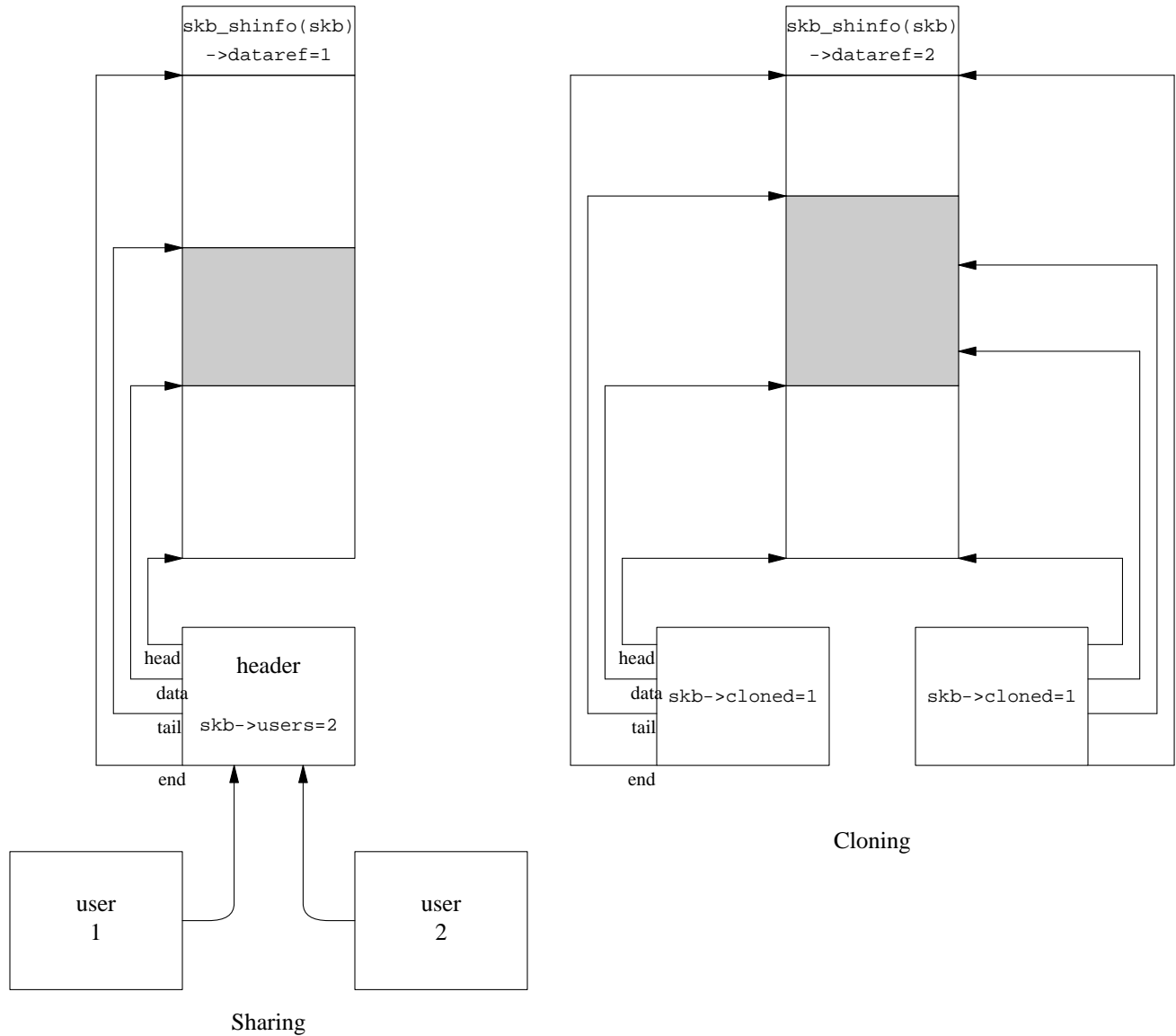
skb_shinfo(skb)
->dataref=1

header

skb->users=2

head
data
tail
end

user
1

user
2

Sharing

skb_shinfo(skb)
->dataref=2

skb->cloned=1

skb->cloned=1

head
data
tail
end

Cloning

*Figure 18 Sharing and cloning skbuffs*

Sharing is the least expensive, but doesnt allow any change. Cloning requires copying the header, but then allows the flexibility to modify the header, without copying the data (for instance the pointers skb->data and skb->tail can be changed to strip away headers). Sharing is denoted by atomically incrementing the skb->users counter for each user holding a ref to the skbuff, cloning is denoted by setting the skb->cloned flag in the header and incrementing the skb_shinfo(skb)->dataref field in the shared_info struct.

An skbuff can be shared calling the skb_get function over an existing skbuff. This will atomically increment the skb->users counter and return a pointer to the skbuff. For instance :

### 8.2. Linear skbuffs

This was the original and only way in which skbuffs were used. In this case the variable length data area was formed by a physically contiguous area of memory. Network operations that resulted in adding/stripping headers or trailers or data, were performed changing the pointers to the data area ( `skb->data`,`skb->tail` ) inside the header. If for any reason the data area could not be expanded in place to add headers, trailers or data, a new linear skbuff is created and the data is copied over. In the case of a linear skbuff, the only info in the skb_shared_info area that could be used is that related to the tso (TCP Segmentation Offloading) and the counter that indicates how many entities share the data area of the skbuff (clones).



*Figure 19 skbuff_three.pic*

### 8.3. Nonlinear skbuffs

For non linear skbuffs the infor

### 8.3.1. Paged skbuffs

In this case some unmapped pages of memory are associated with the skbuff. The number of these pages is stored in the `nr_frags` field of the skb_shared_info. And an array of page pointers, of size to allow at least 64 KB of data (it was just 6 pages in kernel 2.4) is used. This kind of fragmentation was introduced to make use of the scatter/gather capability of modern interfaces which can be given a vector of buffers to deal with by themselves. An important difference with BSD should be stressed : the Linux code allows in principle to have an skbuff using the linear part, and being fragmented in both the possible ways at the same time. On BSD if an `mbuf` uses an external page, then no data is put in its linear part. A hook to allocate paged skbuff is provided by the function `sock_alloc_send_pskb`. This function would allocated a

paged skbuff having at least data_len bytes of space in its external pages.

————————————————————————————————————————————————————— *[net/core/sock.c]*

```
773            finish_wait(sk->sk_sleep, &wait);
774            return timeo;
775     }
776
777
778     /*
779      *   Generic send/receive buffer handlers
780      */
781
782     struct sk_buff *sock_alloc_send_pskb(struct sock *sk, unsigned long header_len,
783                              unsigned long data_len, int noblock, int *errcode)
784     {
785            struct sk_buff *skb;
786            unsigned int gfp_mask;
787            long timeo;
788            int err;
789
790            gfp_mask = sk->sk_allocation;
791            if (gfp_mask & __GFP_WAIT)
792                 gfp_mask |= __GFP_REPEAT;
793
794            timeo = sock_sndtimeo(sk, noblock);
795            while (1) {
796                 err = sock_error(sk);
797                 if (err != 0)
798                      goto failure;
799
800                 err = -EPIPE;
801                 if (sk->sk_shutdown & SEND_SHUTDOWN)
802                      goto failure;
803
804                 if (atomic_read(&sk->sk_wmem_alloc) < sk->sk_sndbuf) {
805                      skb = alloc_skb(header_len, sk->sk_allocation);
806                      if (skb) {
807                           int npages;
808                           int i;
809
810                           /* No pages, we're done... */
811                           if (!data_len)
812                                break;
813
814                           npages = (data_len + (PAGE_SIZE - 1)) >> PAGE_SHIFT;
815                           skb->truesize += data_len;
816                           skb_shinfo(skb)->nr_frags = npages;
817                           for (i = 0; i < npages; i++) {
818                                struct page *page;
819                                skb_frag_t *frag;
820
821                                page = alloc_pages(sk->sk_allocation, 0);
822                                if (!page) {
823                                     err = -ENOBUFS;
```

```
824                                    skb_shinfo(skb)->nr_frags = i;
825                                    kfree_skb(skb);
826                                    goto failure;
827                               }
828
829                               frag = &skb_shinfo(skb)->frags[i];
830                               frag->page = page;
831                               frag->page_offset = 0;
832                               frag->size = (data_len >= PAGE_SIZE ?
833                                          PAGE_SIZE :
834                                          data_len);
835                               data_len -= PAGE_SIZE;
836                          }
837
838                          /* Full success... */
839                          break;
840                     }
841               err = -ENOBUFS;
842               goto failure;
843          }
844          set_bit(SOCK_ASYNC_NOSPACE, &sk->sk_socket->flags);
845          set_bit(SOCK_NOSPACE, &sk->sk_socket->flags);
846          err = -EAGAIN;
847          if (!timeo)
848               goto failure;
849          if (signal_pending(current))
850               goto interrupted;
851          timeo = sock_wait_for_wmem(sk, timeo);
852     }
853
```
——————————————————————————————————————————————— *[net/core/sock.c]*

But in fact this function is never called directly by the current kernel but only through the
`sock_alloc_send_skb` function that doesnt request any space in external pages.

——————————————————————————————————————————————— *[net/core/sock.c]*
```
855          return skb;
856
857     interrupted:
858          err = sock_intr_errno(timeo);
859     failure:
```
——————————————————————————————————————————————— *[net/core/sock.c]*

In the inet_stream_ops table there is a .sendpage method, this is initialized to tcp_sendpage when the table
is defined in af_inet.c. Tcp_sendpage if the scatter/gather support is not listed in the route capabilities, will
call sock_no_sendpage, otherwise it calls do_tcp_sendpages. This seems to be the only place that gives rise
to a paged skbuff. [tcp_alloc_pskb and tcp_alloc_skb only differ because the 1st can add to the truesize of
the allocated skbuff the size of the memory that will be allocated in external pages, but it doesnt allocate
them at all. The second arg of tcp_alloc_pskb is size=linear part of skb, while the 3d arg is mem=paged
part of skb] This paged skbuff is requested to have a paged part of at least tp->mss_cache bytes and a linear
part of 0 bytes (the place for TCP headers is automatically added. So the resulting pskb will have the head-
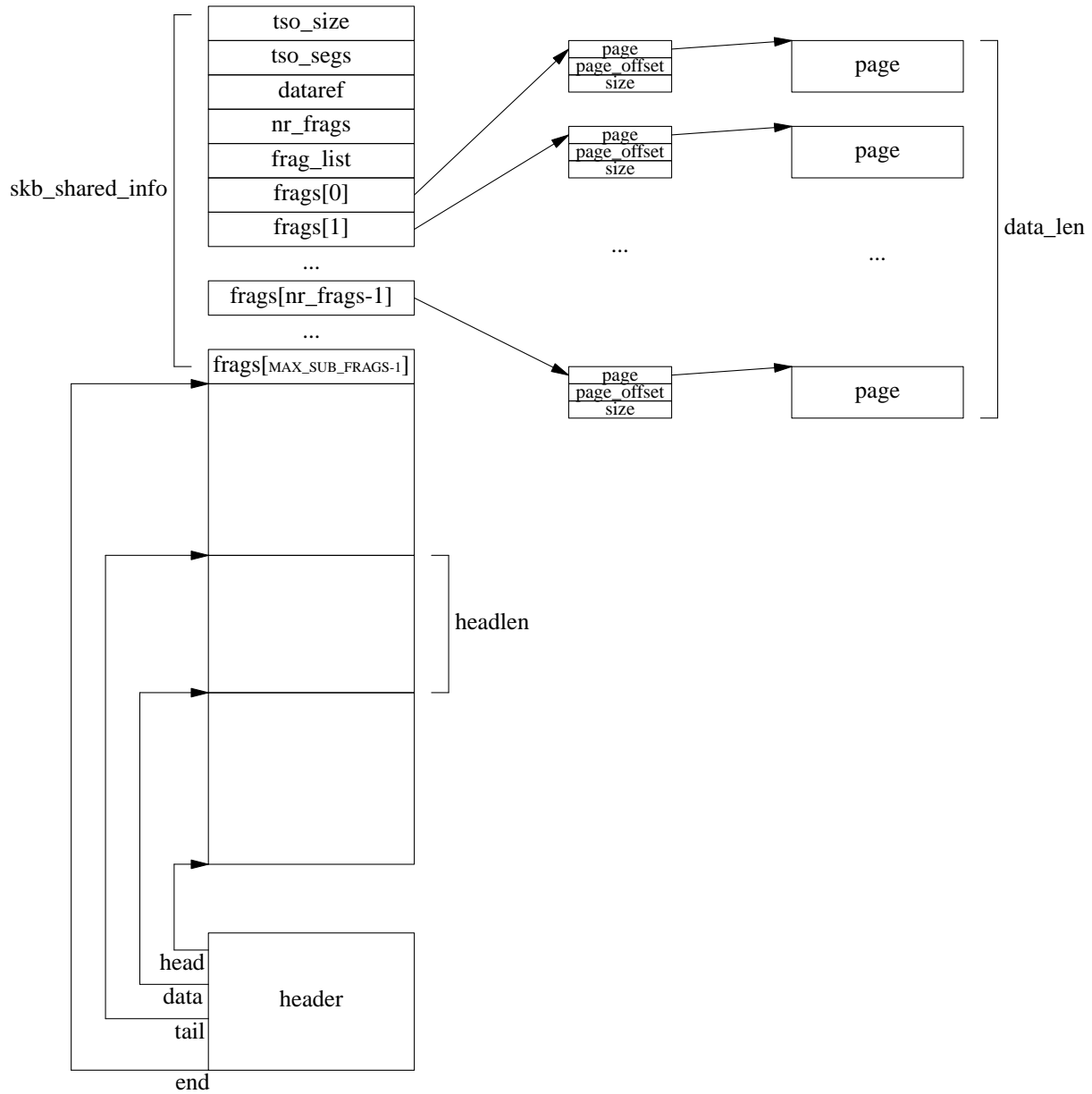ers in the linear part of the skb and the data in the external pages.

*Figure 20 Paged skbuff (use of nr_frags and frags[])*

### 8.3.2. Skbuff chains

Skbuff chains can be treated like a single skbuff and passed around the network stack. This organization is revealed by the `frag_list` pointer of the first skbuff. This pointer points to the second skbuff. The following skbuffs are simply linked through their next pointer. In this way to transform a list of skbuffs, in an skbuff chain it is only necessary to copy the next pointer of the first skbuff over the frag_list pointer. This is what is done inside the ip_local_deliver->ip_defrag->ip_frag_reasm function which transforms a queue of IP fragments stored in an IP fragment queue into an skbuff chain that can be passed around the network stack unitarily. Fragments that are to be forwarded are instead dispatched long before and without any delay in the input process from the ip_rcv_finish function to the ip_forward function.

———————————————————————————————— *[net/ipv4/ip_fragment.c]*

```
590              skb_shinfo(head)->frag_list = head->next;
```

———————————————————————————————— *[net/ipv4/ip_fragment.c]*

In the recent addition (kernel version 2.6 ) of the SCTP protocol this kind of fragmented skbuff arises when reassembling the messages of an SCTP stream.

———————————————————————————————— *[net/sctp/ulpqueue.c]*

```
314                  skb_shinfo(f_frag)->frag_list = pos;
```

———————————————————————————————— *[net/sctp/ulpqueue.c]*

(SCTP is a transport protocol offering reliable, multiple stream, multihoming transport service that is TCP friendly and without head-of-line blocking RFC 3286)



*Figure 21 Skbuff chain (use of frag_list pointer)*

## 9. Queues of skbuffs

During output skbuffs are queued first on the socket, and then when the output interface is determined the skbuffs are moved to the device queue. In input the skbuffs are queued on the device queue and then when the owner socket is found are moved to the owner socket queue. These queues are doubly linked lists of skbuffs, formed using the next and prev pointers in the skbuff header. Also IP fragments are stored in such

a kind of queues, just that fragments queues are kept in offset order and so insertions can happen also in the middle of the queue.
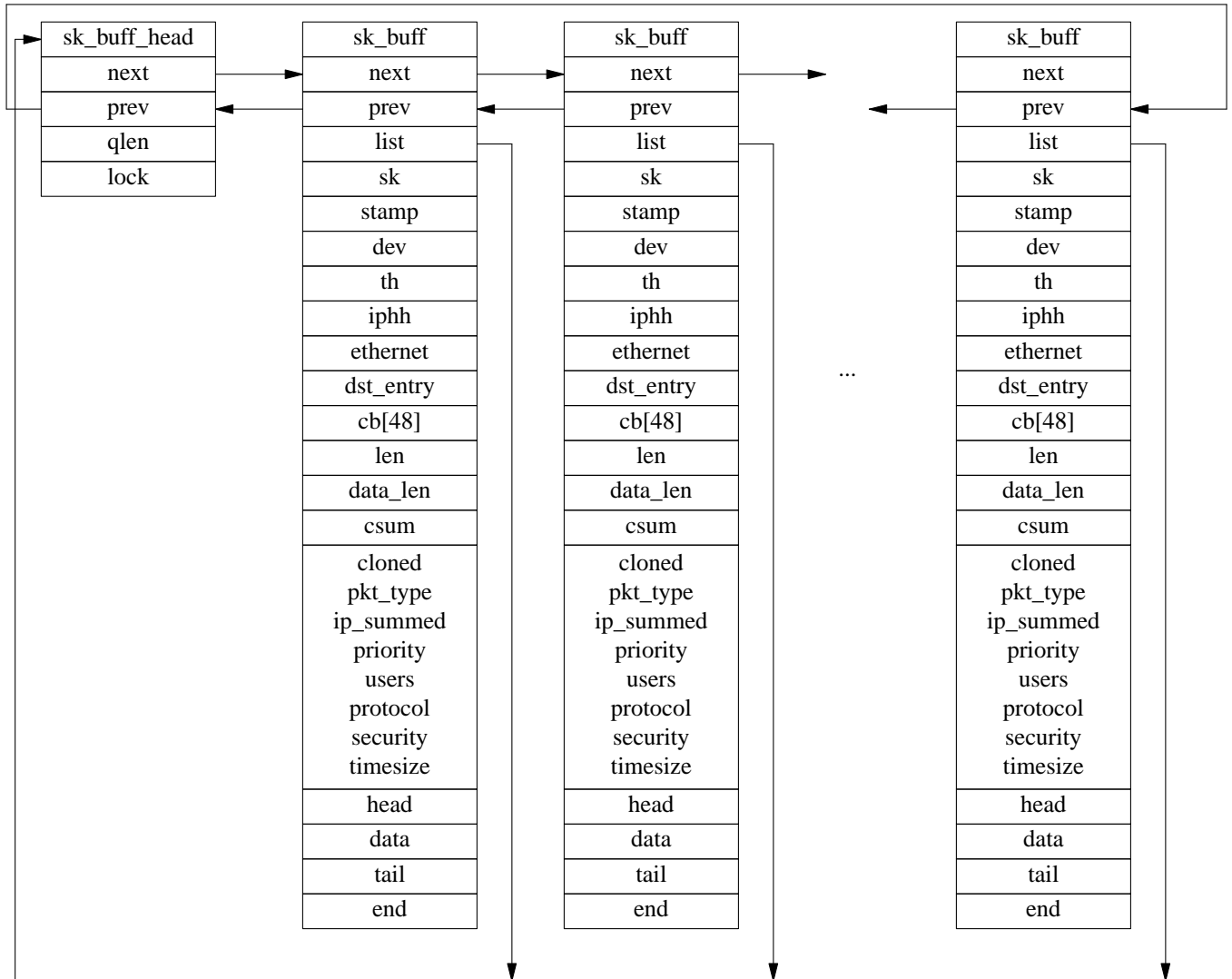


*Figure 22 skb_buff.pic*

## 9.1.  Functions to manage skbuff queues

### 9.1.1.  skb_queue_len

This function returns the number of skbuffs queued on the list that you pass as an argument. The socket write and read queues for instance ( sk->sk_write_queue , sk->sk_receive_queue ) or the device queues dev->qdisc->q or the TCP queue of out of order segments tp->out_of_order_queue or the tp->ucopy.prequeue used to hold skbuffs to pass data to the user in big chunks for efficiency.

*File :* [net/sctp/ulpqueue.c]

——————————————————————————————————————————— *[include/linux/skbuff.h]*

```
485     static inline __u32 skb_queue_len(const struct sk_buff_head *list_)
486     {
487            return list_->qlen;
488     }
```

——————————————————————————————————————————— *[include/linux/skbuff.h]*


### 9.1.2.  skb_queue_head_init

This function initializes a queue of skbuffs : it initializes the spinlock structure inside the sk_buff_head and
sets the `prev` and `next` pointers to the list head and sets the qlen field to zero.

*File :* `[include/linux/skbuff.h]`

——————————————————————————————————————————— *[include/linux/skbuff.h]*

```
490     static inline void skb_queue_head_init(struct sk_buff_head *list)
491     {
492            spin_lock_init(&list->lock);
493            list->prev = list->next = (struct sk_buff *)list;
494            list->qlen = 0;
495     }
```

——————————————————————————————————————————— *[include/linux/skbuff.h]*


### 9.1.3.  skb_queue_head and __skb_queue_head

These two functions queue an skbuff  buffer at the start of a list.  The `skb_queue_head` function estab-
lishes a spin lock on the queue, using its spinlock structure and so it is safe to call it  with interrupts enable,
it then calls the `__skb_queue_head` function to queue the buffer.  The `__skb_queue_head` function
can be called by itself only with interrupts disabled.

*File :* `[include/linux/skbuff.h]`

——————————————————————————————————————————— *[include/linux/skbuff.h]*

```
507     *    @newsk: buffer to queue
508     *
509     *    Queue a buffer at the start of a list. This function takes no locks
510     *    and you must therefore hold required locks before calling it.
511     *
512     *    A buffer cannot be placed on two lists at the same time.
513     */
514     static inline void __skb_queue_head(struct sk_buff_head *list,
515                             struct sk_buff *newsk)
516     {
517            struct sk_buff *prev, *next;
518
519            newsk->list = list;
```

——————————————————————————————————————————— *[include/linux/skbuff.h]*


——————————————————————————————————————————— *[include/linux/skbuff.h]*

```
533     *
```

```
534      *    Queue a buffer at the start of the list. This function takes the
535      *    list lock and can be used safely with other locking &sk_buff functions
536      *    safely.
537      *
538      *    A buffer cannot be placed on two lists at the same time.
539      */
540      static inline void skb_queue_head(struct sk_buff_head *list,
541                              struct sk_buff *newsk)
```
——————————————————————————————————————— *[include/linux/skbuff.h]*


### 9.1.4. skb_queue_tail and __skb_queue_tail

These two functions queue an skbuff at the tail  of an skbuff queue. The skb_queue_tail function
establishes a spin lock on the queue and so it is safe to call it  with interrupts enabled, it then calls the
__skb_queue_tail function to queue the buffer. The __skb_queue_tail function can be called
by itself only with interrupts disabled.


*File :* [include/linux/skbuff.h]

——————————————————————————————————————— *[include/linux/skbuff.h]*
```
585    static inline void skb_queue_tail(struct sk_buff_head *list,
586                            struct sk_buff *newsk)
587    {
588        unsigned long flags;
589
590        spin_lock_irqsave(&list->lock, flags);
591        __skb_queue_tail(list, newsk);
592        spin_unlock_irqrestore(&list->lock, flags);
593    }
```
——————————————————————————————————————— *[include/linux/skbuff.h]*


——————————————————————————————————————— *[include/linux/skbuff.h]*
```
560    static inline void __skb_queue_tail(struct sk_buff_head *list,
561                            struct sk_buff *newsk)
562    {
563        struct sk_buff *prev, *next;
564
565        newsk->list = list;
566        list->qlen++;
567        next = (struct sk_buff *)list;
568        prev = next->prev;
569        newsk->next = next;
570        newsk->prev = prev;
571        next->prev  = prev->next = newsk;
572    }
```
——————————————————————————————————————— *[include/linux/skbuff.h]*

### 9.1.5.  skb_entail

*File :* [include/linux/skbuff.h]
This function is part of the TCP code because it is used only by TCP.  It enqueues an skbuff at the tail of an
skbuff queue using the __skb_queue_tail function, but it also does some housekeeping for the TCP
protocol.  It updates the per socket counters of bytes queued and forwarded.  It stores for example some tcp
status information on the general 48-bytes control buffer in the skbuff.  This information stored on the con-
trol buffer comprises TCP start and end sequence, ack flag, sacked.

```
———————————————————————————————————————— [net/ipv4/tcp.c]
791     static inline void skb_entail(struct sock *sk, struct tcp_opt *tp,
792                            struct sk_buff *skb)
793     {
794         skb->csum = 0;
795         TCP_SKB_CB(skb)->seq = tp->write_seq;
796         TCP_SKB_CB(skb)->end_seq = tp->write_seq;
797         TCP_SKB_CB(skb)->flags = TCPCB_FLAG_ACK;
798         TCP_SKB_CB(skb)->sacked = 0;
799         __skb_queue_tail(&sk->sk_write_queue, skb);
800         tcp_charge_skb(sk, skb);
801         if (!tp->send_head)
802             tp->send_head = skb;
803         else if (tp->nonagle&TCP_NAGLE_PUSH)
804             tp->nonagle &= ~TCP_NAGLE_PUSH;
805     }
———————————————————————————————————————— [net/ipv4/tcp.c]
```

### 9.1.6.  skb_dequeue and __skb_dequeue

These two functions dequeue an skbuff  buffer from the head  of an skbuff queue.  The skb_dequeue
function establishes a spin lock on the queue and therefore it is safe to call it  with interrupts enable, it then
calls the __skb_dequeue function to dequeue the buffer.  The __skb_dequeue function can be called
by itself only with interrupts disabled.

*File :* [net/ipv4/tcp.c]

```
———————————————————————————————————————— [include/linux/skbuff.h]
631     static inline struct sk_buff *skb_dequeue(struct sk_buff_head *list)
632     {
633         unsigned long flags;
634         struct sk_buff *result;
635
636         spin_lock_irqsave(&list->lock, flags);
637         result = __skb_dequeue(list);
638         spin_unlock_irqrestore(&list->lock, flags);
639         return result;
640     }
———————————————————————————————————————— [include/linux/skbuff.h]
```

```
———————————————————————————————————————— [include/linux/skbuff.h]
603     static inline struct sk_buff *__skb_dequeue(struct sk_buff_head *list)
```

```
604     {
605         struct sk_buff *next, *prev, *result;
606
607         prev = (struct sk_buff *) list;
608         next = prev->next;
609         result = NULL;
610         if (next != prev) {
611             result       = next;
612             next      = next->next;
613             list->qlen--;
614             next->prev  = prev;
615             prev->next  = next;
616             result->next = result->prev = NULL;
617             result->list = NULL;
618         }
619         return result;
620     }
```
                                                                    ——————— *[include/linux/skbuff.h]*

### 9.1.7.  skb_insert and __skb_insert

These functions insert an skbuff in a queue of skbuffs before a given skbuff.  The __skb_insert function fixes the pointers to do the queue operations, while the skb_insert function wraps it with code to set and unset a spinlock on the list.

————————————————————————————————————————————— *[include/linux/skbuff.h]*
```
646     static inline void __skb_insert(struct sk_buff *newsk,
647                     struct sk_buff *prev, struct sk_buff *next,
648                     struct sk_buff_head *list)
649     {
650         newsk->next = next;
651         newsk->prev = prev;
652         next->prev  = prev->next = newsk;
653         newsk->list = list;
654         list->qlen++;
655     }
```
————————————————————————————————————————————— *[include/linux/skbuff.h]*

### 9.1.8.  skb_append and __skb_append

These function use __skb_insert to append an skbuff on a queue of skbuffs after a given skbuff.

————————————————————————————————————————————— *[include/linux/skbuff.h]*
```
680     static inline void __skb_append(struct sk_buff *old, struct sk_buff *newsk)
681     {
682         __skb_insert(newsk, old, old->next, old->list);
683     }
```
————————————————————————————————————————————— *[include/linux/skbuff.h]*

### 9.1.9. skb_unlink and __skb_unlink

These function remove an skbuff from a list. As usual the `skb_unlink` function wraps the __skb_unlink function with code to set and unset a spinlock on the list.

——————————————————————————————————————— *[include/linux/skbuff.h]*

```
709     static inline void __skb_unlink(struct sk_buff *skb, struct sk_buff_head *list)
710     {
711             struct sk_buff *next, *prev;
712
713             list->qlen--;
714             next    = skb->next;
715             prev    = skb->prev;
716             skb->next   = skb->prev = NULL;
717             skb->list   = NULL;
718             next->prev = prev;
719             prev->next = next;
720     }
```

——————————————————————————————————————— *[include/linux/skbuff.h]*

### 9.1.10. skb_dequeue_tail and __skb_dequeue_tail

These function dequeue the last skbuff on a list.

——————————————————————————————————————— *[include/linux/skbuff.h]*

```
758     static inline struct sk_buff *__skb_dequeue_tail(struct sk_buff_head *list)
759     {
760             struct sk_buff *skb = skb_peek_tail(list);
761             if (skb)
762                     __skb_unlink(skb, list);
763             return skb;
764     }
```

——————————————————————————————————————— *[include/linux/skbuff.h]*

## 10. Skbuff Functions

Many functions have similar names. Two  often used conventions on their names are that `funct` is the standard function, and the __`funct` is the streamlined function `funct` without some consistency checks, while `pfunct` is the function to be applied if the skbuff is nonlinear (from a form of fragmented skbuff : the paged skbuff).

The following functions distinguish between the three kind of skbuffs : linear, paged skbuff, skbuff chain.

## 10.1. Support functions

### 10.1.1. SKB_LINEAR_ASSERT and skb_is_nonlinear

*File :* `[include/linux/skbuff.h]`
the SKB_LINEAR_ASSERT macro will raise a bug if the skb is nonlinear, this condition is checked looking at the data_len field that reports the size of the data in the nonlinear part of the skbuff.

——————————————————————————————————————— *[include/linux/skbuff.h]*

```
815     #define SKB_LINEAR_ASSERT(skb)  BUG_ON(skb_is_nonlinear(skb))
```
———————————————————————————————————————————————— *[include/linux/skbuff.h]*


———————————————————————————————————————————————— *[include/linux/skbuff.h]*
```
785     static inline int skb_is_nonlinear(const struct sk_buff *skb)
786     {
787           return skb->data_len;
788     }
```
———————————————————————————————————————————————— *[include/linux/skbuff.h]*


### 10.1.2.  SKB_PAGE_ASSERT

*File :* [include/linux/skbuff.h]
This macro will raise a bug if the skbuff is a paged skbuff. We have already discussed that if the skbuff is
fragmented in pages then the number of pages used is kept in the skb_shared_info structure, nr_frags vari-
able.

———————————————————————————————————————————————— *[include/linux/skbuff.h]*
```
813     #define SKB_PAGE_ASSERT(skb)    BUG_ON(skb_shinfo(skb)->nr_frags)
```
———————————————————————————————————————————————— *[include/linux/skbuff.h]*


### 10.1.3.  SKB_FRAG_ASSERT

*File :* [include/linux/skbuff.h]
This macro will raise a bug if the skbuff is an skbuff chain.  We have already discussed that if the skbuff is
fragmented in a list of skbuffs then the pointer to the next skbuff is kept in the skb_shared_info structure,
frag_list variable.

———————————————————————————————————————————————— *[include/linux/skbuff.h]*
```
814     #define SKB_FRAG_ASSERT(skb)    BUG_ON(skb_shinfo(skb)->frag_list)
```
———————————————————————————————————————————————— *[include/linux/skbuff.h]*


### 10.1.4.  skb_headlen and skb_pagelen

*File :* [include/linux/skbuff.h]
The skb_headlen function returns the size of the data occupied in the linear part of the skbuff. This is the
total data stored in the skbuff len, minus the data stored in the nonlinear part of the skbuff : data_len.

———————————————————————————————————————————————— *[include/linux/skbuff.h]*
```
790     static inline unsigned int skb_headlen(const struct sk_buff *skb)
791     {
792           return skb->len - skb->data_len;
793     }
```
———————————————————————————————————————————————— *[include/linux/skbuff.h]*

The skb_pagelen function returns the size of the data stored in the array of pages in which the skbuff is
fragmented.

———————————————————————————————————————————————— *[include/linux/skbuff.h]*

```
795     static inline int skb_pagelen(const struct sk_buff *skb)
796     {
797           int i, len = 0;
798
799           for (i = (int)skb_shinfo(skb)->nr_frags - 1; i >= 0; i--)
800                 len += skb_shinfo(skb)->frags[i].size;
801           return len + skb_headlen(skb);
802     }
```
——————————————————————————————————— *[include/linux/skbuff.h]*


### 10.1.5. skb_fill_page_desc

*File :* `[include/linux/skbuff.h]`
This function add a new page to a paged skbuff, filling the relevant information in the shared_info structure
of the skbuff.

——————————————————————————————————— *[include/linux/skbuff.h]*
```
  804     static inline void skb_fill_page_desc(struct sk_buff *skb, int i, struct page
*page, int off, int size)
  805     {
  806           skb_frag_t *frag = &skb_shinfo(skb)->frags[i];
  807           frag->page = page;
  808           frag->page_offset = off;
  809           frag->size = size;
  810           skb_shinfo(skb)->nr_frags = i+1;
  811     }
```
——————————————————————————————————— *[include/linux/skbuff.h]*


### 10.1.6. pskb_may_pull

*File :* `[include/linux/skbuff.h]`
This function is used for instance with fragmented skbuffs resulting from IP fragments reassembly. It
checks if there are at least len bytes in the skbuff, otherwise it returns false (0). Then if there are less than
len bytes in the first skbuff of the list, the skbuff is linearized calling __pskb_pull_tail, in such a way that
the required bytes are in the first skbuff, and the pull operation for stripping a header can be performed sim-
ply changing the data pointer.

——————————————————————————————————— *[include/linux/skbuff.h]*
```
  912     static inline int pskb_may_pull(struct sk_buff *skb, unsigned int len)
  913     {
  914           if (likely(len <= skb_headlen(skb)))
  915                 return 1;
  916           if (unlikely(len > skb->len))
  917                 return 0;
  918           return __pskb_pull_tail(skb, len-skb_headlen(skb)) != NULL;
  919     }
```
——————————————————————————————————— *[include/linux/skbuff.h]*

## 10.2.  Initialization

### 10.2.1.  skb_init
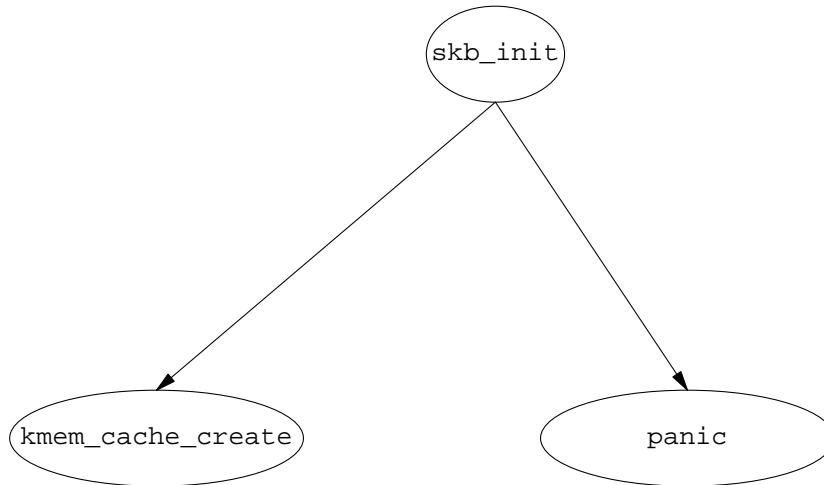
*File :* `[include/net/socket.h]`



*Figure 23 skb_init.pic*

This function initializes the skbuff system. It is called by the `sock_init()` function in the `[net/socket.c]` at socket initialization time. It creates a slab named `skbuff_head_cache` for skbuff head objects. It doesnt specify any specific constructor or destructor for them.

———————————————————————————————————————————— *[include/net/socket.h]*

```
1096      *    Tune the memory allocator for a new MTU size.
1097      */
1098     void skb_add_mtu(int mtu)
1099     {
1100          /* Must match allocation in alloc_skb */
1101          mtu = SKB_DATA_ALIGN(mtu) + sizeof(struct skb_shared_info);
1102
1103          kmem_add_cache_size(mtu);
1104     }
1105     #endif
```
———————————————————————————————————————————— *[include/net/socket.h]*

## 10.3.  Allocation of skbuffs

### 10.3.1.  alloc_skb, tcp_alloc_skb

*File :* `[include/net/tcp.h]`

These functions allocate a network buffer for output. The alloc_skb takes 2 arguments the size in bytes of the data area requested and the set of flags that tell the memory allocator how to behave. For the most part the network code calls the memory allocator with the `GFP_ATOMIC` set of flags: do not return without completing the task (for the moment this is equivalent to the `__GFP_HIGH` flag : can use emergency pools). The memory allocated for the data area, of course, is contiguous physical memory. The current

slab allocator provides size-N caches in power of 2 sizes from 32 bytes to 128 KB.



───────────────────────────────────────────────────────────────── *[net/core/skbuff.c]*

```
125     struct sk_buff *alloc_skb(unsigned int size, int gfp_mask)
126     {
127          struct sk_buff *skb;
128          u8 *data;
129
```
───────────────────────────────────────────────────────────────── *[net/core/skbuff.c]*

An skbuff head is allocated from the `skbuff_head_cache` slab. DMA suitable memory is not needed for the skbuff header, because no I/O is performed over its data, thus the __GFP_DMA flag is reset in the call for the allocation in the case the alloc_skb function was called with the flag set. If the allocation fails the function returns NULL .

───────────────────────────────────────────────────────────────── *[net/core/skbuff.c]*

```
130          /* Get the HEAD */
131          skb = kmem_cache_alloc(skbuff_head_cache,
132                              gfp_mask & ~__GFP_DMA);
133          if (!skb)
134              goto out;
```
───────────────────────────────────────────────────────────────── *[net/core/skbuff.c]*

If it succeeds then it allocates the skbuff data area from one of the size-N slabs using the `kmalloc()` function. The size of the data area requested through `kmalloc()` is augmented with the size of the `skb_shared_info` area that stores the additional information needed by fragmented skbuffs and a reference counter needed by the sharing mechanism. The SKB_DATA_ALIGN macro adds enough bytes to the requested size so that the shared_info area can be aligned with a level 1 cache line ( on P4 for example the X86_L1_CACHE_SHIFT is 7 and L1_CACHE_BYTES is then $2^7$=128 bytes and so the size is rounded up to the next 128 multiple ).

───────────────────────────────────────────────────────────────── *include/linux/cache.h*

```
12     #ifndef SMP_CACHE_BYTES
13     #define SMP_CACHE_BYTES L1_CACHE_BYTES
14     #endif
```

*——————————————————————————————— include/linux/cache.h*

*——————————————————————————————— include/linux/skbuff.h*

```
39      #define SKB_DATA_ALIGN(X) (((X) + (SMP_CACHE_BYTES - 1)) & \
40                        ~(SMP_CACHE_BYTES - 1))
```

*——————————————————————————————— include/linux/skbuff.h*

*——————————————————————————————— net/core/skbuff.c*

```
137           size = SKB_DATA_ALIGN(size);
138           data = kmalloc(size + sizeof(struct skb_shared_info), gfp_mask);
139           if (!data)
140                goto nodata;
```

*——————————————————————————————— net/core/skbuff.c*

If it fails in allocating the data area it gives back the area for the skbuff head and returns NULL.

*——————————————————————————————— net/core/skbuff.c*

```
154           skb_shinfo(skb)->frag_list = NULL;
155     out:
156           return skb;
157     nodata:
158           kmem_cache_free(skbuff_head_cache, skb);
159           skb = NULL;
160           goto out;
```

*——————————————————————————————— net/core/skbuff.c*

Then it initializes to 0 all bytes of the skbuff head up to the truesize field. The remaining bytes are not zeroed because they will be immediately initialized with the appropriate values (pointers to the data area of the skbuff and size). The skb truesize is initialized to the total allocated size : the requested data size plus the size of the skbuff header.

*——————————————————————————————— net/core/skbuff.c*

```
141
142           memset(skb, 0, offsetof(struct sk_buff, truesize));
```

*——————————————————————————————— net/core/skbuff.c*

Then the skbuff pointers inside the data area are initialized to a 0 size area :

*——————————————————————————————— net/core/skbuff.c*

```
144           atomic_set(&skb->users, 1);
145           skb->head = data;
146           skb->data = data;
147           skb->tail = data;
148           skb->end  = data + size;
```

*——————————————————————————————— net/core/skbuff.c*

*Figure 24 skbuff header and data area immediately after allocation*

The reference count is set to 1 as the header currently allocated is the only one referencing the data area

───────────────────────────────────────────────────────────── *net/core/skbuff.c*

```
150          atomic_set(&(skb_shinfo(skb)->dataref), 1);
151          skb_shinfo(skb)->nr_frags  = 0;
```
───────────────────────────────────────────────────────────── *net/core/skbuff.c*

and the counter for the external pages is set to 0 which indicates this is not a paged skbuff.

───────────────────────────────────────────────────────────── *net/core/skbuff.c*

```
152          skb_shinfo(skb)->tso_size = 0;
153          skb_shinfo(skb)->tso_segs = 0;
154          skb_shinfo(skb)->frag_list = NULL;
```
───────────────────────────────────────────────────────────── *net/core/skbuff.c*

The tso fields (for the TCP Segmentation Offloading support) are set to 0 because the network buffer is allocated for a generic interface and eventually only later in the process of sending the segment over an interface they can be set differently if the specific interface supports TSO . This initial allocation creates a simple skbuff, and not an skbuff chain and this is indicated setting the fragment list to NULL.



*Figure 25 Immediately after alloc_skb all the room in an skbuff is in the tail*

No space is reserved as headroom with alloc_skb. In case a headroom is needed it has to be reserved separately.

The tcp code instead uses its own skbuff allocator `tcp_alloc_skb` to request additional `MAX_TCP_HEADER` bytes to accomodate a headroom sufficient for the TCP/IP header (usually this is 128+ MAX_HEADER 32= MAX_TCP_HEADER = 160 bytes).This function allocates only a linear skbuff. The

tcp_alloc_skb function is simply a 1-line wrapper around the more general tcp_alloc_pskb function that can allocate paged skbuffs when the mem argument is different from 0 (in this case the size argument is only the size of the linear part of the skbuff and mem is the size of the paged part).

*————————————————————————————————————————————————— net/core/skbuff.c*

```
1841    static inline struct sk_buff *tcp_alloc_pskb(struct sock *sk, int size, int
mem, int gfp)
1842    {
1843          struct sk_buff *skb = alloc_skb(size+MAX_TCP_HEADER, gfp);
1844
1845          if (skb) {
1846                skb->truesize += mem;
1847                if (sk->sk_forward_alloc >= (int)skb->truesize ||
1848                    tcp_mem_schedule(sk, skb->truesize, 0)) {
1849                      skb_reserve(skb, MAX_TCP_HEADER);
1850                      return skb;
1851                }
1852                __kfree_skb(skb);
1853          } else {
1854                tcp_enter_memory_pressure();
1855                tcp_moderate_sndbuf(sk);
1856          }
1857          return NULL;
1858    }
1859
1860    static inline struct sk_buff *tcp_alloc_skb(struct sock *sk, int size, int gfp)
1861    {
1862          return tcp_alloc_pskb(sk, size, 0, gfp);
1863    }
```

*————————————————————————————————————————————————— net/core/skbuff.c*

The tcp_alloc functions reserve as headroom with skb_reserve the additional space allocated for a maximum TCP header (MAX_TCP_HEADER = 128 usually).

| Head Room | Tail Room |
|---|---|

*Figure 26 tcp_alloc_skb allocates additional space sufficient for headers, and reserves it as headroom*

Another allocation of skbuffs in the output process that requires some size adjustment is in the IP fragmentation case. In this situation the IP datagram has to be split and additional bytes for the IP header and the link layer header are requested for each of the chunks following the first one (the link layer size requested is rounded up to allow an efficient alignment of the IP header) :

*————————————————————————————————————————————————— [net/ipv4/ip_output.c]*

```
591             if ((skb2 = alloc_skb(len+hlen+LL_RESERVED_SPACE(rt->u.dst.dev),
GFP_ATOMIC)) == NULL) {
592                   NETDEBUG(printk(KERN_INFO "IP: frag: no memory for new frag-
ment!0));
593                   err = -ENOMEM;
594                   goto fail;
```

```
595                 }
596
597                 /*
598                  *    Set up data on packet
599                  */
600
601                 ip_copy_metadata(skb2, skb);
602                 skb_reserve(skb2, LL_RESERVED_SPACE(rt->u.dst.dev));
```
———————————————————————————————————————————————— *[net/ipv4/ip_output.c]*


( `LL_RESERVED_SPACE` is the link layer header size + space to properly align the IP header  to mod
`HH_DATA_MOD` = 16 bytes).


| Head Room | Tail Room |
|-----------|-----------|
| | |

*Figure 27 the IP fragmentation code allocates additional space for the current IP header and the link
header, and reserves as headroom the space for the link header*


### 10.3.2. sock_alloc_send_pskb,sock_alloc_send_skb

These functions allocate an skbuff for a socket.  There is no receive equivalent because when a packet
arrives, the socket is not known immediately.  sock_alloc_send_skb can allocate only linear skbuffs, and so
it is simply a 1-line wrapper around the sock_alloc_send_pskb that can allocate also paged skbuffs.
header_len is the size to be allocated in the linear part of the skbuff, while data_len as usual is the size to be
allocated using external pages.  The function sock_alloc_send_pskb checks if the memory limit allows the
allocation.  Then it allocates a linear skbuff for header_len bytes.  If paged memory is requested (data_len
!= 0), then it computes the number of pages necessary, allocates the pages, and initializes the pointers to
them in the skbuff shared_info area.  It then sets the owner of the skbuff to the socket and returns the
skbuff.  This function can wait until memory is available : when it starts it sets a timeout on the socket, if
memory is not available it sleeps on memory until it is available or the timeout expires and it reports an
error.

———————————————————————————————————————————————— *[net/core/sock.c]*
```
782    struct sk_buff *sock_alloc_send_pskb(struct sock *sk, unsigned long header_len,
783                            unsigned long data_len, int noblock, int *errcode)
784    {
785         struct sk_buff *skb;
786         unsigned int gfp_mask;
787         long timeo;
788         int err;
789
790         gfp_mask = sk->sk_allocation;
791         if (gfp_mask & __GFP_WAIT)
792             gfp_mask |= __GFP_REPEAT;
793
794         timeo = sock_sndtimeo(sk, noblock);
795         while (1) {
796             err = sock_error(sk);
```

```
797                    if (err != 0)
798                         goto failure;
799
800                    err = -EPIPE;
801                    if (sk->sk_shutdown & SEND_SHUTDOWN)
802                         goto failure;
803
804                    if (atomic_read(&sk->sk_wmem_alloc) < sk->sk_sndbuf) {
805                         skb = alloc_skb(header_len, sk->sk_allocation);
806                         if (skb) {
807                              int npages;
808                              int i;
809
810                              /* No pages, we're done... */
811                              if (!data_len)
812                                   break;
813
814                              npages = (data_len + (PAGE_SIZE - 1)) >> PAGE_SHIFT;
815                              skb->truesize += data_len;
816                              skb_shinfo(skb)->nr_frags = npages;
817                              for (i = 0; i < npages; i++) {
818                                   struct page *page;
819                                   skb_frag_t *frag;
820
821                                   page = alloc_pages(sk->sk_allocation, 0);
822                                   if (!page) {
823                                        err = -ENOBUFS;
824                                        skb_shinfo(skb)->nr_frags = i;
825                                        kfree_skb(skb);
826                                        goto failure;
827                                   }
828
829                                   frag = &skb_shinfo(skb)->frags[i];
830                                   frag->page = page;
831                                   frag->page_offset = 0;
832                                   frag->size = (data_len >= PAGE_SIZE ?
833                                             PAGE_SIZE :
834                                             data_len);
835                                   data_len -= PAGE_SIZE;
836                              }
837
838                              /* Full success... */
839                              break;
840                         }
841                         err = -ENOBUFS;
842                         goto failure;
843                    }
844                    set_bit(SOCK_ASYNC_NOSPACE, &sk->sk_socket->flags);
845                    set_bit(SOCK_NOSPACE, &sk->sk_socket->flags);
846                    err = -EAGAIN;
847                    if (!timeo)
848                         goto failure;
849                    if (signal_pending(current))
850                         goto interrupted;
```

```
851                    timeo = sock_wait_for_wmem(sk, timeo);
852          }
853
854          skb_set_owner_w(skb, sk);
855          return skb;
856
857   interrupted:
858          err = sock_intr_errno(timeo);
859   failure:
860          *errcode = err;
861          return NULL;
862   }
863
864   struct sk_buff *sock_alloc_send_skb(struct sock *sk, unsigned long size,
865                            int noblock, int *errcode)
866   {
867          return sock_alloc_send_pskb(sk, size, 0, noblock, errcode);
868   }
```
———————————————————————————————————————————————————————— *[net/core/sock.c]*


### 10.3.3. sock_wmalloc,sock_rmalloc

sk_sndbuf and sk_rcvbuf are two socket variables initially set to the global variables syctl configurable syctl_wmem_default,sysctl_rmem_default. Their values can later be tuned, for instance when a tcp socket enters the established state. These 2 functions allocate an skbuff of the requested size checking, unless the argument force is set, that the memory for reading or writing owned by the socket is less than the memory allowed.

———————————————————————————————————————————————————————— *[net/core/sock.c]*
```
694     struct sk_buff *sock_wmalloc(struct sock *sk, unsigned long size, int force,
int priority)
695     {
696          if (force || atomic_read(&sk->sk_wmem_alloc) < sk->sk_sndbuf) {
697              struct sk_buff * skb = alloc_skb(size, priority);
698              if (skb) {
699                  skb_set_owner_w(skb, sk);
700                  return skb;
701              }
702          }
703          return NULL;
704     }
705
706     /*
707      * Allocate a skb from the socket's receive buffer.
708      */
709     struct sk_buff *sock_rmalloc(struct sock *sk, unsigned long size, int force,
int priority)
710     {
711          if (force || atomic_read(&sk->sk_rmem_alloc) < sk->sk_rcvbuf) {
712              struct sk_buff *skb = alloc_skb(size, priority);
713              if (skb) {
714                  skb_set_owner_r(skb, sk);
715                  return skb;
```

```
716                       }
717              }
718              return NULL;
719      }
```
———————————————————————————————————————————————————— *[net/core/sock.c]*

bt_skb_alloc [include/net/bluetooth/bluetooth.h] is a special
skb allocation function, it simply allocates and reserves as headroom
additional 8 bytes, and initializes to 0 a control variable kept
in the skb control buffer.

dn_alloc_skb [net/decnet/dn_nsp_out.c] is a special
allocation function for decnet. it
allocates additional 64 bytes for header and reserves them as headroom


### 10.3.4. dev_alloc_skb and __dev_alloc_skb

These function are used to allocate an skbuff for an incoming packet at the device driver level. The
dev_alloc_skb function is a 1-line wrapper that immediately calls the __dev_alloc_skb function. The only
argument for dev_alloc_skb is the size in bytes of the requested skb. This function calls the double under-
lined function with the GFP_ATOMIC specification. This __dev_alloc_skb function doesn't use a parame-
terized link layer header size, it simply allocates 16 more bytes than those requested to leave space for an
ethernet header (14 bytes) and keep the IP header 16 bytes aligned for efficiency reasons (cache). This
additional space is reserved as headroom calling skb_reserve.

———————————————————————————————————————————————— *[include/linux/skbuff.h]*
```
1056    static inline struct sk_buff *__dev_alloc_skb(unsigned int length,
1057                                      int gfp_mask)
1058    {
1059            struct sk_buff *skb = alloc_skb(length + 16, gfp_mask);
1060            if (likely(skb))
1061                    skb_reserve(skb, 16);
1062            return skb;
1063    }
1064
1065    /**
1066     *   dev_alloc_skb - allocate an skbuff for sending
1067     *   @length: length to allocate
1068     *
1069     *   Allocate a new &sk_buff and assign it a usage count of one. The
1070     *   buffer has unspecified headroom built in. Users should allocate
1071     *   the headroom they think they need without accounting for the
1072     *   built in space. The built in space is used for optimisations.
1073     *
1074     *   %NULL is returned in there is no free memory. Although this function
1075     *   allocates memory it can be called from an interrupt.
1076     */
1077    static inline struct sk_buff *dev_alloc_skb(unsigned int length)
1078    {
1079            return __dev_alloc_skb(length, GFP_ATOMIC);
1080    }
1081
```
———————————————————————————————————————————————— *[include/linux/skbuff.h]*

The __dev_alloc_skb function allocates the skb requesting 16 additional bytes of data space using the standard alloc_skb (this allocates an skb having a data area of the size requested plus the size of the shared_info struct). Then _dev_alloc_skb shifts the data and tail pointers of 16 bytes with skb_reserve.

| Head Room | Tail Room |
|-----------|-----------|

*Figure 28 these functions allocate and reserve as headroom 16 bytes for the link header*

## 10.4. Data Pointer Manipulation

### 10.4.1. skb_reserve

*File :* `[include/linux/skbuff.h]`
When an skbuff is allocated all the free room is allocated in the tail, after the data and tail pointers position.



*Figure 29 Skbuff before any headroom reservation*

this function adjusts the skbuff headroom (at the beginning there is no headroom and all tailroom) just after the creation (the skbuff should be empty). This is done moving the data and tail pointers, that just after creation point to the `skb->head,` by `len` bytes.

_____ *[include/linux/skbuff.h]*

```
944      *    skb_reserve - adjust headroom
945      *    @skb: buffer to alter
946      *    @len: bytes to move
947      *
948      *    Increase the headroom of an empty &sk_buff by reducing the tail
949      *    room. This is only allowed for an empty buffer.
```
*―――――――――――――――――――――――――――――――――――――――― [include/linux/skbuff.h]*



*Figure 30 Skbuff after reserving some headroom*

### 10.4.2.  skb_put and __skb_put

*File :* [include/linux/skbuff.h]
As for other functions there are 2 versions of the skb_put function.  The __skb_put() function saves just a consistency check on the sufficency of data space (tail > end).  After calling the skb_reserve() function to move the data pointer, skb_put/__skb_put is then usually called to move the tail pointer and prepare the space for copying over the data.  It updates the len field of the skbuff header, it doesn't perform any copy. This operation can be applied only on linear skbuffs.

```
                                    ( skb_put )

      SKB_LINEAR_ASSERT          unlikely              skb_over_panic


                                  ( __skb_put )


                             SKB_LINEAR_ASSERT
```

_____ *[include/linux/skbuff.h]*

```
815    #define SKB_LINEAR_ASSERT(skb)  BUG_ON(skb_is_nonlinear(skb))
```

_____ *[include/linux/skbuff.h]*

```
820    static inline unsigned char *__skb_put(struct sk_buff *skb, unsigned int len)
821    {
822        unsigned char *tmp = skb->tail;
823        SKB_LINEAR_ASSERT(skb);
824        skb->tail += len;
825        skb->len  += len;
826        return tmp;
827    }
```

_____ *[include/linux/skbuff.h]*

```
838    static inline unsigned char *skb_put(struct sk_buff *skb, unsigned int len)
839    {
840        unsigned char *tmp = skb->tail;
```

```
841          SKB_LINEAR_ASSERT(skb);
842          skb->tail += len;
843          skb->len  += len;
844          if (unlikely(skb->tail>skb->end))
845              skb_over_panic(skb, len, current_text_addr());
846          return tmp;
847      }
```

————————————————————————————————————————— *[include/linux/skbuff.h]*



*Figure 31 An skbuff after calling skb_put*

### 10.4.3. skb_add_data

*File :* `[net/ipv4/tcp.c]`
This function is inside the tcp code because it is used only by TCP. It uses the skb_put function to make room (advancing the tail pointer) for copying `copy` bytes of data. Differently from skb_put, this function also copies the data with the csum_and_copy_from_user or copy_from_user function accordingly if the interface doesn't support or supports the checksum in hardware.

—————————————————————————————————————————— *[net/ipv4/tcp.c]*

```
995      static inline int skb_add_data(struct sk_buff *skb, char *from, int copy)
996      {
997            int err = 0;
998            unsigned int csum;
999            int off = skb->len;
1000
1001            if (skb->ip_summed == CHECKSUM_NONE) {
1002                  csum = csum_and_copy_from_user(from, skb_put(skb, copy),
1003                                      copy, 0, &err);
1004                  if (!err) {
1005                        skb->csum = csum_block_add(skb->csum, csum, off);
1006                        return 0;
1007                  }
1008            } else {
1009                  if (!copy_from_user(skb_put(skb, copy), from, copy))
1010                        return 0;
1011            }
1012
1013            __skb_trim(skb, off);
1014            return -EFAULT;
1015      }
```

—————————————————————————————————————————— *[net/ipv4/tcp.c]*

### 10.4.4. skb_push

*File :* `[include/linux/skbuff.h]`



This function is usually called to prepare the space where to prepend protocol headers. For example in the net/ipv4/tcp_output.c file the skbuff is adjusted for the tcp header space with

—————————————————————————————————————————— *[include/linux/skbuff.h]*

```
227                   }
228                   th = (struct tcphdr *) skb_push(skb, tcp_header_size);
```

—————————————————————————————————————————— *[include/linux/skbuff.h]*

after calling this function the result is the new `skb->data` pointer (the new beginning of the data area) and the header is then copied from there on.

*———————————————————————— [include/linux/skbuff.h]*

```
849     static inline unsigned char *__skb_push(struct sk_buff *skb, unsigned int len)
850     {
851             skb->data -= len;
852             skb->len  += len;
853             return skb->data;
854     }
```

*———————————————————————— [include/linux/skbuff.h]*

```
865     static inline unsigned char *skb_push(struct sk_buff *skb, unsigned int len)
866     {
867             skb->data -= len;
868             skb->len  += len;
869             if (unlikely(skb->data<skb->head))
870                     skb_under_panic(skb, len, current_text_addr());
871             return skb->data;
872     }
```

*———————————————————————— [include/linux/skbuff.h]*

The `skb->data` pointer pointing at the beginning of the used data area is shrunk of len bytes and the len of the data area used in the skbuff is increased of the same number. In the unlikely case in which the `skb->data` pointer with this operation goes outside the skbuff available data area the kernel panics with an "skput: under .. " message. At the end the function returns the updated `skb->data` pointer. The skb_push function simply doesn't perform the consistency check.

| Head Room | push area | Data Area | Tail Room |
|-----------|-----------|-----------|-----------|

| Head Room | Data Area | Tail Room |
|-----------|-----------|-----------|

*Figure 32 Skbuff before and after a push operation*

### 10.4.5.  skb_pull()

*File :* `[include/linux/skbuff.h]`

The skb_pull function is used to strip away the headers during input processing. In the pull operation if the len by which you ask to decrease the used data area is larger then the actual skb->len then the function returns a NULL. Otherwise the actual skb->len is decreased by the requested amount and the skb->data pointer is augmented by the same amount. A consistency check is performed inside __skb_pull to check that skb->len is greater than or equal to skb->data_len. If that is not true, this would mean that the beginning of the data would be inside the nonlinear part of the skbuff, and so it would need to be rearranged, for this reason a bug is raised.

——————————————————————————————————————————— *[include/linux/skbuff.h]*

```
891     static inline unsigned char *skb_pull(struct sk_buff *skb, unsigned int len)
892     {
893             return unlikely(len > skb->len) ? NULL : __skb_pull(skb, len);
894     }
```

——————————————————————————————————————————— *[include/linux/skbuff.h]*

```
874     static inline unsigned char *__skb_pull(struct sk_buff *skb, unsigned int len)
875     {
876             skb->len -= len;
877             BUG_ON(skb->len < skb->data_len);
878             return skb->data += len;
879     }
```
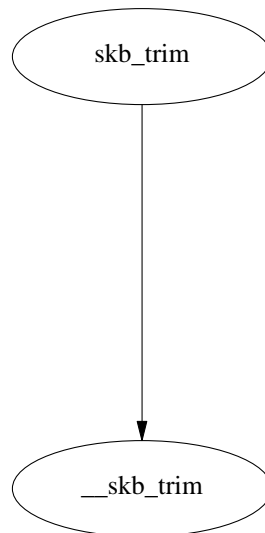
——————————————————————————————————————————— *[include/linux/skbuff.h]*

| Head Room | Data Area | Tail Room |
|---|---|---|

| Head Room | skb_pulled area | Data Area | Tail Room |
|---|---|---|---|

*Figure 33 Skbuff before and after a pull operation*
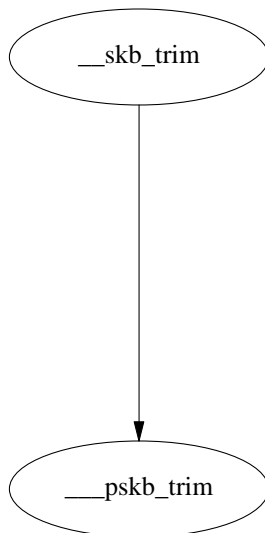
### 10.4.6. skb_trim, __skb_trim, __pskb_trim

*File :* `include/linux/skbuff.h`

```
              ╭─────────────╮
             (   skb_trim   )
              ╰──────┬──────╯
                     │
                     ▼
              ╭─────────────╮
             (  __skb_trim  )
              ╰─────────────╯
```

These functions trim the data area at its end. skb_trim is the highest level function. It is a wrapper function
that makes a consistency check to see if the total data in the skbuff is sufficient to satisfy the request and
then calls __skb_trim()

───────────────────────────────────────────────────────────────── *include/linux/skbuff.h*

```
969      *   skb_trim - remove end from a buffer
970      *   @skb: buffer to alter
971      *   @len: new length
972      *
973      *   Cut the length of a buffer down by removing data from the tail. If
974      *   the buffer is already under the length specified it is not modified.
```
───────────────────────────────────────────────────────────────── *include/linux/skbuff.h*

If the skbuff is linear (skb->data_len == 0) __skb_trim simply sets the total data length to len and trims the skb->tail pointer. Otherwise if the skbuff is nonlinear it calls the paged skbuff version ___pskb_trim().

———————————————————————————————————————— *include/linux/skbuff.h*

```
951     static inline void skb_reserve(struct sk_buff *skb, unsigned int len)
952     {
953          skb->data += len;
954          skb->tail += len;
955     }
956
957     extern int ___pskb_trim(struct sk_buff *skb, unsigned int len, int realloc);
958
959     static inline void __skb_trim(struct sk_buff *skb, unsigned int len)
960     {
```

———————————————————————————————————————— *include/linux/skbuff.h*



This function is called to trim at the end a generic skbuff (be it linear or not). The len argument will be the

new size of the data area. The operation is used for instance to remove any padding added by other network layers to an IP datagram after having determined the correct size from the IP header :

*—————————————————————————————————————————— net/ipv4/ip_input.c*

```
418                     if (skb->len > len) {
419                             __pskb_trim(skb, len);
420                         if (skb->ip_summed == CHECKSUM_HW)
421                             skb->ip_summed = CHECKSUM_NONE;
422                     }
```

*—————————————————————————————————————————— net/ipv4/ip_input.c*

or to trim back an skbuff to its initial size after having attempted unsuccessfully to add some data to it (see skb_add_data) .

An skbuff can have an array of external pages ( skb->frags[] ) and/or be an skbuff chain ( skb->fraglist ).

*—————————————————————————————————————————— net/core/skbuff.c*

```
649    int ___pskb_trim(struct sk_buff *skb, unsigned int len, int realloc)
650    {
651          int offset = skb_headlen(skb);
652          int nfrags = skb_shinfo(skb)->nr_frags;
653          int i;
654
655          for (i = 0; i < nfrags; i++) {
656                  int end = offset + skb_shinfo(skb)->frags[i].size;
657                  if (end > len) {
658                          if (skb_cloned(skb)) {
659                                  if (!realloc)
660                                      BUG();
661                                  if (pskb_expand_head(skb, 0, 0, GFP_ATOMIC))
662                                      return -ENOMEM;
663                          }
664                          if (len <= offset) {
665                                  put_page(skb_shinfo(skb)->frags[i].page);
666                                  skb_shinfo(skb)->nr_frags--;
667                          } else {
668                                  skb_shinfo(skb)->frags[i].size = len - offset;
669                          }
670                  }
671                  offset = end;
672          }
673
674          if (offset < len) {
675                  skb->data_len -= skb->len - len;
676                  skb->len       = len;
677          } else {
678                  if (len <= skb_headlen(skb)) {
679                          skb->len       = len;
680                          skb->data_len = 0;
681                          skb->tail     = skb->data + len;
682                          if (skb_shinfo(skb)->frag_list && !skb_cloned(skb))
683                              skb_drop_fraglist(skb);
684                  } else {
```

```
685                         skb->data_len -= skb->len - len;
686                         skb->len        = len;
687                 }
688         }
689
690         return 0;
691     }
```
———————————————————————————————————————————————— *net/core/skbuff.c*


We know that an skbuff has multiple data pages associated with it if the number in the
`skb_shared_info` structure `skb_shinfo(skb)->nr_frags` is different from zero. In this case
we run through the pages until eventually their total size reaches the requested len. If this happen before
the end of the array the page is relinquished and the number of fragments is decresed by 1.


———————————————————————————————————————————————— *include/linux/skbuff.h*
```
662                             return -ENOMEM;
663                     }
664                     if (len <= offset) {
665                             put_page(skb_shinfo(skb)->frags[i].page);
666                             skb_shinfo(skb)->nr_frags--;
667                     } else {
668                             skb_shinfo(skb)->frags[i].size = len - offset;
669                     }
670                 }
671             offset = end;
672         }
673
674         if (offset < len) {
675             skb->data_len -= skb->len - len;
676             skb->len        = len;
677         } else {
678             if (len <= skb_headlen(skb)) {
679                 skb->len      = len;
680                 skb->data_len = 0;
```
———————————————————————————————————————————————— *include/linux/skbuff.h*

if the data in the pages associated with this skbuff is not enough to satisfy the request then we simply reset
the total data len of the skbuff to len and we decrease the length of data in the remaining skbuffs
(skb->data_len) by the proper amount. offset is here the total data in all the array of pages associated with
the 1st skbuff (while headlen is the data in this skbuff) Otherwise there are 2 possibilities : - if the data in
the skbuff is enough we reset the skbuff to
 a linear one, we set the len to the requested one, we just
 trime the tail of this skbuff, and eventually we drop all the
 other fragments in the fraglist - we still need some fragments .. in this case we reset the length
 to len and we decrease the skb->data_len of the proper amount

### 10.4.7. tcp_trim_head, __pskb_trim_head
These functions are used only by TCP. They cut len bytes from the beginning of the data area.
tcp_trim_head is the more general of these functions. If the skbuff is cloned (2 headers using the same
data) then it is copied over to a newly allocated linear skbuff (pskb_expand_head). If there are len bytes or
more in the linear part of the skbuff (for sure this happens if it is a linear skbuff) then this operation reduces
to a pull operation (__skb_pull), otherwise the paged version of this function is called (__pskb_trim_head).
__pskb_trim_head runs through the external pages of the paged skbuff and releases all the pages comprised

in the first len bytes. Then it adjusts the number of pages left and the page_offset of the first page, and it collapses the data area in the linear part ( skb->tail = skb->data ).

*———————————————————————————————— [net/ipv4/tcp_output.c]*
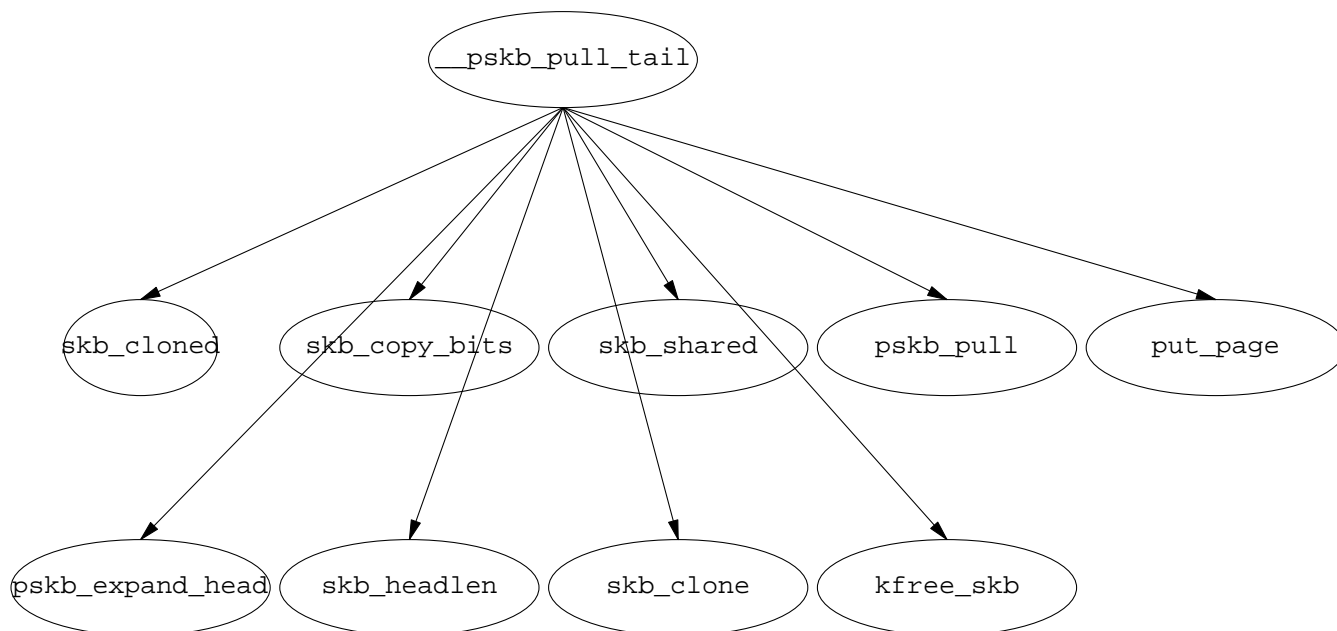
```
486    unsigned char * __pskb_trim_head(struct sk_buff *skb, int len)
487    {
488          int i, k, eat;
489
490          eat = len;
491          k = 0;
492          for (i=0; i<skb_shinfo(skb)->nr_frags; i++) {
493                if (skb_shinfo(skb)->frags[i].size <= eat) {
494                      put_page(skb_shinfo(skb)->frags[i].page);
495                      eat -= skb_shinfo(skb)->frags[i].size;
496                } else {
497                      skb_shinfo(skb)->frags[k] = skb_shinfo(skb)->frags[i];
498                      if (eat) {
499                            skb_shinfo(skb)->frags[k].page_offset += eat;
500                            skb_shinfo(skb)->frags[k].size -= eat;
501                            eat = 0;
502                      }
503                      k++;
504                }
505          }
506          skb_shinfo(skb)->nr_frags = k;
507
508          skb->tail = skb->data;
509          skb->data_len -= len;
510          skb->len = skb->data_len;
511          return skb->tail;
512    }
513
514    static int tcp_trim_head(struct sock *sk, struct sk_buff *skb, u32 len)
515    {
516          if (skb_cloned(skb) &&
517              pskb_expand_head(skb, 0, 0, GFP_ATOMIC))
518                return -ENOMEM;
519
520          if (len <= skb_headlen(skb)) {
521                __skb_pull(skb, len);
522          } else {
523                if (__pskb_trim_head(skb, len-skb_headlen(skb)) == NULL)
524                      return -ENOMEM;
525          }
526
527          TCP_SKB_CB(skb)->seq += len;
528          skb->ip_summed = CHECKSUM_HW;
529          return 0;
530    }
```

*———————————————————————————————— [net/ipv4/tcp_output.c]*

### 10.4.8. __pskb_pull_tail

*File :* [net/ipv4/tcp_output.c]



The purpose of this function is to put the headers of a packet in the linear part of an skbuff, this is required to parse the headers with the usual pointer arithmetic that is possible only on a contiguous data area. This function expands the data area in the linear skbuff part of a fragmented skbuff copying the data from the remaining fragments. If that space is not sufficient (delta - (skb->end - skb->tail) > 0 ), then it allocates a new data area and copies the data from the old to the new one and updates pointers in the descriptor. delta are the more bytes requested in the linear skbuff part. eat is the part of them that is not possible to allocate in the current linear part of the skbuff. If eat <=0 then we can keep the current skbuff data area and just update the pointers. If eat > 0 then we have to allocate a new linear part and in this case we will request 128 additional bytes to accomodate eventual future requests. It allocates a new linear part also in the case the skbuff has been cloned since that data area needs to be changed. Then delta bytes are copied from the fragmented tail of the old skbuff (those after headlen) to the the tail. To update the skbuff now if there is no frag_list then it is sufficient to pull the array of pages, otherwise  the frag_list needs to be scanned. If the array of pages associated with this skbuff is large enough then again it is sufficient to pull the array. Then going through the frag_list and eating (kfree_skb) all the complete skbuffs that can be eaten. When you are here it means that you are on an skbuff that you cant eat completely . For this you go through the array of pages and you free all those that you can completely eat (put_page). For the last page you update the page_offset and size values in the frags[k] structure.

*—————————————————————————————————————————————— [net/ipv4/tcp_output.c]*

```
718     unsigned char *__pskb_pull_tail(struct sk_buff *skb, int delta)
719     {
720             /* If skb has not enough free space at tail, get new one
721              * plus 128 bytes for future expansions. If we have enough
722              * room at tail, reallocate without expansion only if skb is cloned.
723              */
724             int i, k, eat = (skb->tail + delta) - skb->end;
725
726             if (eat > 0 || skb_cloned(skb)) {
727                     if (pskb_expand_head(skb, 0, eat > 0 ? eat + 128 : 0,
```

```
728                                  GFP_ATOMIC))
729                     return NULL;
730         }
731
732         if (skb_copy_bits(skb, skb_headlen(skb), skb->tail, delta))
733             BUG();
734
735         /* Optimization: no fragments, no reasons to preestimate
736          * size of pulled pages. Superb.
737          */
738         if (!skb_shinfo(skb)->frag_list)
739             goto pull_pages;
740
741         /* Estimate size of pulled pages. */
742         eat = delta;
743         for (i = 0; i < skb_shinfo(skb)->nr_frags; i++) {
744             if (skb_shinfo(skb)->frags[i].size >= eat)
745                 goto pull_pages;
746             eat -= skb_shinfo(skb)->frags[i].size;
747         }
748
749         /* If we need update frag list, we are in troubles.
750          * Certainly, it possible to add an offset to skb data,
751          * but taking into account that pulling is expected to
752          * be very rare operation, it is worth to fight against
753          * further bloating skb head and crucify ourselves here instead.
754          * Pure masohism, indeed. 8)8)
755          */
756         if (eat) {
757             struct sk_buff *list = skb_shinfo(skb)->frag_list;
758             struct sk_buff *clone = NULL;
759             struct sk_buff *insp = NULL;
760
761             do {
762                 if (!list)
763                     BUG();
764
765                 if (list->len <= eat) {
766                     /* Eaten as whole. */
767                     eat -= list->len;
768                     list = list->next;
769                     insp = list;
770                 } else {
771                     /* Eaten partially. */
772
773                     if (skb_shared(list)) {
774                         /* Sucks! We need to fork list. :-( */
775                         clone = skb_clone(list, GFP_ATOMIC);
776                         if (!clone)
777                             return NULL;
778                         insp = list->next;
779                         list = clone;
780                     } else {
781                         /* This may be pulled without
```

```
782                            * problems. */
783                          insp = list;
784                    }
785                    if (!pskb_pull(list, eat)) {
786                          if (clone)
787                                kfree_skb(clone);
788                          return NULL;
789                    }
790                    break;
791              }
792        } while (eat);
793
794        /* Free pulled out fragments. */
795        while ((list = skb_shinfo(skb)->frag_list) != insp) {
796              skb_shinfo(skb)->frag_list = list->next;
797              kfree_skb(list);
798        }
799        /* And insert new clone at head. */
800        if (clone) {
801              clone->next = list;
802              skb_shinfo(skb)->frag_list = clone;
803        }
804    }
805    /* Success! Now we may commit changes to skb data. */
806
807 pull_pages:
808    eat = delta;
809    k = 0;
810    for (i = 0; i < skb_shinfo(skb)->nr_frags; i++) {
811          if (skb_shinfo(skb)->frags[i].size <= eat) {
812                put_page(skb_shinfo(skb)->frags[i].page);
813                eat -= skb_shinfo(skb)->frags[i].size;
814          } else {
815                skb_shinfo(skb)->frags[k] = skb_shinfo(skb)->frags[i];
816                if (eat) {
817                      skb_shinfo(skb)->frags[k].page_offset += eat;
818                      skb_shinfo(skb)->frags[k].size -= eat;
819                      eat = 0;
820                }
821                k++;
822          }
823    }
824    skb_shinfo(skb)->nr_frags = k;
825
826    skb->tail     += delta;
827    skb->data_len -= delta;
828
829    return skb->tail;
830 }
```
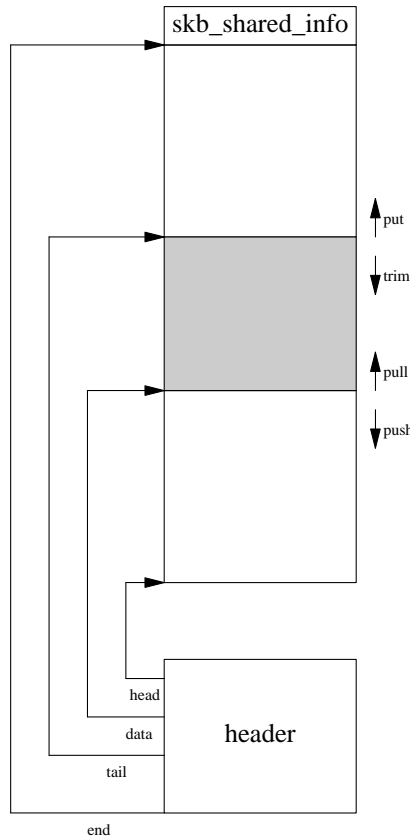
_____ *[net/ipv4/tcp_output.c]*

------------------ The pskb_.. and __pskb_.. functions refer to the fragmented skbuffs. As usual the __

functions perform less or no check at all.  -------------------

So the head and end pointers are fixed for an skbuff. The head points to the very beginning of the data area as obtained from kmalloc while the end points to the last useable area by the data. After it the skb_shared_info structure is kept.



The data and tail instead can be moved with the following operations :
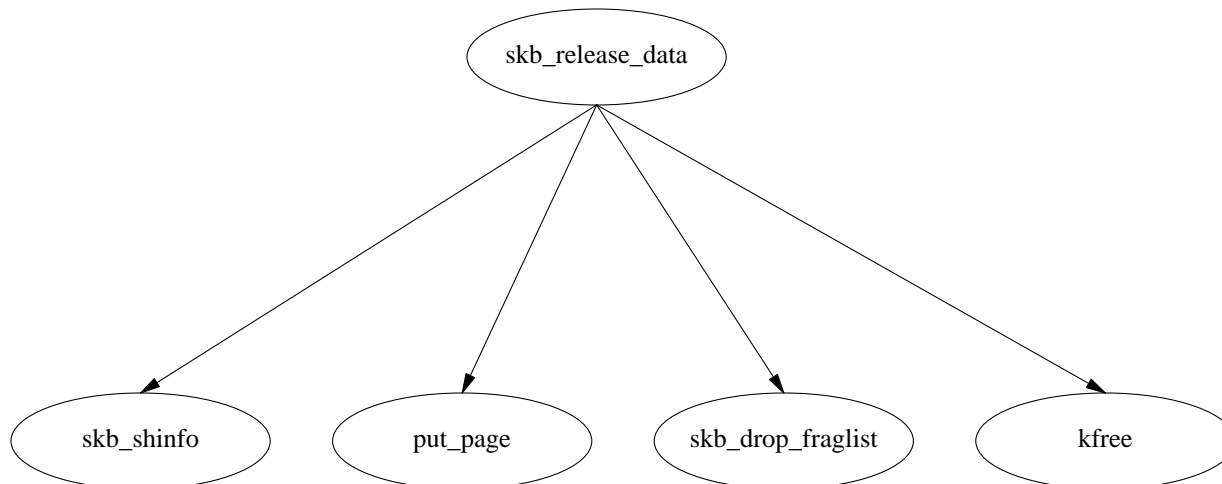
```
skb->data : push .. extends the used data area towards the beginning of
                    the buffer skb->head
            pull .. shrinks the beginning of the used data area
skb->tail : trim .. shrinks the end of the used data area
            put .. extends the end of the used data area towards skb->end
```

There are 2 implementations of each of these operations on skbuffs. One with consistency checks named ude/linux/skbiff.h] skb_put() kb_push() .. and so on. And one named with a prepended double underscore ( `__skb_push()`,... ) that doesnt apply any consistency check, this is used for eficiency reasons when it is clear that the checks are not needed.

## 10.5.  Releasing skbuffs

### 10.5.1.  skb_release_data

*File :* `[net/core/skbuff.c]`

```
                              skb_release_data
```

```
   skb_shinfo          put_page          skb_drop_fraglist          kfree
```

This function releases all the data areas associated with an skbuff if this skbuff was not cloned or the number of users is 0. In that case the function goes through the array of fragments and puts back to the page allocator the pages associated. Then if there is a frag_list it drops all the fragments (skb_drop_fraglist). and finally it frees the data area associated with the skbuff and consequently the skb_shared_info area.

———————————————————————————————————————————————— *[net/core/skbuff.c]*

```
184
185     void skb_release_data(struct sk_buff *skb)
186     {
187          if (!skb->cloned ||
188              atomic_dec_and_test(&(skb_shinfo(skb)->dataref))) {
189            if (skb_shinfo(skb)->nr_frags) {
190                 int i;
191                 for (i = 0; i < skb_shinfo(skb)->nr_frags; i++)
192                      put_page(skb_shinfo(skb)->frags[i].page);
193            }
194
195            if (skb_shinfo(skb)->frag_list)
196                 skb_drop_fraglist(skb);
197
198            kfree(skb->head);
199          }
```

———————————————————————————————————————————————— *[net/core/skbuff.c]*

The `skb->data_len` field is different from zero only on nonlinear skbuffs. In fact the function `skb_is_nonlinear()` returns that field. It represents the number of bytes in the nonlinear part of the skbuff. The function skb_headlen() returns the bytes in the linear part of the skbuff : `skb->len - skb->data_len`.


### 10.5.2. skb_drop_fraglist

*File :* [net/core/skbuff.c]

this function drops ( `kfree_skb(skb)` ) all the fragments of this skbuff and resets to null the skb->frag_list pointer. This function is called when all the data can be discarded or all the data in the skbuff sbk->len is in this skbuff and so the skbuff is reset to a linear one.
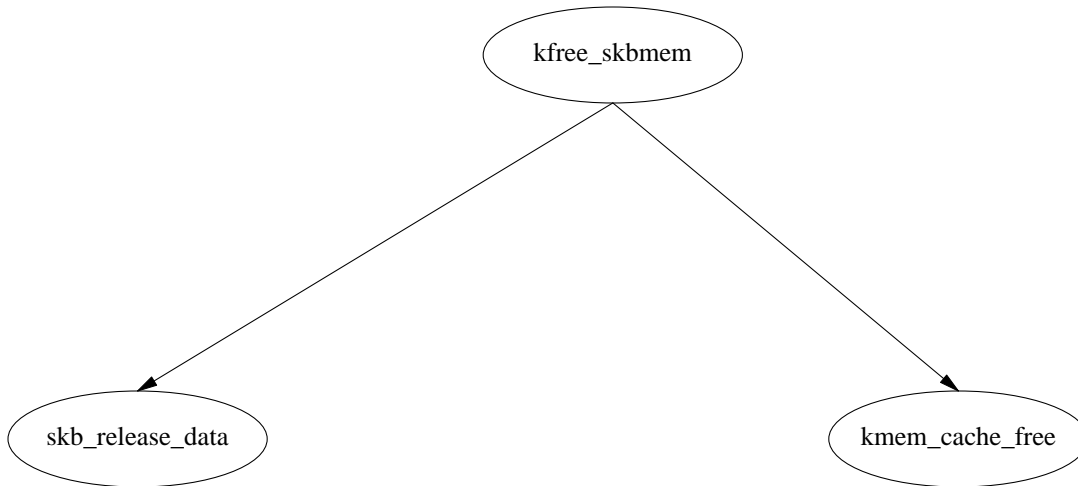
———————————————————————————————————————————————— *[net/core/skbuff.c]*

```
163
164     static void skb_drop_fraglist(struct sk_buff *skb)
165     {
166            struct sk_buff *list = skb_shinfo(skb)->frag_list;
167
168            skb_shinfo(skb)->frag_list = NULL;
169
170            do {
171                  struct sk_buff *this = list;
172                  list = list->next;
173                  kfree_skb(this);
174            } while (list);
```
———————————————————————————————————————————————— *[net/core/skbuff.c]*

### 10.5.3. kfree_skbmem

*File :* `[net/core/skbuff.c]`

```
                          ┌─────────────────┐
                          │  kfree_skbmem   │
                          └─────────────────┘
                         ╱                     ╲
                        ╱                       ╲
           ┌─────────────────┐         ┌─────────────────┐
           │ skb_release_data│         │ kmem_cache_free │
           └─────────────────┘         └─────────────────┘
```

This function releases all the memory associated with an skbuff. It doesnt clean its state. First it tries to release all data areas (this is not done if the data areas are in use). Then it frees the skbuff header from the appropriate slab.

——————————————————————————————————————————————— *[net/core/skbuff.c]*

```
201
202     /*
203      *   Free an skbuff by memory without cleaning the state.
204      */
205     void kfree_skbmem(struct sk_buff *skb)
206     {
207             skb_release_data(skb);
208             kmem_cache_free(skbuff_head_cache, skb);
209     }
```
——————————————————————————————————————————————— *[net/core/skbuff.c]*

### 10.5.4. kfree_skb()

*File :* [include/linux/skbuff.h]

```
                          ╭───────────────╮
                          │   kfree_skb   │
                          ╰───────┬───────╯
                                  │
                                  │
                                  ▼
                          ╭───────────────╮
                          │  __kfree_skb  │
                          ╰───────────────╯
```

Look at this code !!!!!!!  It means :
  if there is only 1 user (then this user of the skbuff is freeing it)
  of the skbuff call __kfree_skb(skb) and return.
  otherwise just decrement the number of users and return.

──────────────────────────────────────────────────────────── *[include/linux/skbuff.h]*

```
332    /*
333     * If users == 1, we are the only owner and are can avoid redundant
334     * atomic change.
335     */
336
337    /**
338     *   kfree_skb - free an sk_buff
339     *   @skb: buffer to free
340     *
341     *   Drop a reference to the buffer and free it if the usage count has
342     *   hit zero.
343     */
344    static inline void kfree_skb(struct sk_buff *skb)
```

──────────────────────────────────────────────────────────── *[include/linux/skbuff.h]*

### 10.5.5.  __kfree_skb

*File :* [net/core/skbuff.c]

This function cleans the state of the skbuff and releases any data area associated with it. Something went wrong if we came here and the skbuff is still on a list ( a socket or device list), print a kernel warning msg. Release the dst entry in the dst cache. If there is a destructor defined for the skb call it and eventually print a warning if we are executing out of an IRQ. Release all the memory associated with the skb calling kfree_skbmem.

──────────────────────────────────────────────────────────────────────── *[net/core/skbuff.c]*

```
201
202     /*
203      *   Free an skbuff by memory without cleaning the state.
204      */
205     void kfree_skbmem(struct sk_buff *skb)
206     {
207             skb_release_data(skb);
208             kmem_cache_free(skbuff_head_cache, skb);
209     }
210
211     /**
212      *   __kfree_skb - private function
213      *   @skb: buffer
214      *
215      *   Free an sk_buff. Release anything attached to the buffer.
216      *   Clean the state. This is an internal helper function. Users should
217      *   always call kfree_skb
218      */
219
220     void __kfree_skb(struct sk_buff *skb)
221     {
222             if (skb->list) {
223                     printk(KERN_WARNING "Warning: kfree_skb passed an skb still "
224                             "on a list (from %p).0, NET_CALLER(skb));
225                     BUG();
226             }
227
228             dst_release(skb->dst);
229     #ifdef CONFIG_XFRM
```

```
230         secpath_put(skb->sp);
231     #endif
232         if(skb->destructor) {
233             if (in_irq())
234                 printk(KERN_WARNING "Warning: kfree_skb on "
235                         "hard IRQ %p0, NET_CALLER(skb));
236             skb->destructor(skb);
237         }
238     #ifdef CONFIG_NETFILTER
239         nf_conntrack_put(skb->nfct);
240     #ifdef CONFIG_BRIDGE_NETFILTER
241         nf_bridge_put(skb->nf_bridge);
242     #endif
243     #endif
244         kfree_skbmem(skb);
245     }
```
——————————————————————————————————————————————————— *[net/core/skbuff.c]*

## 10.6. skb_split

*File :* [net/core/skbuff.c]
This function is in the [ipv4/tcp_output.c] file because it is used only by the TCP code in case it needs to fragment a segment (tcp_fragment).

——————————————————————————————————————————————————— *[net/core/skbuff.c]*
```
355     static void skb_split(struct sk_buff *skb, struct sk_buff *skb1, u32 len)
356     {
357         int i;
358         int pos = skb_headlen(skb);
359
360         if (len < pos) {
361             /* Split line is inside header. */
362             memcpy(skb_put(skb1, pos-len), skb->data + len, pos-len);
363
364             /* And move data appendix as is. */
365             for (i = 0; i < skb_shinfo(skb)->nr_frags; i++)
366                 skb_shinfo(skb1)->frags[i] = skb_shinfo(skb)->frags[i];
367
368             skb_shinfo(skb1)->nr_frags = skb_shinfo(skb)->nr_frags;
369             skb_shinfo(skb)->nr_frags = 0;
370
371             skb1->data_len = skb->data_len;
372             skb1->len += skb1->data_len;
373             skb->data_len = 0;
374             skb->len = len;
375             skb->tail = skb->data+len;
376         } else {
377             int k = 0;
378             int nfrags = skb_shinfo(skb)->nr_frags;
379
380             /* Second chunk has no header, nothing to copy. */
381
382             skb_shinfo(skb)->nr_frags = 0;
```

```
383                 skb1->len = skb1->data_len = skb->len - len;
384                 skb->len = len;
385                 skb->data_len = len - pos;
386
387                 for (i=0; i<nfrags; i++) {
388                     int size = skb_shinfo(skb)->frags[i].size;
389                     if (pos + size > len) {
390                         skb_shinfo(skb1)->frags[k] = skb_shinfo(skb)->frags[i];
391
392                         if (pos < len) {
393                             /* Split frag.
394                              * We have to variants in this case:
395                              * 1. Move all the frag to the second
396                              *    part, if it is possible. F.e.
397                              *    this approach is mandatory for TUX,
398                              *    where splitting is expensive.
399                              * 2. Split is accurately. We make this.
400                              */
401                             get_page(skb_shinfo(skb)->frags[i].page);
402                             skb_shinfo(skb1)->frags[0].page_offset += (len-pos);
403                             skb_shinfo(skb1)->frags[0].size -= (len-pos);
404                             skb_shinfo(skb)->frags[i].size = len-pos;
405                             skb_shinfo(skb)->nr_frags++;
406                         }
407                         k++;
408                     } else {
409                         skb_shinfo(skb)->nr_frags++;
410                     }
411                     pos += size;
412                 }
413                 skb_shinfo(skb1)->nr_frags = k;
414         }
415     }
```

———————————————————————————————————————————————————— *[net/core/skbuff.c]*

## 10.7. skb_checksum

*File :* [net/core/skbuff.c]

There are 3 different functions to perform the checksum. One that simply performs the checksum, one that performs the checksum while copying the data ( it is required by efficiency reasons) and one that takes advantage from the hardware support offered by the card. Partial checksums over contiguous data areas are performed using the csum_partial function, and then added together using the csum_block_add function. skb_checksum computes in this way the checksum for the skbuff also in the case the skbuff is fragmented in an array of pages or a list of skbuffs.

———————————————————————————————————————————————————— *[net/core/skbuff.c]*

```
911     unsigned int skb_checksum(const struct sk_buff *skb, int offset,
912                     int len, unsigned int csum)
913     {
```

```
914           int start = skb_headlen(skb);
915           int i, copy = start - offset;
916           int pos = 0;
917
918           /* Checksum header. */
919           if (copy > 0) {
920                 if (copy > len)
921                       copy = len;
922                 csum = csum_partial(skb->data + offset, copy, csum);
923                 if ((len -= copy) == 0)
924                       return csum;
925                 offset += copy;
926                 pos  = copy;
927           }
928
929           for (i = 0; i < skb_shinfo(skb)->nr_frags; i++) {
930                 int end;
931
932                 BUG_TRAP(start <= offset + len);
933
934                 end = start + skb_shinfo(skb)->frags[i].size;
935                 if ((copy = end - offset) > 0) {
936                       unsigned int csum2;
937                       u8 *vaddr;
938                       skb_frag_t *frag = &skb_shinfo(skb)->frags[i];
939
940                       if (copy > len)
941                             copy = len;
942                       vaddr = kmap_skb_frag(frag);
943                       csum2 = csum_partial(vaddr + frag->page_offset +
944                                       offset - start, copy, 0);
945                       kunmap_skb_frag(vaddr);
946                       csum = csum_block_add(csum, csum2, pos);
947                       if (!(len -= copy))
948                             return csum;
949                       offset += copy;
950                       pos   += copy;
951                 }
952                 start = end;
953           }
954
955           if (skb_shinfo(skb)->frag_list) {
956                 struct sk_buff *list = skb_shinfo(skb)->frag_list;
957
958                 for (; list; list = list->next) {
959                       int end;
960
961                       BUG_TRAP(start <= offset + len);
962
963                       end = start + list->len;
964                       if ((copy = end - offset) > 0) {
965                             unsigned int csum2;
966                             if (copy > len)
967                                   copy = len;
```

```
968                         csum2 = skb_checksum(list, offset - start,
969                                         copy, 0);
970                         csum = csum_block_add(csum, csum2, pos);
971                         if ((len -= copy) == 0)
972                                 return csum;
973                         offset += copy;
974                         pos    += copy;
975                 }
976             start = end;
977         }
978     }
979     if (len)
980             BUG();
981
982     return csum;
983 }
```
———————————————————————————————————————————————————— *[net/core/skbuff.c]*

skb_copy_and_csum_bits performs at the same time the function of the skb_copy_bits and the skb_checksum functions.
In case the device supports checksumming in hardware (CHECKSUM_HW flag), the skb_copy_and_csum_dev function is used.

### 10.8.  Copying functions

### 10.8.1.  skb_copy()

Makes a complete copy of an skbuff header and its data. It eventually converts a nonlinear skbuff to a linear one. The headroom of the old skbuff is computed and reserved also in the new one (the variable for it is inappropriately called headerlen). If the data doesn't need to be modified then the use of pskb_copy that uses the reference count mechanisms for the nonlinear part is recommended.

A linear skbuff capable of storing both the data and the headroom of the old skb is allocated with alloc_skb, this function allocates both the header and a contiguous  data area.  The data area size is equal to the size of the linear part of the old skbuff plus the size of data in the nonlinear parts (external pages and rest of the chain of skbuffs).  If it is not possible to allocate such an skbuff it returns NULL.  A headroom equal to the one in the old skbuff is reserved in the new skbuff calling skb_reserve.  It sets the tail pointer in the new skbuff to accomodate for the total length of the old skbuff(calling skb_put).  It copies the checksum and the ip_summed flag.  Then it calls the skb_copy_bits function to copy everything from the very beginning of the old skbuff (from skb->head not skb->data) to the end of the data. This means that it copies also the headroom. The second argument of skb_copy_bits is the displacement from the standard skb->data pointer from where to begin the copy.  Finally it copies the skbuff header, and returns a pointer to the new skbuff.

_____ *net/core/skbuff.c*

```
391     struct sk_buff *skb_copy(const struct sk_buff *skb, int gfp_mask)
392     {
393          int headerlen = skb->data - skb->head;
394          /*
395           *   Allocate the copy buffer
396           */
397          struct sk_buff *n = alloc_skb(skb->end - skb->head + skb->data_len,
398                           gfp_mask);
399          if (!n)
400              return NULL;
401
402          /* Set the data pointer */
403          skb_reserve(n, headerlen);
404          /* Set the tail pointer and length */
405          skb_put(n, skb->len);
406          n->csum        = skb->csum;
407          n->ip_summed = skb->ip_summed;
408
409          if (skb_copy_bits(skb, -headerlen, n->head, headerlen + skb->len))
410              BUG();
411
412          copy_skb_header(n, skb);
413          return n;
414     }
```

_____ *net/core/skbuff.c*


### 10.8.2.  pskb_copy()

*File :* `net/core/skbuff.c`



It allocates an skbuff having the same linear size of the old one (because of this it doesnt use skb->len, but skb->end - skb->head).  It reserves the same headroom available in the old skbuff moving the data and tail pointers (headroom is the space between the beginning of the linear part of the skbuff and the skb->data pointer).  Then it adjusts the tail pointer to accomodate for the data in the linear part of the skbuff (skb_headlen is the size of such data).  It copies headlen bytes of data from the old skbuff to the new one. Copies checksum, ip_summed flag,data_len (the size of the data not in the linear part), len from old to new. Runs through the array of pages, if it exists, copies fragment descriptors from the old to the new skbuff. For each page it increments the usage count without copying it.  Then if it is an skbuff chain,  it copies the frag_list pointer, and it goes through the list of skbuffs and increases the usage count (skb_clone_fraglist increases the skb->users counter on each of the following skbuffs).  Finally it copies the skb header (copy_skb_header).  And it returns a pointer to the new skbuff head.

_____ *net/core/skbuff.c*

```
430     struct sk_buff *pskb_copy(struct sk_buff *skb, int gfp_mask)
431     {
432             /*
433              *    Allocate the copy buffer
434              */
435             struct sk_buff *n = alloc_skb(skb->end - skb->head, gfp_mask);
436
437             if (!n)
438                     goto out;
439
440             /* Set the data pointer */
441             skb_reserve(n, skb->data - skb->head);
442             /* Set the tail pointer and length */
443             skb_put(n, skb_headlen(skb));
444             /* Copy the bytes */
445             memcpy(n->data, skb->data, n->len);
446             n->csum        = skb->csum;
447             n->ip_summed = skb->ip_summed;
```

```
448
449          n->data_len  = skb->data_len;
450          n->len            = skb->len;
451
452          if (skb_shinfo(skb)->nr_frags) {
453                int i;
454
455                for (i = 0; i < skb_shinfo(skb)->nr_frags; i++) {
456                      skb_shinfo(n)->frags[i] = skb_shinfo(skb)->frags[i];
457                      get_page(skb_shinfo(n)->frags[i].page);
458                }
459                skb_shinfo(n)->nr_frags = i;
460          }
461          skb_shinfo(n)->tso_size = skb_shinfo(skb)->tso_size;
462          skb_shinfo(n)->tso_segs = skb_shinfo(skb)->tso_segs;
463
464          if (skb_shinfo(skb)->frag_list) {
465                skb_shinfo(n)->frag_list = skb_shinfo(skb)->frag_list;
466                skb_clone_fraglist(n);
467          }
468
469          copy_skb_header(n, skb);
470     out:
471          return n;
472     }
```
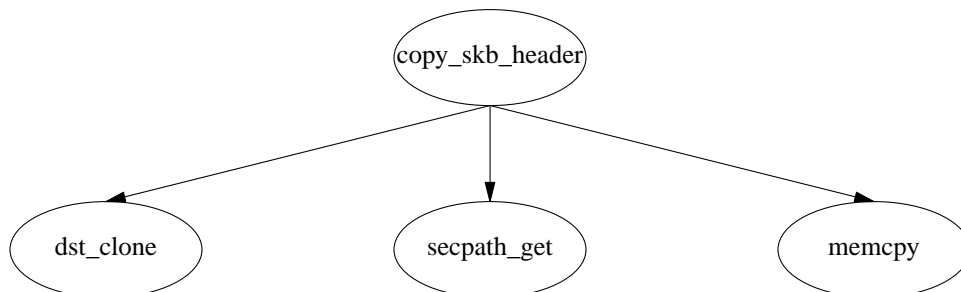
——————————————————————————————————————————————————————————————— *net/core/skbuff.c*

### 10.8.3. copy_skb_header

*File :* `net/core/skbuff.c`



This function supposes that a copy of the data area of the old skb has already being done and the standard pointers of the skbuff are already initialized to the new data area (head, data, tail, end). Most of the remaining fields are just copied over from the old skbuff. But the list pointer (pointer to queue on which the skbuff is linked) is set to NULL, as sk (pointer to the socket the skbuff belongs). As the forwarding mechanism uses a reference count, the dst field cannot be simply copied, but this is done through the dst_clone function that updates the reference counter atomically. The security mechanism also uses a reference count, and so instead of simply copy the sp pointer, it calls the secpath_get function that atomically increases the reference counter associated with the path used. It initializes the pointers to the different layer headers (transport,network,mac .. ) in the new skb to the same relative position as in the old skb. It copies over with memcpy the cb control buffer. The skbuff destructor is set to NULL. Two of the fields used by

netfilter, that use the reference count mechanism, nfct (connection track) and nf_bridge, are copied using functions that atomically increase the reference counters : nf_conntrack_get and nf_bridge_get. Finally it atomically sets the number of users of the new skbuff to 1.

*——————————————————————————————————————————— net/core/skbuff.c*

```
329     static void copy_skb_header(struct sk_buff *new, const struct sk_buff *old)
330     {
331         /*
332          *    Shift between the two data areas in bytes
333          */
334         unsigned long offset = new->data - old->data;
335
336         new->list = NULL;
337         new->sk         = NULL;
338         new->dev  = old->dev;
339         new->real_dev   = old->real_dev;
340         new->priority   = old->priority;
341         new->protocol   = old->protocol;
342         new->dst  = dst_clone(old->dst);
343     #ifdef CONFIG_INET
344         new->sp         = secpath_get(old->sp);
345     #endif
346         new->h.raw = old->h.raw + offset;
347         new->nh.raw     = old->nh.raw + offset;
348         new->mac.raw    = old->mac.raw + offset;
349         memcpy(new->cb, old->cb, sizeof(old->cb));
350         new->local_df   = old->local_df;
351         new->pkt_type   = old->pkt_type;
352         new->stamp = old->stamp;
353         new->destructor = NULL;
354         new->security   = old->security;
355     #ifdef CONFIG_NETFILTER
356         new->nfmark     = old->nfmark;
357         new->nfcache    = old->nfcache;
358         new->nfct = old->nfct;
359         nf_conntrack_get(old->nfct);
360     #ifdef CONFIG_NETFILTER_DEBUG
361         new->nf_debug   = old->nf_debug;
362     #endif
363     #ifdef CONFIG_BRIDGE_NETFILTER
364         new->nf_bridge = old->nf_bridge;
365         nf_bridge_get(old->nf_bridge);
366     #endif
367     #endif
368     #ifdef CONFIG_NET_SCHED
369         new->tc_index   = old->tc_index;
370     #endif
371         atomic_set(&new->users, 1);
372     }
```

*——————————————————————————————————————————— net/core/skbuff.c*
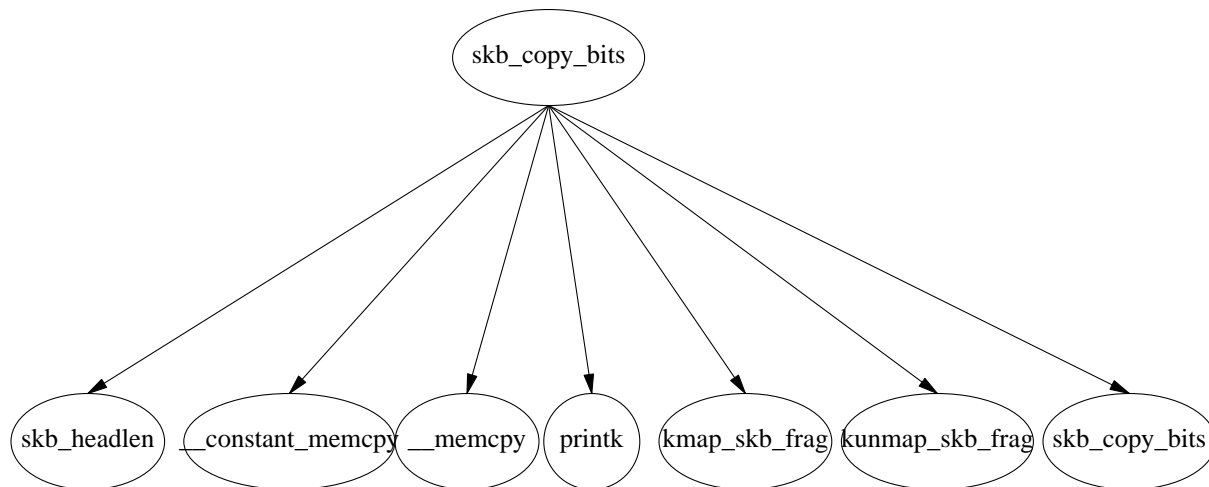
### 10.8.4.  skb_headlen

*File :* `[include/linux/skbuff.h]`



This function returns the number of data bytes in the linear part of the first skbuff.  It is equal to (skb->tail - skb->data).

```
static inline unsigned int skb_headlen(const struct sk_buff *skb)
{
        return skb->len - skb->data_len;
}
```

### 10.8.5.  skb_copy_bits

*File :* `[net/core/skbuff.c]`



This function copies bytes of data from an skbuff to another area of memory and doing so it eventually linearize nonlinear skbuffs.  The first argument is a pointer to an skbuff header from which data area the data should be copied, the third argument is a pointer to an area of memory where the data should be copied to. The offset argument is the quantity that is added  to the skb->data pointer to obtain the address from which the copy will start: skb->data + offset, it is the number of bytes to skip from data pointer position.  But it can be also negative and in that case it represents the bytes to copy from the headroom in addition to those starting from the data pointer.  The len argument is the number of bytes that should be copied.  This function is able to treat paged skbuffs and chains of  skbuffs.  Initially it copies from the linear part of the first skbuff.  It copies copy = headlen - offset bytes.  It has code to copy all the external pages with memcpy.  If a page is in high memory, then it is necessary to find a free slot in the low memory area dedicated to mapping, where to map it (kmap_skb_frag), the page can then be copied and unmapped again (kunmap_skb_frag).  If it is an skbuff chain, then the following skbuffs are copied calling itself recursively for each skbuff in the chain (frag_list).

_____ *[net/core/skbuff.c]*

```
  832     /* Copy some data bits from skb to kernel buffer. */
```

```
833
834     int skb_copy_bits(const struct sk_buff *skb, int offset, void *to, int len)
835     {
836           int i, copy;
837           int start = skb_headlen(skb);
838
839           if (offset > (int)skb->len - len)
840                 goto fault;
841
842           /* Copy header. */
843           if ((copy = start - offset) > 0) {
844                 if (copy > len)
845                       copy = len;
846                 memcpy(to, skb->data + offset, copy);
847                 if ((len -= copy) == 0)
848                       return 0;
849                 offset += copy;
850                 to     += copy;
851           }
852
853           for (i = 0; i < skb_shinfo(skb)->nr_frags; i++) {
854                 int end;
855
856                 BUG_TRAP(start <= offset + len);
857
858                 end = start + skb_shinfo(skb)->frags[i].size;
859                 if ((copy = end - offset) > 0) {
860                       u8 *vaddr;
861
862                       if (copy > len)
863                             copy = len;
864
865                       vaddr = kmap_skb_frag(&skb_shinfo(skb)->frags[i]);
866                       memcpy(to,
867                             vaddr + skb_shinfo(skb)->frags[i].page_offset+
868                             offset - start, copy);
869                       kunmap_skb_frag(vaddr);
870
871                       if ((len -= copy) == 0)
872                             return 0;
873                       offset += copy;
874                       to     += copy;
875                 }
876                 start = end;
877           }
878
879           if (skb_shinfo(skb)->frag_list) {
880                 struct sk_buff *list = skb_shinfo(skb)->frag_list;
881
882                 for (; list; list = list->next) {
883                       int end;
884
885                       BUG_TRAP(start <= offset + len);
886
```

```
887                         end = start + list->len;
888                         if ((copy = end - offset) > 0) {
889                                 if (copy > len)
890                                         copy = len;
891                                 if (skb_copy_bits(list, offset - start,
892                                                 to, copy))
893                                         goto fault;
894                                 if ((len -= copy) == 0)
895                                         return 0;
896                                 offset += copy;
897                                 to     += copy;
898                         }
899                         start = end;
900                 }
901         }
902         if (!len)
903                 return 0;
904
905 fault:
906         return -EFAULT;
907 }
908
```

*———————————————————————————————————————————————— [net/core/skbuff.c]*

### 10.8.6. skb_copy_expand

*File :* [`net/core/skbuff.c`]
It makes a copy of a skbuff and its data expanding the free data areas at the beginning and at the end of it.
A new skbuff is allocated with `alloc_skb` of size equal to the size of the data in the old skbuff plus the required `newheadroom` and `newtailroom`. Then the pointers of the skbuff are displaced with `skb_reserve` of the requested new headroom, and the skbuff and its data are copied over with `skb_copy_bits` (if the skbuff was fragmented it is linearized). The TSO info are copied too.

*———————————————————————————————————————————————— [net/core/skbuff.c]*

```
578     struct sk_buff *skb_copy_expand(const struct sk_buff *skb,
579                         int newheadroom, int newtailroom, int gfp_mask)
580     {
581         /*
582          *   Allocate the copy buffer
583          */
584         struct sk_buff *n = alloc_skb(newheadroom + skb->len + newtailroom,
585                             gfp_mask);
586         int head_copy_len, head_copy_off;
587
588         if (!n)
589             return NULL;
590
591         skb_reserve(n, newheadroom);
592
593         /* Set the tail pointer and length */
594         skb_put(n, skb->len);
595
596         head_copy_len = skb_headroom(skb);
```

```
597            head_copy_off = 0;
598            if (newheadroom <= head_copy_len)
599                    head_copy_len = newheadroom;
600            else
601                    head_copy_off = newheadroom - head_copy_len;
602
603            /* Copy the linear header and data. */
604            if (skb_copy_bits(skb, -head_copy_len, n->head + head_copy_off,
605                          skb->len + head_copy_len))
606                    BUG();
607
608            copy_skb_header(n, skb);
609            skb_shinfo(n)->tso_size = skb_shinfo(skb)->tso_size;
610            skb_shinfo(n)->tso_segs = skb_shinfo(skb)->tso_segs;
611
612            return n;
613    }
```
———————————————————————————————————————————————— *[net/core/skbuff.c]*


### 10.8.7.  skb_copy_datagram, skb_copy_datagram_iovec

*File :* [net/core/datagram.c]
The first function is simply a call to the second after having set up an iovec with only one entry.  The second copies a datagram to an existing iovec: if the skbuff is fragmented (in a list out of IP defragmentation or in an array of pages) the data is linearized (as it can be an iovec .. it can be also splitted !).  This function is called in tcp_input.c and also in udp.c and raw.c to deliver the data to the user.

———————————————————————————————————————————————— *[net/core/datagram.c]*
```
202    /*
203     *   Copy a datagram to a linear buffer.
204     */
205    int skb_copy_datagram(const struct sk_buff *skb, int offset, char *to, int
size)
206    {
207            struct iovec iov = {
208                    .iov_base = to,
209                    .iov_len =size,
210            };
211
212            return skb_copy_datagram_iovec(skb, offset, &iov, size);
213    }
214
215    /**
216     *   skb_copy_datagram_iovec - Copy a datagram to an iovec.
217     *   @skb - buffer to copy
218     *   @offset - offset in the buffer to start copying from
219     *   @iovec - io vector to copy to
220     *   @len - amount of data to copy from buffer to iovec
221     *
222     *   Note: the iovec is modified during the copy.
223     */
224    int skb_copy_datagram_iovec(const struct sk_buff *skb, int offset,
```

```
225                        struct iovec *to, int len)
226     {
227          int start = skb_headlen(skb);
228          int i, copy = start - offset;
229
230          /* Copy header. */
231          if (copy > 0) {
232               if (copy > len)
233                    copy = len;
234               if (memcpy_toiovec(to, skb->data + offset, copy))
235                    goto fault;
236               if ((len -= copy) == 0)
237                    return 0;
238               offset += copy;
239          }
240
241          /* Copy paged appendix. Hmm... why does this look so complicated? */
242          for (i = 0; i < skb_shinfo(skb)->nr_frags; i++) {
243               int end;
244
245               BUG_TRAP(start <= offset + len);
246
247               end = start + skb_shinfo(skb)->frags[i].size;
248               if ((copy = end - offset) > 0) {
249                    int err;
250                    u8  *vaddr;
251                    skb_frag_t *frag = &skb_shinfo(skb)->frags[i];
252                    struct page *page = frag->page;
253
254                    if (copy > len)
255                         copy = len;
256                    vaddr = kmap(page);
257                    err = memcpy_toiovec(to, vaddr + frag->page_offset +
258                                    offset - start, copy);
259                    kunmap(page);
260                    if (err)
261                         goto fault;
262                    if (!(len -= copy))
263                         return 0;
264                    offset += copy;
265               }
266               start = end;
267          }
268
269          if (skb_shinfo(skb)->frag_list) {
270               struct sk_buff *list = skb_shinfo(skb)->frag_list;
271
272               for (; list; list = list->next) {
273                    int end;
274
275                    BUG_TRAP(start <= offset + len);
276
277                    end = start + list->len;
278                    if ((copy = end - offset) > 0) {
```

```
279                          if (copy > len)
280                                copy = len;
281                          if (skb_copy_datagram_iovec(list,
282                                              offset - start,
283                                              to, copy))
284                                goto fault;
285                          if ((len -= copy) == 0)
286                                return 0;
287                          offset += copy;
288                     }
289                  start = end;
290            }
291        }
292        if (!len)
293              return 0;
294
295    fault:
296        return -EFAULT;
297    }
```
———————————————————————————————————————————————— *[net/core/datagram.c]*


## 10.9. Cloning and sharing

### 10.9.1. clone_fraglist()

*File :* `net/core/skbuff.c`



clone_fraglist

This function traverse the list of skbuffs and invokes skb_get for each of them. Skb_get simply increments the number of users of the skbuff. So this clone function works through a copy-on-write mechanism : nothing is really copied now, it is the responsability of those who wants to write on the skbuffs to copy them. This can save some not needed copies.

```
176
177    static void skb_clone_fraglist(struct sk_buff *skb)
178    {
179        struct sk_buff *list;
180
181        for (list = skb_shinfo(skb)->frag_list; list; list = list->next)
182            skb_get(list);
183    }
```

### 10.9.2. skb_shared, skb_clone, skb_share_check, skb_unshare

skb_shared atomically tests if the skbuff is shared, that is there is more than 1 user of the same header and data.

———————————————————————————————————————————————— *net/core/skbuff.c*

```
377     static inline int skb_shared(const struct sk_buff *skb)
378     {
379             return atomic_read(&skb->users) != 1;
380     }
```
——————————————————————————————————————————————————— *net/core/skbuff.c*

It is used to panic if a shared skbuff is passed to functions that do not expect it or in the case a user wants to modify the header, to check if cloning is necessary. skb_cloned checks if an skbuff has been cloned (multiple headers have pointers to the same data area).

——————————————————————————————————————————————————— *net/core/skbuff.c*
```
365     static inline int skb_cloned(const struct sk_buff *skb)
366     {
367             return skb->cloned && atomic_read(&skb_shinfo(skb)->dataref) != 1;
368     }
```
——————————————————————————————————————————————————— *net/core/skbuff.c*

skb_share_check checks if an skbuff is shared, in that case it is cloned, the users counter is decremented on the old skbuff and a pointer to the new is returned. The new skbuff has skb->users == 1 and skb->cloned == 1, the old skbuff has skb->cloned == 1 and skb->users decremented by one.

——————————————————————————————————————————————————— *net/core/skbuff.c*
```
395     static inline struct sk_buff *skb_share_check(struct sk_buff *skb, int pri)
396     {
397             might_sleep_if(pri & __GFP_WAIT);
398             if (skb_shared(skb)) {
399                     struct sk_buff *nskb = skb_clone(skb, pri);
400                     kfree_skb(skb);
401                     skb = nskb;
402             }
403             return skb;
404     }
```
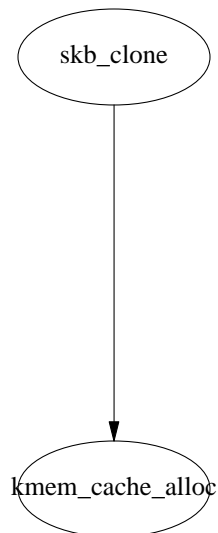——————————————————————————————————————————————————— *net/core/skbuff.c*

skb_unshare checks if an skbuff is cloned. In such a case makes a complete copy of the skbuff (header and data), decrements skb->users and returns a pointer to the new skbuff.

### 10.9.3. skb_clone

*File :* `net/core/skbuff.c`

```
                    ╭───────────────╮
                    │   skb_clone   │
                    ╰───────┬───────╯
                            │
                            │
                            │
                            ▼
                ╭─────────────────────╮
                │  kmem_cache_alloc   │
                ╰─────────────────────╯
```

You clone an skbuff allocating a new skbuff header from the slab and initilizing its fields from the values of the old one. The data area is not copied, it is shared, so you increment the counter skb_shared_info->data_ref in the shared info area. The number of users of the new header (n->users) is set to 1, and you put the cloned flag in the old and new header to 1. The pointers of the doubly linked list to which the skbuff can be linked are initialized to NULL in the new header. And also the destructor in the new header is initialized to NULL.

_____ *net/core/skbuff.c*

```
246
247    /**
248     *    skb_clone  -    duplicate an sk_buff
249     *    @skb: buffer to clone
250     *    @gfp_mask: allocation priority
251     *
252     *    Duplicate an &sk_buff. The new one is not owned by a socket. Both
253     *    copies share the same packet data but not structure. The new
254     *    buffer has a reference count of 1. If the allocation fails the
255     *    function returns %NULL otherwise the new buffer is returned.
256     *
257     *    If this function is called from an interrupt gfp_mask() must be
258     *    %GFP_ATOMIC.
259     */
260
261    struct sk_buff *skb_clone(struct sk_buff *skb, int gfp_mask)
262    {
263            struct sk_buff *n = kmem_cache_alloc(skbuff_head_cache, gfp_mask);
264
265            if (!n)
266                    return NULL;
267
268    #define C(x) n->x = skb->x
269
270            n->next = n->prev = NULL;
271            n->list = NULL;
272            n->sk = NULL;
273            C(stamp);
```

```
274          C(dev);
275          C(real_dev);
276          C(h);
277          C(nh);
278          C(mac);
279          C(dst);
280          dst_clone(skb->dst);
281          C(sp);
282     #ifdef CONFIG_INET
283          secpath_get(skb->sp);
284     #endif
285          memcpy(n->cb, skb->cb, sizeof(skb->cb));
286          C(len);
287          C(data_len);
288          C(csum);
289          C(local_df);
290          n->cloned = 1;
291          C(pkt_type);
292          C(ip_summed);
293          C(priority);
294          C(protocol);
295          C(security);
296          n->destructor = NULL;
297     #ifdef CONFIG_NETFILTER
298          C(nfmark);
299          C(nfcache);
300          C(nfct);
301          nf_conntrack_get(skb->nfct);
302     #ifdef CONFIG_NETFILTER_DEBUG
303          C(nf_debug);
304     #endif
305     #ifdef CONFIG_BRIDGE_NETFILTER
306          C(nf_bridge);
307          nf_bridge_get(skb->nf_bridge);
308     #endif
309     #endif /*CONFIG_NETFILTER*/
310     #if defined(CONFIG_HIPPI)
311          C(private);
312     #endif
313     #ifdef CONFIG_NET_SCHED
314          C(tc_index);
315     #endif
316          C(truesize);
317          atomic_set(&n->users, 1);
318          C(head);
319          C(data);
320          C(tail);
321          C(end);
322
323          atomic_inc(&(skb_shinfo(skb)->dataref));
324          skb->cloned = 1;
325
326          return n;
327     }
```

### 10.9.4.  skb_get()

*File :* `include/linux/skbuff.h`
It atomically increments the number of users of the skbuff returning back a pointer to it. The complete skbuff (header and data area) is shared in this way and if it needs to be modified, then a copy should be done.

```
319     /**
320      *   skb_get - reference buffer
321      *   @skb: buffer to reference
322      *
323      *   Makes another reference to a socket buffer and returns a pointer
324      *   to the buffer.
325      */
326     static inline struct sk_buff *skb_get(struct sk_buff *skb)
327     {
328             atomic_inc(&skb->users);
329             return skb;
330     }
331
```

### 10.9.5.  Cloning and udp multicasting/broadcasting

Udp multicast/broadcast data have to be delivered to each listener. Each listener has a separate receive queue, so an skbuff containing a multicast/broadcast udp packet has to be cloned for each listener, and enqueued on its own socket receive queue (the enqueuing of the skbuff makes necessary the adjustment of their next and prev pointers and so each listener needs its unique copy of the header). The udp_v4_mcast_deliver function runs through the list of listeners (using udp_v4_mcast_next) and it enqueues a clone of the skbuff in the socket receive queue of it (udp_queue_rcv_skb).

```
1082    static int udp_v4_mcast_deliver(struct sk_buff *skb, struct udphdr *uh,
1083                               u32 saddr, u32 daddr)
1084    {
1085         struct sock *sk;
1086         int dif;
1087
1088         read_lock(&udp_hash_lock);
1089         sk = sk_head(&udp_hash[ntohs(uh->dest) & (UDP_HTABLE_SIZE - 1)]);
1090         dif = skb->dev->ifindex;
1091         sk = udp_v4_mcast_next(sk, uh->dest, daddr, uh->source, saddr, dif);
1092         if (sk) {
1093              struct sock *sknext = NULL;
1094
1095              do {
1096                   struct sk_buff *skb1 = skb;
1097
1098                   sknext = udp_v4_mcast_next(sk_next(sk), uh->dest, daddr,
1099                                       uh->source, saddr, dif);
1100                   if(sknext)
1101                        skb1 = skb_clone(skb, GFP_ATOMIC);
```

```
1102
1103                      if(skb1) {
1104                            int ret = udp_queue_rcv_skb(sk, skb1);
1105                            if (ret > 0)
1106                                  /* we should probably re-process instead
1107                                   * of dropping packets here. */
1108                                  kfree_skb(skb1);
1109                      }
1110                      sk = sknext;
1111               } while(sknext);
1112         } else
1113               kfree_skb(skb);
1114         read_unlock(&udp_hash_lock);
1115         return 0;
1116    }
```
————————————————————————————————————————————————— *include/linux/skbuff.h*

### 10.9.6. Sharing of fragment list

When you copy a list fragmented skbuff you copy the header and you just increase the users counter on the following skbuffs, because usually you dont need to perform any operation on the remaining skbuffs ( you add/strip headers only on the first skbuff). This seems to be the only sensible use of skbuff sharing.

### 10.10. Miscellaneous functions

### 10.10.1. pskb_expand_head

*File :* `include/linux/skbuff.h`



This function allocates a new linear skbuff data area with enough space to provide the specified headroom and tailroom. Then it copies all the data from the old skbuff to this one, eventually reducing a fragmented skbuff to a linear one. The skbuff header remains the same, just the pointers to the data area are changed. The other skbuffs in the skbuff chain are shared through calls to skb_get that increments the dataref field.

_____ *include/linux/skbuff.h*

```
474     /**
475      *    pskb_expand_head - reallocate header of &sk_buff
476      *    @skb: buffer to reallocate
477      *    @nhead: room to add at head
478      *    @ntail: room to add at tail
479      *    @gfp_mask: allocation priority
480      *
481      *    Expands (or creates identical copy, if &nhead and &ntail are zero)
482      *    header of skb. &sk_buff itself is not changed. &sk_buff MUST have
483      *    reference count of 1. Returns zero in the case of success or error,
484      *    if expansion failed. In the last case, &sk_buff is not changed.
485      *
486      *    All the pointers pointing into skb header may change and must be
487      *    reloaded after call to this function.
488      */
489
490     int pskb_expand_head(struct sk_buff *skb, int nhead, int ntail, int gfp_mask)
491     {
492          int i;
493          u8 *data;
494          int size = nhead + (skb->end - skb->head) + ntail;
495          long off;
496
497          if (skb_shared(skb))
498               BUG();
499
500          size = SKB_DATA_ALIGN(size);
501
502          data = kmalloc(size + sizeof(struct skb_shared_info), gfp_mask);
503          if (!data)
504               goto nodata;
505
506          /* Copy only real data... and, alas, header. This should be
507           * optimized for the cases when header is void. */
508          memcpy(data + nhead, skb->head, skb->tail - skb->head);
509          memcpy(data + size, skb->end, sizeof(struct skb_shared_info));
510
511          for (i = 0; i < skb_shinfo(skb)->nr_frags; i++)
512               get_page(skb_shinfo(skb)->frags[i].page);
513
514          if (skb_shinfo(skb)->frag_list)
515               skb_clone_fraglist(skb);
516
517          skb_release_data(skb);
518
519          off = (data + nhead) - skb->head;
520
521          skb->head    = data;
522          skb->end     = data + size;
523          skb->data   += off;
524          skb->tail   += off;
525          skb->mac.raw += off;
526          skb->h.raw   += off;
```

```
527          skb->nh.raw  += off;
528          skb->cloned  = 0;
529          atomic_set(&skb_shinfo(skb)->dataref, 1);
530          return 0;
531
532   nodata:
533          return -ENOMEM;
534      }
```

*include/linux/skbuff.h*


### 10.10.2. skb_realloc_headroom

*File :* `include/linux/skbuff.h`
It makes a private copy of the skb checking that there is at least `headroom` bytes of free space at the beginning of the buffer. The copy is made with `pskb_copy` if there is already enough space, otherwise it is done first by cloning the skbuff with `skb_clone` and then copying expanding the header with `pskb_expand_head`.

*include/linux/skbuff.h*

```
536     /* Make private copy of skb with writable head and some headroom */
537
538     struct sk_buff *skb_realloc_headroom(struct sk_buff *skb, unsigned int head-
room)
539     {
540          struct sk_buff *skb2;
541          int delta = headroom - skb_headroom(skb);
542
543          if (delta <= 0)
544               skb2 = pskb_copy(skb, GFP_ATOMIC);
545          else {
546               skb2 = skb_clone(skb, GFP_ATOMIC);
547               if (skb2 && pskb_expand_head(skb2, SKB_DATA_ALIGN(delta), 0,
548                                    GFP_ATOMIC)) {
549                    kfree_skb(skb2);
550                    skb2 = NULL;
551               }
552          }
553          return skb2;
554     }
555
```

*include/linux/skbuff.h*


### 10.10.3. skb_pad

*File :* `include/linux/skbuff.h`
In some cases it is possible that network cards using DMA transfer data beyond the buffer end because of size of the unit of transfer. It is therefore necessary sometimes to fill some space after the buffer with zeroes. This function can be expensive because if there is not enough tailroom in the skbuff it has to copy and expand it. This function is called to be assured there are `pad` bytes of zeroes after the data. This function is only called by the skb_padto function.

─────────────────────────────────────────────── *include/linux/skbuff.h*

```
627     struct sk_buff *skb_pad(struct sk_buff *skb, int pad)
628     {
629            struct sk_buff *nskb;
630
631            /* If the skbuff is non linear tailroom is always zero.. */
632            if (skb_tailroom(skb) >= pad) {
633                  memset(skb->data+skb->len, 0, pad);
634                  return skb;
635            }
636
637            nskb = skb_copy_expand(skb, skb_headroom(skb), skb_tailroom(skb) + pad,
GFP_ATOMIC);
638            kfree_skb(skb);
639            if (nskb)
640                  memset(nskb->data+nskb->len, 0, pad);
641            return nskb;
642     }
```

─────────────────────────────────────────────── *include/linux/skbuff.h*


### 10.10.4.  skb_padto

*File :* `include/linux/skbuff.h`
This function uses the `skb_pad` function to pad the data in the skbuff with zeroes up to `len` bytes.
It is used to pad ethernet frames to the minimum size prescribed by the standard. 60 bytes excluded the
frame check sequence (FCS) that is computed by the card (4 bytes).  Inside ethernet drivers it is normally
used this way : skb_padto(skb,skb= where ETH_ZLEN is 60.

─────────────────────────────────────────────── *include/linux/skbuff.h*

```
1118    static inline struct sk_buff *skb_padto(struct sk_buff *skb, unsigned int len)
1119    {
1120           unsigned int size = skb->len;
1121           if (likely(size >= len))
1122                   return skb;
1123           return skb_pad(skb, len-size);
1124    }
```

─────────────────────────────────────────────── *include/linux/skbuff.h*


### 10.10.5.  skb_set_owner_r,skb_set_owner_w

These two function are called when the owner socket of an skbuff is known. During output the owner
socket is known since the beginning and thus they are called inside the allocation functions
(sock_alloc_skb,sock_alloc_pskb ..). In input packets need to be processed to discover the owner socket.

─────────────────────────────────────────────── *include/net/sock.h*

```
886     static inline void skb_set_owner_w(struct sk_buff *skb, struct sock *sk)
887     {
888            sock_hold(sk);
889            skb->sk = sk;
890            skb->destructor = sock_wfree;
```

```
891          atomic_add(skb->truesize, &sk->sk_wmem_alloc);
892      }
893
894      static inline void skb_set_owner_r(struct sk_buff *skb, struct sock *sk)
895      {
896          skb->sk = sk;
897          skb->destructor = sock_rfree;
898          atomic_add(skb->truesize, &sk->sk_rmem_alloc);
899      }
```
*—————————————————————————————————————————— include/net/sock.h*


### 10.10.6. get_page

This function is part of the memory management functions. It simply atomically increases the reference
counter on the page, relying on the standard COW (copy on write) kernel mechanism to copy the page
when it is needed. If the kernel was compiled with huge page support (CONFIG_HUGETLB_PAGE) then
this function is defined with a further check. If the page is a large/huge page (named compound in the code
because it is seen also as a set of small pages), then the lru.next field points to the head (the first small page
of the large page). It is the refcount of that head page that is incremented in that case.

*—————————————————————————————————————————— include/linux/mm.h*
```
240      static inline void get_page(struct page *page)
241      {
242          if (PageCompound(page))
243              page = (struct page *)page->lru.next;
244          atomic_inc(&page->count);
245      }
```
*—————————————————————————————————————————— include/linux/mm.h*


---------------------- comment : it is better to call descriptor the fixed size part of the skbuff that keeps point-
ers to data areas and control info (the sk_buff). and skbuff list head the sk_buff_head structure.
--------------------------------------- Since the paper of Cox[2] many things have changed in the skbuff system.
The two most important changes are : - now skbuffs can be nonlinear and stored  in two different ways :
   - in fragments as an array of pages using the skb_shared_info.frags
     array (now it can have enough pages for 64KB + 2pages)
   - in fragments as a list of sk_buff using the skb_shared_info.frag_list
     pointer - the interface with the kernel memory allocator has been completely
 changed and now you can use named slabs for structures that
 can get benefits from caching freed objects of the same kind


--------------------------------------------


### 10.11. IRIX

With some devices a bus adapter is interposed between the device and the system bus, in this case the bus
adapter can translate between bus addresses used by the devices (for instance PCI cards on a PCI bus) and
system addresses. 64 bits addresses are divided accroding to their most significant 2 bits into 4 segments :
3 of them access memory throught the TLB and so are mapped (user process space xkuseg=00, supervisor
mode space xksseg=01, kernel virtual space xkseg=00)

and one accesses memory through physical addresses (xkphys=10).
Using dma maps with pci devices

Allocate a map (on the pci address space): pccio_dmamap_alloc(device identifier:bus and slot, device descriptor,bytes,falgs)

( You can get a default device descriptor calling device_desc_default_get with the device identifier) You activate the map calling pciio_dmamap_addr() or pciio_dmamap_list(). These operations set up a translation table on the bus adapter that converts bus to/from system addresses.

3, 4, 5, 6, 7, 8, 9, 10, 11 # end

# Table of Contents

## References

1.    K.C.Knowlton, "A Fast Storage Allocator," *Communications of ACM,* Vol.8, October 1965 (1965).

2.    Alan Cox, "Network Buffers and Memory Management," *Linux Journal,* 29 (September 29, 1996).

3.    Richard Stevens, *TCP/IP Illustrated vol 1, The protocols,* p. Addison Wesley (1994).

4.    G. Wright, R. Stevens, *TCP/IP Illustrated vol. 2, The implementation,* p. Addison Wesley (1995).

5.    Daniel P. Bovet, Marco Cesati, *Understanding the Linux kernel, 2nd edition,* p. O'Reilly (December 2002).

6.    Marshall Kirk McKusick, Keith Bostic , Michael J. Karels, John S. Quarterman, *The Design and Implementation of the 4.4 BSD Operating System, 2nd edition,* p. Addison Wesley (1996).

7.    Mel Gorman, *Understanding the Linux Virtual Memory Manager,* p. available on the web (July, 2003).

8.    Jeff Bonwick (Sun Microsystems), "The Slab Allocator: An Object-Caching Kernel Memory Allocator," *USENIX Summer Technical Conference* (1994).

9.    Brad Fitzgibbons (gatech), *The Linux Slab Allocator,* p. on the web (October 30, 2000).

10.   Mika Karlstedt (University of Helsinki), *Linux 2.4 memory management,* p. on the web (2002).

11.   B.Goodheart,J.Cox,J.R.Mashey, *The Magic Garden Explained: The internals of Unix System V Release 4,* p. Prentice Hall (1994).