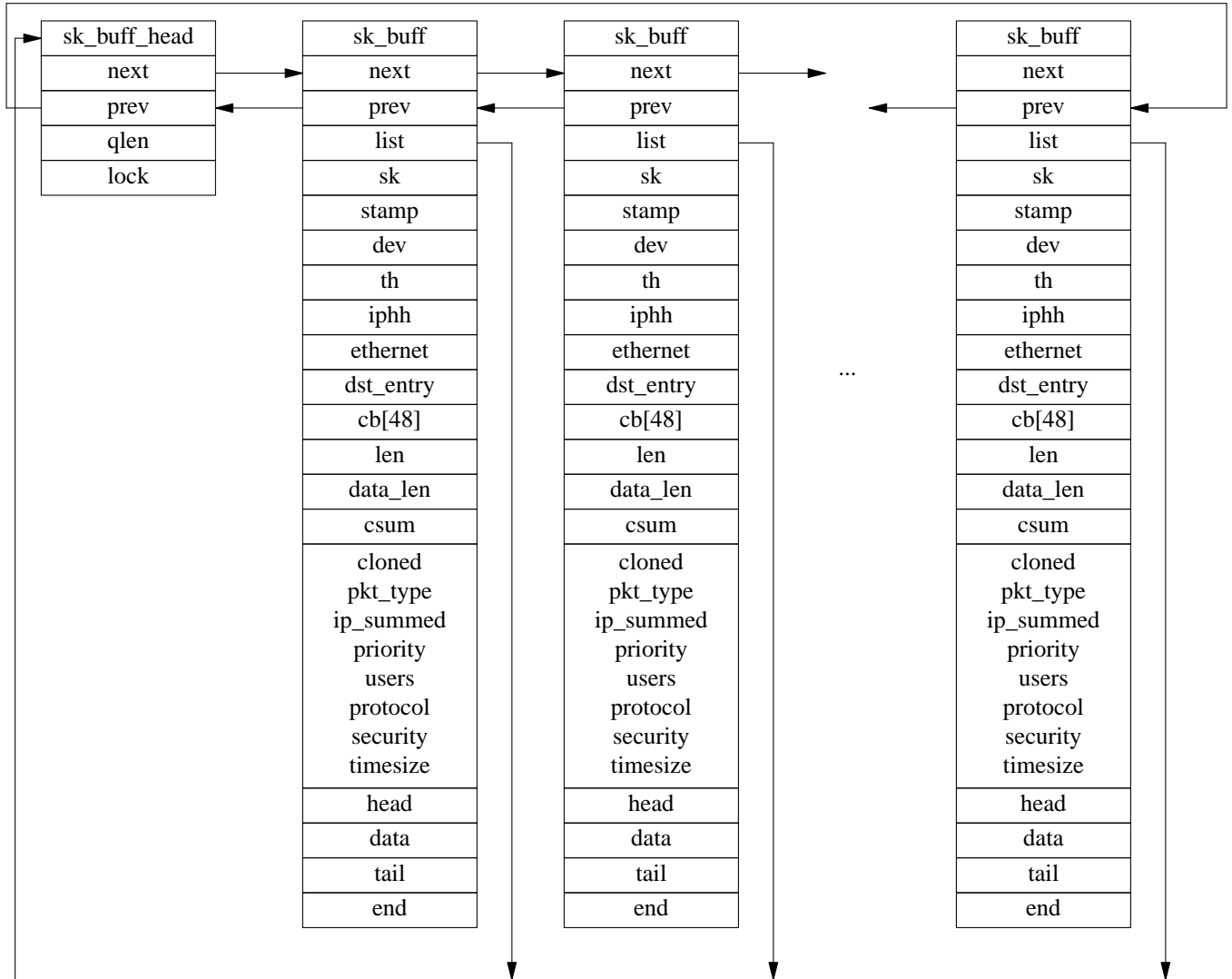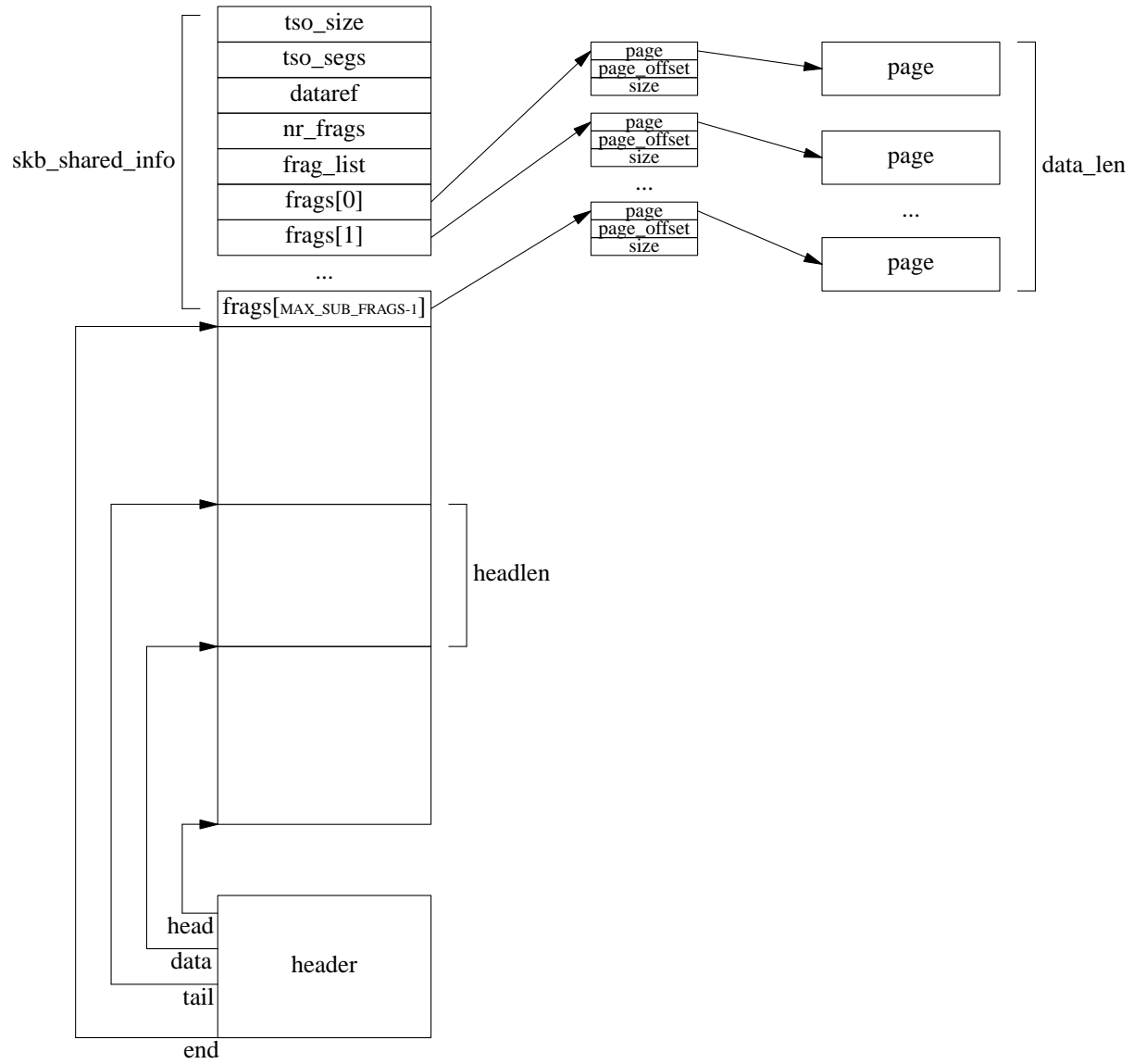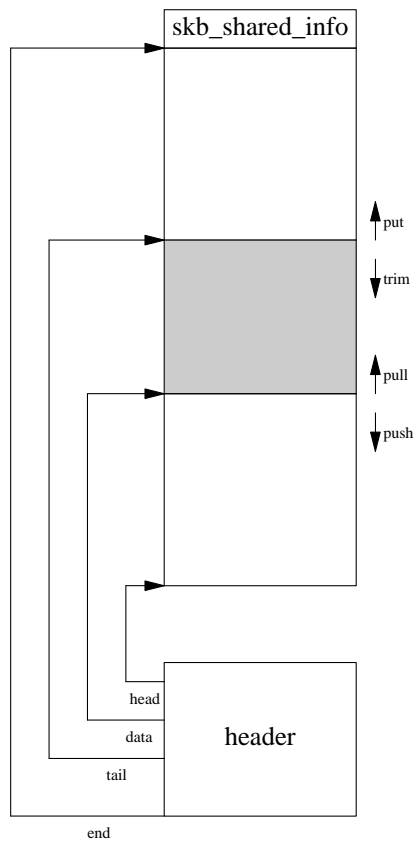# Skbuffs - A tutorial

*ri-oa*

sissa

## 1. Introduction

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 kernel 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.

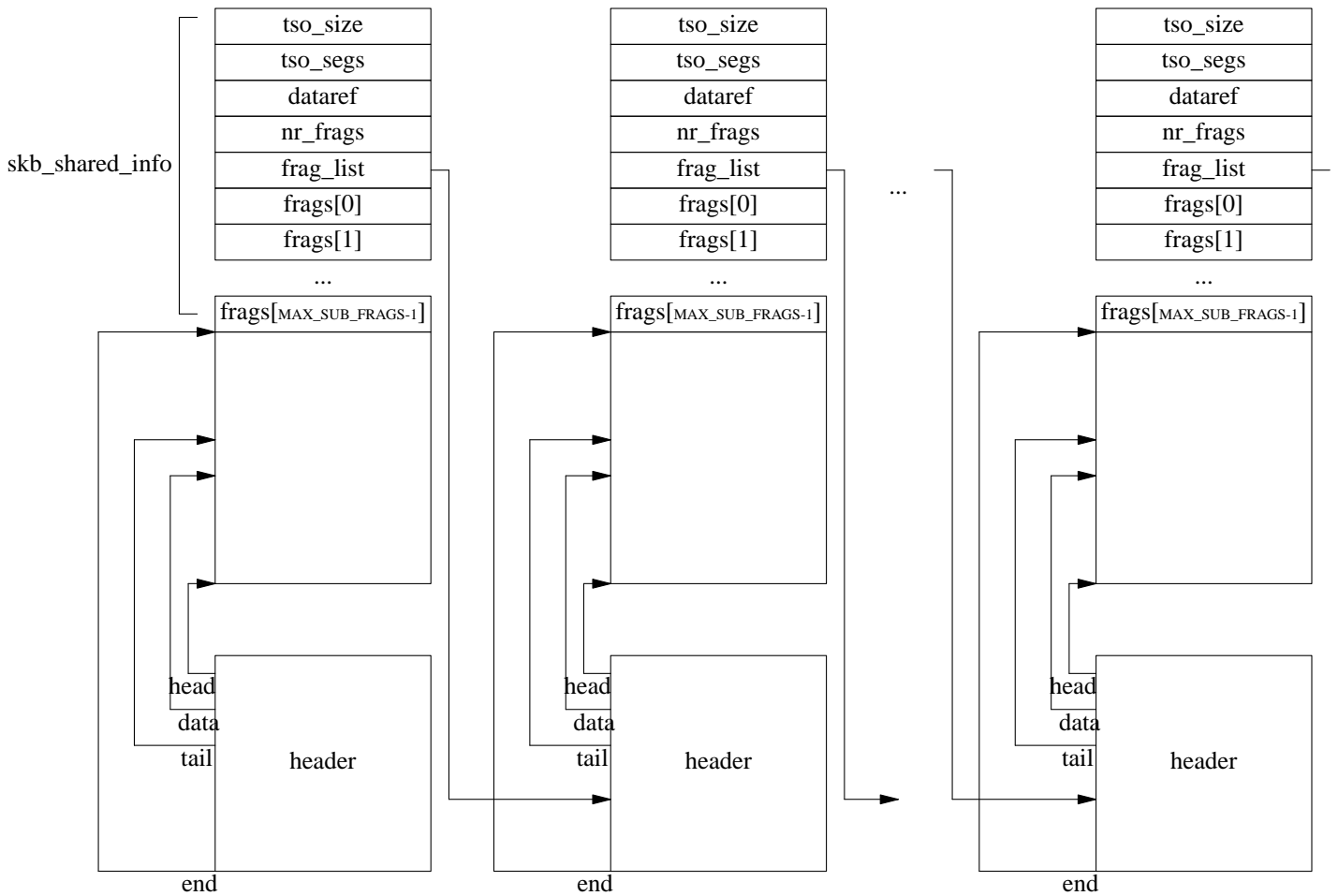| sk_buff_head | | sk_buff | | sk_buff | | sk_buff |
| --- | --- | --- | --- | --- | --- | --- |
| next | | next | | next | | next |
| prev | | prev | | prev | | prev |
| qlen | | list | | list | | list |
| lock | | sk | | sk | | sk |
| | | stamp | | stamp | | stamp |
| | | dev | | dev | | dev |
| | | th | | th | | th |
| | | iphh | | iphh | | iphh |
| | | ethernet | | ethernet | | ethernet |
| | | dst_entry | | dst_entry | | dst_entry |
| | | cb[48] | | cb[48] | | cb[48] |
| | | len | | len | | len |
| | | data_len | | data_len | | data_len |
| | | csum | | csum | | csum |
| | | cloned<br>pkt_type<br>ip_summed<br>priority<br>users<br>protocol<br>security<br>timesize | | cloned<br>pkt_type<br>ip_summed<br>priority<br>users<br>protocol<br>security<br>timesize | | cloned<br>pkt_type<br>ip_summed<br>priority<br>users<br>protocol<br>security<br>timesize |
| | | head | | head | | head |
| | | data | | data | | data |
| | | tail | | tail | | tail |
| | | end | | end | | end |

...

skb_shared_info

| tso_size |
| tso_segs |
| dataref |
| nr_frags |
| frag_list |
| frags[0] |
| frags[1] |
| ... |
| frags[MAX_SUB_FRAGS-1] |

| page |
| page_offset |
| size |

page

| page |
| page_offset |
| size |

page

...

| page |
| page_offset |
| size |

...

page

data_len

headlen

| head |
| data |
| tail |
| end |

header

skb_shared_info

↕

↕

head

data

header

tail

end

skb_shared_info

put

trim

pull

push

head

data

header

tail

end

The skbuff system uses in two ways the slab allocator.[1] The skbuff heads are allocated from a named slab called `skbuff_head_cache`.

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

```
void __init skb_init(void)
{
        skbuff_head_cache = kmem_cache_create("skbuff_head_cache",
                                      sizeof(struct sk_buff),
                                      0,
                                      SLAB_HWCACHE_ALIGN,
                                      NULL, NULL);
        if (!skbuff_head_cache)
                panic("cannot create skbuff cache");
}
```

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


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

```
struct sk_buff *alloc_skb(unsigned int size, int gfp_mask)
{
        struct sk_buff *skb;
        u8 *data;

        /* Get the HEAD */
        skb = kmem_cache_alloc(skbuff_head_cache,
                               gfp_mask & ~__GFP_DMA);
        if (!skb)
                goto out;

        /* Get the DATA. Size must match skb_add_mtu(). */
        size = SKB_DATA_ALIGN(size);
        data = kmalloc(size + sizeof(struct skb_shared_info), gfp_mask);
        if (!data)
                goto nodata;
```

---------------------- comment : i think is better to call descriptor the fixed size part of the skbuff that keeps pointers to data areas and control info (the sk_buff).  and skbuff list head the sk_buff_head structure. --------------------------------------- Since the paper 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


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


## 2.  Constrains imposed by Network APIs and hardware interfaces

- Copy semantic of standard Unix network calls - Network cards usually unable to perform gather/scatter One feature that distinguish some network cards from others is the possibility to perform scatter/gather operations.  In this case the frame to be transmitted can lay fragmented in noncontiguous areas of memory and the card is able to collect them based on a list of pointers and to transmit it (gather).  The same can happen when a frame is received, even if this feature is less used. The different headers of the frame can be deposited in different noncontiguous areas (scatter).  Most of the cheap PC cards never supported this features.  Therefore it is required to prepare the complete frame to transmit in a physical contiguous area of memory.


## 3.  Fundamental data structures

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


### 3.1.  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.  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]*

```
1406    /* Deliver skb to an old protocol, which is not threaded well
1407       or which do not understand shared skbs.
1408     */
1409    static int deliver_to_old_ones(struct packet_type *pt,
1410                        struct sk_buff *skb, int last)
1411    {
1412         int ret = NET_RX_DROP;
1413
1414         if (!last) {
```

—————————————————————————————————————————————————————— *[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 block ( 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 data_lenlen- 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]*

```
22    /* Packet types */
23
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 UNNECESSARY :

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

```
34
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
141            memset(skb, 0, offsetof(struct sk_buff, truesize));
142            skb->truesize = size + sizeof(struct sk_buff);
143            atomic_set(&skb->users, 1);
```

*[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     struct sk_buff {
186             /* These two members must be first. */
187             struct sk_buff      *next;
188             struct sk_buff      *prev;
189
190             struct sk_buff_head *list;
191             struct sock         *sk;
192             struct timeval      stamp;
193             struct net_device   *dev;
194             struct net_device   *real_dev;
195
196             union {
197                     struct tcphdr   *th;
198                     struct udphdr   *uh;
199                     struct icmphdr  *icmph;
200                     struct igmphdr  *igmph;
201                     struct iphdr    *ipiph;
202                     unsigned char   *raw;
203             } h;
204
205             union {
206                     struct iphdr    *iph;
207                     struct ipv6hdr  *ipv6h;
208                     struct arphdr   *arph;
209                     unsigned char   *raw;
210             } nh;
211
212             union {
213                     struct ethhdr   *ethernet;
214                     unsigned char   *raw;
215             } mac;
216
217             struct  dst_entry   *dst;
218             struct    sec_path  *sp;
219
220             /*
221              * This is the control buffer. It is free to use for every
222              * layer. Please put your private variables there. If you
223              * want to keep them across layers you have to do a skb_clone()
224              * first. This is owned by whoever has the skb queued ATM.
225              */
226             char            cb[48];
227
228             unsigned int        len,
229                         data_len,
230                         csum;
231             unsigned char       local_df,
232                         cloned,
233                         pkt_type,
234                         ip_summed;
235             __u32           priority;
```

```
236          unsigned short          protocol,
237                           security;
238
239          void              (*destructor)(struct sk_buff *skb);
240     #ifdef CONFIG_NETFILTER
241            unsigned long           nfmark;
242          __u32            nfcache;
243          struct nf_ct_info    *nfct;
244     #ifdef CONFIG_NETFILTER_DEBUG
245            unsigned int        nf_debug;
246     #endif
247     #if defined(CONFIG_BRIDGE) || defined(CONFIG_BRIDGE_MODULE)
248          struct nf_bridge_info     *nf_bridge;
249     #endif
250     #endif /* CONFIG_NETFILTER */
251     #if defined(CONFIG_HIPPI)
252          union {
253                __u32      ifield;
254          } private;
255     #endif
256     #ifdef CONFIG_NET_SCHED
257            __u32                tc_index;                /* traffic control index */
258     #endif
259
260          /* These elements must be at the end, see alloc_skb() for details.  */
261          unsigned int        truesize;
262          atomic_t        users;
263          unsigned char        *head,
264                           *data,
265                           *tail,
266                           *end;
267     };
268
```

### 3.2. 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. 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     /* To allow 64K frame to be packed as single skb without frag_list */
125     #define MAX_SKB_FRAGS (65536/PAGE_SIZE + 2)
```

*[include/linux/skbuff.h]*

---

*[include/linux/skbuff.h]*

```
138     struct skb_shared_info {
139         atomic_t   dataref;
140         unsigned int    nr_frags;
141         unsigned short  tso_size;
142         unsigned short  tso_segs;
143         struct sk_buff  *frag_list;
144         skb_frag_t frags[MAX_SKB_FRAGS];
145     };
```

*[include/linux/skbuff.h]*

---

## 4. Skbuff organizations

Until recently the data area of the skbuff was unique and physically contiguous. 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
net packet.

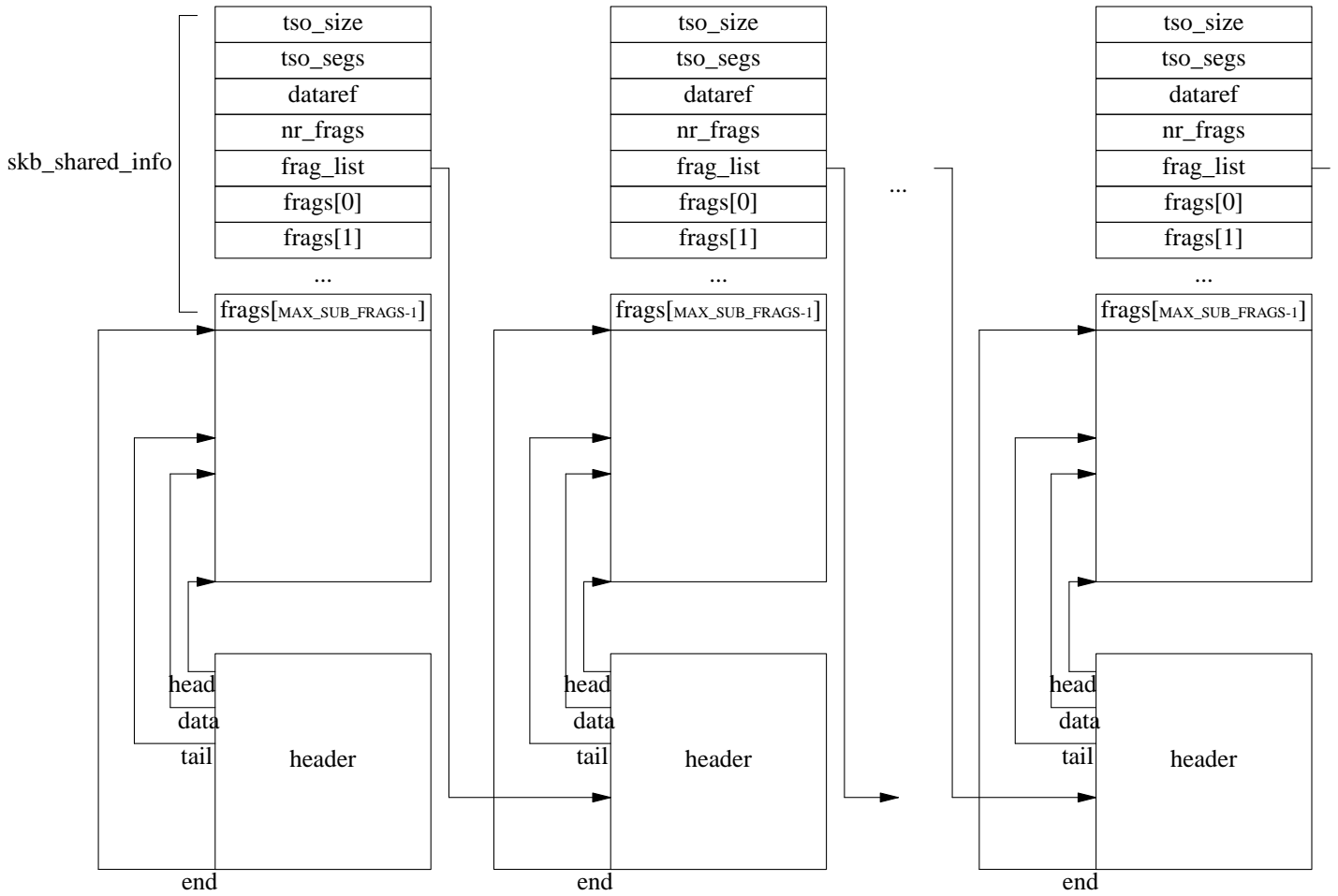### 4.1. Linear skbuffs

## 4.2.  Nonlinear skbuffs

### 4.2.1.  array of pages fragmentation

**4.2.2. skbuff list fragmentation**

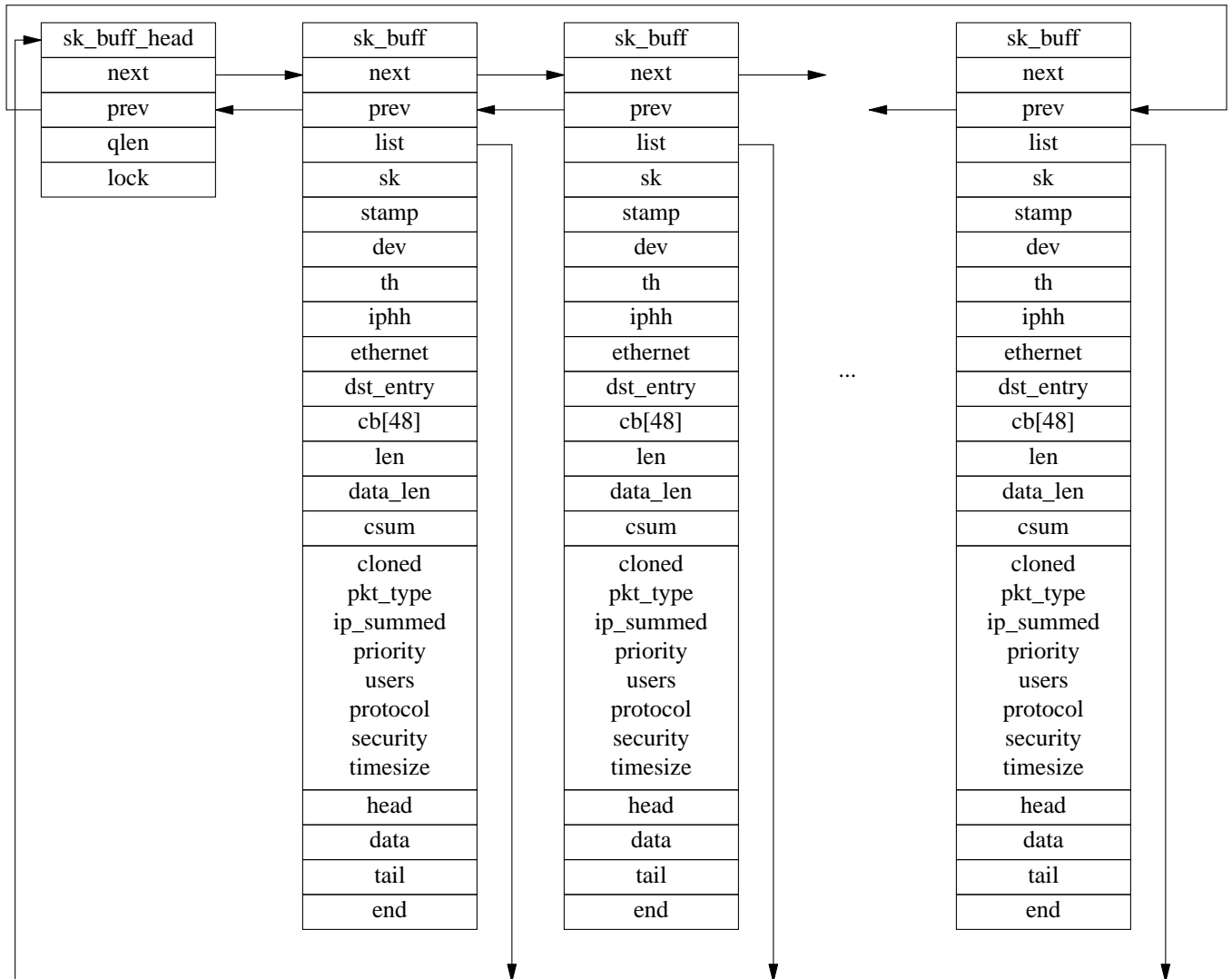## 5. Queues of skbuffs

In output skbuff is queued first on the socket, and then when the output interface is determined the skbuff is moved to the device queue. In input the skbuff is queued on the device queue and then when the owner socket is found it is moved to the owner socket queue.

| sk_buff_head | sk_buff | sk_buff | | sk_buff |
|---|---|---|---|---|
| next | next | next | | next |
| prev | prev | prev | | prev |
| qlen | list | list | | list |
| lock | sk | sk | | sk |
| | stamp | stamp | | stamp |
| | dev | dev | | dev |
| | th | th | | th |
| | iphh | iphh | | iphh |
| | ethernet | ethernet | | ethernet |
| | dst_entry | dst_entry | ... | dst_entry |
| | cb[48] | cb[48] | | cb[48] |
| | len | len | | len |
| | data_len | data_len | | data_len |
| | csum | csum | | csum |
| | cloned pkt_type ip_summed priority users protocol security timesize | cloned pkt_type ip_summed priority users protocol security timesize | | cloned pkt_type ip_summed priority users protocol security timesize |
| | head | head | | head |
| | data | data | | data |
| | tail | tail | | tail |
| | end | end | | end |

## 5.1.  Functions to manage skbuff queues

### 5.1.1.  skb_queue_len

This function returns the number of skbuff queued on the list that you pass as an argument. *File :* `[include/linux/skbuff.h]`

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

```
478     static inline __u32 skb_queue_len(struct sk_buff_head *list_)
479     {
```

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

### 5.1.2. skb_queue_head_init

This function initializes a queue of skbuffs. It locks the list setting the list establishing a spin lock on its lock variable and sets the `prev` and `next` pointers to the list head. *File :* `[include/linux/skbuff.h]`

*[include/linux/skbuff.h]*

```
483    static inline void skb_queue_head_init(struct sk_buff_head *list)
484    {
485        spin_lock_init(&list->lock);
486        list->prev = list->next = (struct sk_buff *)list;
487        list->qlen = 0;
488    }
```

*[include/linux/skbuff.h]*

### 5.1.3. 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 establishes a spin lock on the queue lock variable and so it is safe to be called 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    static inline void __skb_queue_head(struct sk_buff_head *list,
508                         struct sk_buff *newsk)
509    {
510        struct sk_buff *prev, *next;
511
512        newsk->list = list;
513        list->qlen++;
514        prev = (struct sk_buff *)list;
515        next = prev->next;
516        newsk->next = next;
517        newsk->prev = prev;
518        next->prev  = prev->next = newsk;
519    }
```

*[include/linux/skbuff.h]*

*[include/linux/skbuff.h]*

```
533    static inline void skb_queue_head(struct sk_buff_head *list,
534                         struct sk_buff *newsk)
535    {
536        unsigned long flags;
537
538        spin_lock_irqsave(&list->lock, flags);
539        __skb_queue_head(list, newsk);
540        spin_unlock_irqrestore(&list->lock, flags);
541    }
```

*[include/linux/skbuff.h]*

### 5.1.4. skb_queue_tail and __skb_queue_tail

These two functions queue an skbuff buffer at the tail of a list. The `skb_queue_tail` function establishes a spin lock on the queue lock variable and so it is safe to be called with interrupts enable, 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]*

```
553    static inline void __skb_queue_tail(struct sk_buff_head *list,
554                        struct sk_buff *newsk)
555    {
556        struct sk_buff *prev, *next;
557
558        newsk->list = list;
559        list->qlen++;
560        next = (struct sk_buff *)list;
561        prev = next->prev;
562        newsk->next = next;
563        newsk->prev = prev;
564        next->prev  = prev->next = newsk;
565    }
```

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

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

```
578    static inline void skb_queue_tail(struct sk_buff_head *list,
579                        struct sk_buff *newsk)
580    {
581        unsigned long flags;
582
583        spin_lock_irqsave(&list->lock, flags);
584        __skb_queue_tail(list, newsk);
585        spin_unlock_irqrestore(&list->lock, flags);
586    }
```

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

### 5.1.5. skb_dequeue and __skb_dequeue

These two functions dequeue an skbuff buffer from the head of a list. The `skb_dequeue` function establishes a spin lock on the queue lock variable and so it is safe to be called 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 :* `[include/linux/skbuff.h]`

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

```
624    static inline struct sk_buff *skb_dequeue(struct sk_buff_head *list)
625    {
626        unsigned long flags;
627        struct sk_buff *result;
628
629        spin_lock_irqsave(&list->lock, flags);
```

```
630            result = __skb_dequeue(list);
631            spin_unlock_irqrestore(&list->lock, flags);
632            return result;
633    }
```

```
596    static inline struct sk_buff *__skb_dequeue(struct sk_buff_head *list)
597    {
598            struct sk_buff *next, *prev, *result;
599
600            prev = (struct sk_buff *) list;
601            next = prev->next;
602            result = NULL;
603            if (next != prev) {
604                    result          = next;
605                    next        = next->next;
606                    list->qlen--;
607                    next->prev  = prev;
608                    prev->next  = next;
609                    result->next = result->prev = NULL;
610                    result->list = NULL;
611            }
612            return result;
613    }
```

### 5.1.6. skb_insert and __skb_insert

### 5.1.7. skb_append and __skb_append

### 5.1.8. skb_unlink and __skb_unlink

### 5.1.9. skb_dequeue_tail and __skb_dequeue_tail

## 6. Skbuff Functions

The following functions distinguish between the three kind of skbuffs : linear, fragmented in an array of additional pages, fragmented in a list of skbuffs.

### 6.1. 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]*

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

*[include/linux/skbuff.h]*

*[include/linux/skbuff.h]*

```
778    static inline int skb_is_nonlinear(const struct sk_buff *skb)
779    {
780         return skb->data_len;
781    }
```

*[include/linux/skbuff.h]*

### 6.2. SKB_PAGE_ASSERT

*File :* [include/linux/skbuff.h]
This macro will raise a bug if the skbuff is fragmented in additional pages. 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 variable.

*[include/linux/skbuff.h]*

```
806    #define SKB_PAGE_ASSERT(skb)    BUG_ON(skb_shinfo(skb)->nr_frags)
```

*[include/linux/skbuff.h]*

 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 variable.

### 6.3. SKB_FRAG_ASSERT

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

*[include/linux/skbuff.h]*

```
807    #define SKB_FRAG_ASSERT(skb)    BUG_ON(skb_shinfo(skb)->frag_list)
```

*[include/linux/skbuff.h]*

### 6.4. 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]*

```
783    static inline unsigned int skb_headlen(const struct sk_buff *skb)
784    {
785         return skb->len - skb->data_len;
786    }
```

—————————————————————————————————————————————— *[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]*
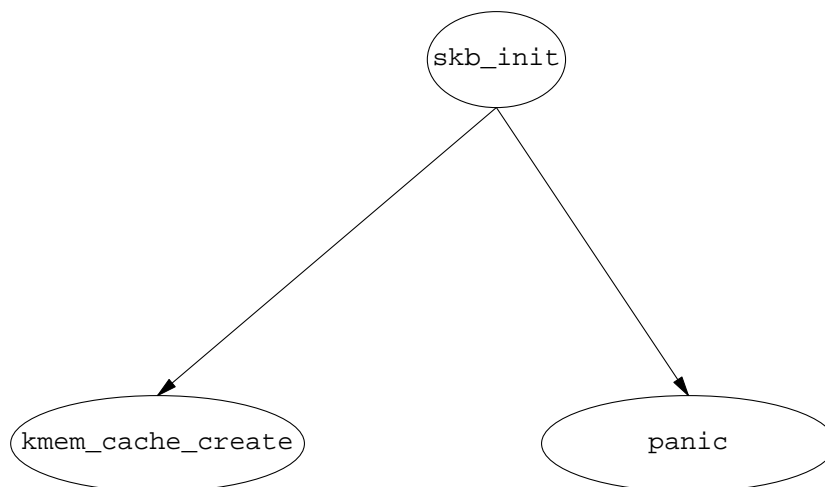
```
788     static inline int skb_pagelen(const struct sk_buff *skb)
789     {
790          int i, len = 0;
791
792          for (i = (int)skb_shinfo(skb)->nr_frags - 1; i >= 0; i--)
793               len += skb_shinfo(skb)->frags[i].size;
794          return len + skb_headlen(skb);
795     }
```

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

## 6.5. skb_init

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



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 specifiy any specific constructor or destructor for them.

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

```
1096    void __init skb_init(void)
1097    {
1098         skbuff_head_cache = kmem_cache_create("skbuff_head_cache",
1099                                     sizeof(struct sk_buff),
1100                                     0,
1101                                     SLAB_HWCACHE_ALIGN,
1102                                     NULL, NULL);
1103         if (!skbuff_head_cache)
```

```
1104                panic("cannot create skbuff cache");
1105    }
```

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

## 6.6. alloc_skb

*File :* [include/net/tcp.h]
This function allocates a network buffer. It 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 tcp code for instance
uses its own skb allocator tcp_alloc_skb to request additional MAX_TCP_HEADER bytes to accomo-
date a headroom sufficient for the header ( usually this is 128+32=160 bytes ).

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

```
1808    static inline struct sk_buff *tcp_alloc_pskb(struct sock *sk, int size, int
mem, int gfp)
1809    {
1810            struct sk_buff *skb = alloc_skb(size+MAX_TCP_HEADER, gfp);
1811
1812            if (skb) {
1813                    skb->truesize += mem;
1814                    if (sk->sk_forward_alloc >= (int)skb->truesize ||
1815                        tcp_mem_schedule(sk, skb->truesize, 0)) {
1816                            skb_reserve(skb, MAX_TCP_HEADER);
1817                            return skb;
1818                    }
1819                    __kfree_skb(skb);
1820            } else {
1821                    tcp_enter_memory_pressure();
1822                    tcp_moderate_sndbuf(sk);
1823            }
1824            return NULL;
1825    }
1826
1827    static inline struct sk_buff *tcp_alloc_skb(struct sock *sk, int size, int gfp)
1828    {
1829            return tcp_alloc_pskb(sk, size, 0, gfp);
1830    }
1831
```

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

In the ip fragmentation case for instance, additional bytes for the ip header and the link layer header properly aligned are requested :

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

```
583                /*
584                 *    Allocate buffer.
585                 */
586
```

```
   587                   if  ((skb2   =   alloc_skb(len+hlen+LL_RESERVED_SPACE(rt->u.dst.dev),
GFP_ATOMIC)) == NULL) {
   588                       NETDEBUG(printk(KERN_INFO "IP:  frag:  no  memory  for  new  frag-
ment!0));
   589                       err = -ENOMEM;
   590                       goto fail;
   591                   }
```
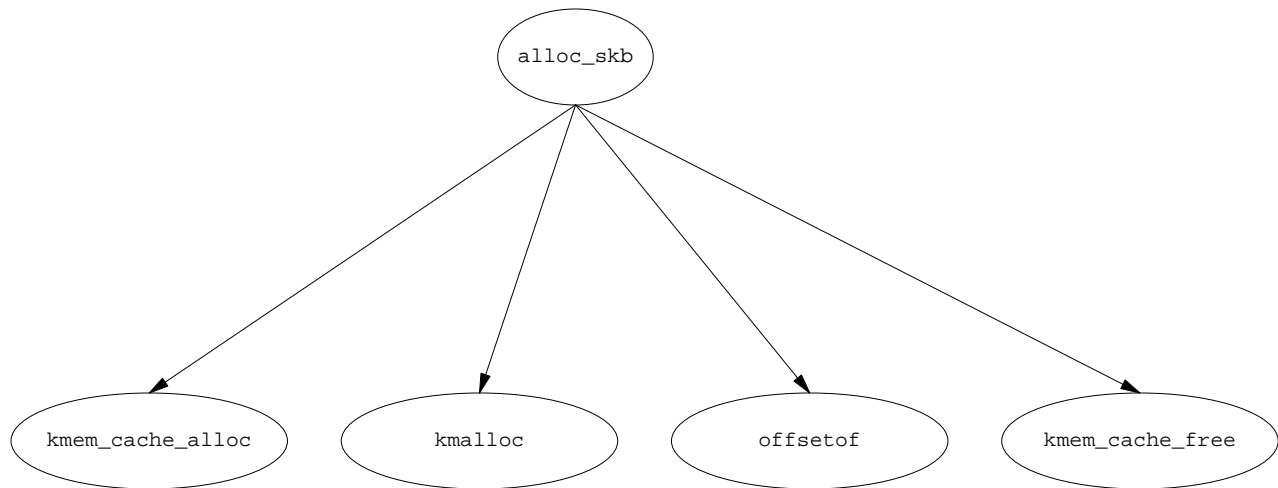
--------------------------------------------------------------------------------- *[net/ipv4/ip_output.c]*

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

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]*

```
   112    /**
   113     *    alloc_skb -    allocate a network buffer
   114     *    @size: size to allocate
   115     *    @gfp_mask: allocation mask
   116     *
   117     *    Allocate a new &sk_buff. The returned buffer has no headroom and a
   118     *    tail room of size bytes. The object has a reference count of one.
   119     *    The return is the buffer. On a failure the return is %NULL.
   120     *
   121     *    Buffers may only be allocated from interrupts using a @gfp_mask of
   122     *    %GFP_ATOMIC.
   123     */
   124    struct sk_buff *alloc_skb(unsigned int size, int gfp_mask)
   125    {
   126          struct sk_buff *skb;
   127          u8 *data;
   128
```

--------------------------------------------------------------------------------- *[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 you dont perform i/o over its data, so we reset that flag 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]*

```
129           /* Get the HEAD */
130           skb = kmem_cache_alloc(skbuff_head_cache,
131                           gfp_mask & ~__GFP_DMA);
132           if (!skb)
133               goto out;
134
```

*[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` that can store the information on the `frag_list` or `frags[]` array of pages used by fragmented skbuffs.  Th SKB_DATA_ALIGN macro adds enough bytes to the requested size so that the skb data area can be aligned with a level 1 cache line ( on P4 for example the X86_L1_CACHE_SHIFT is 2ˆ7=128 bytes and so 127 is added )

*[net/core/skbuff.c]*

```
135           /* Get the DATA. Size must match skb_add_mtu(). */
136           size = SKB_DATA_ALIGN(size);
137           data = kmalloc(size + sizeof(struct skb_shared_info), gfp_mask);
138           if (!data)
139               goto nodata;
140
```

*[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   out:
155           return skb;
156   nodata:
157           kmem_cache_free(skbuff_head_cache, skb);
158           skb = NULL;
159           goto out;
160   }
```

*[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           memset(skb, 0, offsetof(struct sk_buff, truesize));
142           skb->truesize = size + sizeof(struct sk_buff);
```

*[net/core/skbuff.c]*

Then the skbuff pointers inside the data area are initialized to a 0 size area :
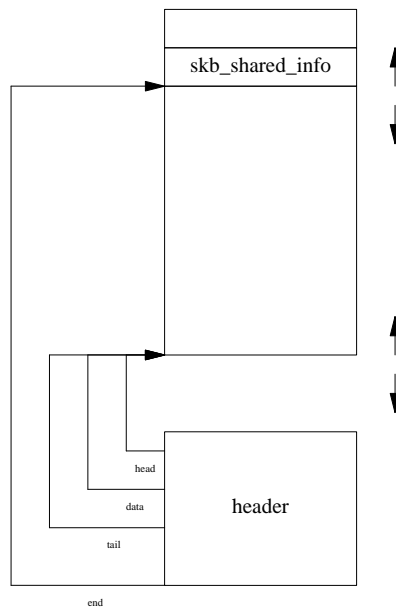
*[net/core/skbuff.c]*

```
144          skb->head = data;
145          skb->data = data;
146          skb->tail = data;
147          skb->end  = data + size;
148
```

And the shared_info are then initialized to an unfragmented skbuff :

```
149          atomic_set(&(skb_shinfo(skb)->dataref), 1);
```

only 1 reference to this skb, itself.

```
150          skb_shinfo(skb)->nr_frags  = 0;
```

no pages in the frags array. The tso_ fields refer to the Tcp Segmentation Offloading experimental kernel feature to support some intellignet network cards that can perform the tcp segmentation (Myricom Gigabit Ethernet, National DP83820,.. ). This feature is described by the NETIF_F_TSO flag in the netdev structure. The tso_size is set to the mtu - hlen and tso_segs is the number of segments required to transmit this skbuff.

```
151          skb_shinfo(skb)->tso_size = 0;
152          skb_shinfo(skb)->tso_segs = 0;
153          skb_shinfo(skb)->frag_list = NULL;
```

and an empty fragment list.

Fig. - After alloc_sk

---------------------------------- The `skb->data_len` field .. it is different from zero only on nonlinear skbuffs. In fact the function `skb_is_nonlinear()` returns that field. It seems to represent the number of bytes in the remaining skbuff (after the 1st). The function `skb_headlen()` returns the bytes in the linear part of the first skbuff : skb->data_len.`skb->len`-

------------------- 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.

### 6.7. 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                  th = (struct tcphdr *) skb_push(skb, tcp_header_size);
228                  skb->h.th = th;
```

*[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]*

```
848
849    /**
850     *   skb_push - add data to the start of a buffer
851     *   @skb: buffer to use
852     *   @len: amount of data to add
853     *
854     *   This function extends the used data area of the buffer at the buffer
855     *   start. If this would exceed the total buffer headroom the kernel will
856     *   panic. A pointer to the first byte of the extra data is returned.
```

```
857      */
858      static inline unsigned char *skb_push(struct sk_buff *skb, unsigned int len)
859      {
860            skb->data -= len;
861            skb->len  += len;
862            if (unlikely(skb->data<skb->head))
863                  skb_under_panic(skb, len, current_text_addr());
864            return skb->data;
865      }
```
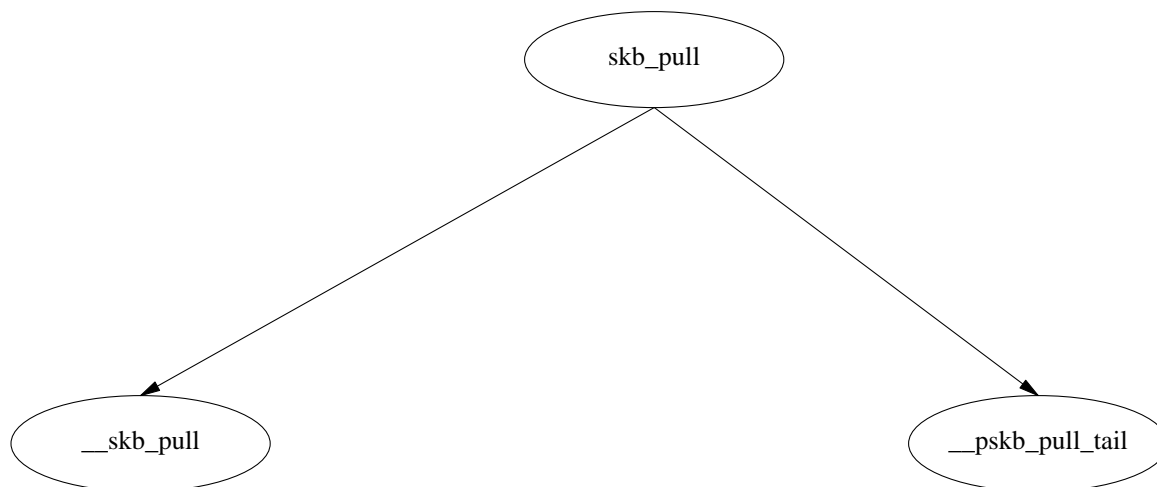
*[include/linux/skbuff.h]*

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. This is the very easily understandable implementation without the check :

*[include/linux/skbuff.h]*

```
841
842      static inline unsigned char *__skb_push(struct sk_buff *skb, unsigned int len)
843      {
844            skb->data -= len;
845            skb->len  += len;
846            return skb->data;
847      }
848
```

*[include/linux/skbuff.h]*

I dont know why this is not called inside skb_push.

## 6.8.  skb_pull()

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



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.

| Tail Room |
|---|

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

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

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

| Head Room | skb_push area | Data Area | skb_put area | Tail Room |
|---|---|---|---|---|

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 here is eventually performed to check that `skb->len` is less than `skb->data_len`.

<div align="right"><em>[include/linux/skbuff.h]</em></div>

```
874    /**
875     *    skb_pull - remove data from the start of a buffer
876     *    @skb: buffer to use
877     *    @len: amount of data to remove
878     *
879     *    This function removes data from the start of a buffer, returning
880     *    the memory to the headroom. A pointer to the next data in the buffer
881     *    is returned. Once the data has been pulled future pushes will overwrite
882     *    the old data.
883     */
884    static inline unsigned char *skb_pull(struct sk_buff *skb, unsigned int len)
885    {
886            return (len > skb->len) ? NULL : __skb_pull(skb, len);
887    }
888
```

<div align="right"><em>[include/linux/skbuff.h]</em></div>

<div align="right"><em>[include/linux/skbuff.h]</em></div>

```
866
867    static inline char *__skb_pull(struct sk_buff *skb, unsigned int len)
```

```
868        {
869              skb->len -= len;
870              BUG_ON(skb->len < skb->data_len);
871              return skb->data += len;
872        }
873
```

*[include/linux/skbuff.h]*

## 6.9.  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     static void skb_drop_fraglist(struct sk_buff *skb)
164     {
165              struct sk_buff *list = skb_shinfo(skb)->frag_list;
166
167              skb_shinfo(skb)->frag_list = NULL;
168
169              do {
170                      struct sk_buff *this = list;
171                      list = list->next;
172                      kfree_skb(this);
173              } while (list);
174     }
```

*[net/core/skbuff.c]*

## 6.10.  ___pskb_trim

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

This function is called to trim a nonlinear skbuff. An skbuff can be constituted by an array of pages ( `skb->frags[]` ) and/or a list of skbuffs ( `skb->fraglist` ).

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

```
632
633     /* Trims skb to length len. It can change skb pointers, if "realloc" is 1.
634      * If realloc==0 and trimming is impossible without change of data,
635      * it is BUG().
636      */
637
638     int ___pskb_trim(struct sk_buff *skb, unsigned int len, int realloc)
639     {
640           int offset = skb_headlen(skb);
641           int nfrags = skb_shinfo(skb)->nr_frags;
642           int i;
643
644           for (i = 0; i < nfrags; i++) {
645                 int end = offset + skb_shinfo(skb)->frags[i].size;
646                 if (end > len) {
647                       if (skb_cloned(skb)) {
648                             if (!realloc)
649                                   BUG();
650                             if (pskb_expand_head(skb, 0, 0, GFP_ATOMIC))
651                                   return -ENOMEM;
652                       }
653                       if (len <= offset) {
654                             put_page(skb_shinfo(skb)->frags[i].page);
655                             skb_shinfo(skb)->nr_frags--;
656                       } else {
657                             skb_shinfo(skb)->frags[i].size = len - offset;
658                       }
659                 }
660                 offset = end;
661           }
```

_____ *[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.

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

```
662
663          if (offset < len) {
664                skb->data_len -= skb->len - len;
665                skb->len        = len;
666          } else {
667                if (len <= skb_headlen(skb)) {
668                      skb->len      = len;
669                      skb->data_len = 0;
670                      skb->tail     = skb->data + len;
671                      if (skb_shinfo(skb)->frag_list && !skb_cloned(skb))
672                            skb_drop_fraglist(skb);
673                } else {
674                      skb->data_len -= skb->len - len;
675                      skb->len        = len;
676                }
677          }
678
679          return 0;
680    }
```
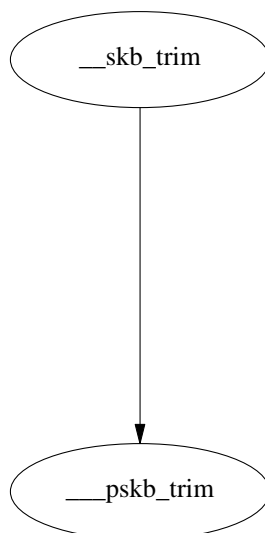
―――――――――――――――――――――――――――――――――――――――――――――――――――――――――― *[net/core/skbuff.c]*

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

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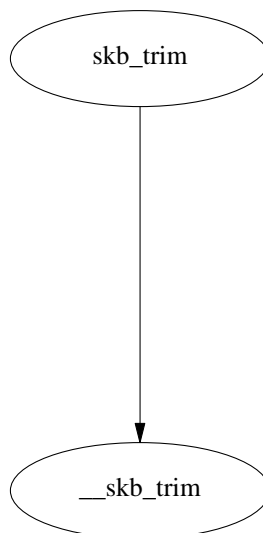
*File :* `[include/linux/skbuff.h]`



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

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

```
951
952     static inline void __skb_trim(struct sk_buff *skb, unsigned int len)
953     {
954          if (!skb->data_len) {
955               skb->len  = len;
956               skb->tail = skb->data + len;
957          } else
958               ___pskb_trim(skb, len, 0);
959     }
960
```

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

## 6.12. skb_trim

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



This wrapper function just 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     static inline void skb_trim(struct sk_buff *skb, unsigned int len)
970     {
971          if (skb->len > len)
972               __skb_trim(skb, len);
973     }
974
```

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

## 6.13. skb_reserve

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

```
                        ╭─────────────────────╮
                        │     skb_reserve      │
                        ╰─────────────────────╯
                                   │
                                   │
                                   ▼
                        ╭─────────────────────╮
                        │     ___pskb_trim     │
                        ╰─────────────────────╯
```

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]*

```
936    /**
937     *   skb_reserve - adjust headroom
938     *   @skb: buffer to alter
939     *   @len: bytes to move
940     *
941     *   Increase the headroom of an empty &sk_buff by reducing the tail
942     *   room. This is only allowed for an empty buffer.
943     */
944    static inline void skb_reserve(struct sk_buff *skb, unsigned int len)
945    {
946          skb->data += len;
947          skb->tail += len;
948    }
949
```

*[include/linux/skbuff.h]*

```
                    end  ┌──────────────┐
                         │              │
                         │              │
                         │              │
                         │              │
               data, tail├──────────────┤
                         │              │
                         │              │
                         │              │
                    head └──────────────┘
```

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

Fig. - After skb_reserve

## 6.14.  skb_release_data

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



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

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

## 6.15. 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      *    Free an skbuff by memory without cleaning the state.
203      */
204     void kfree_skbmem(struct sk_buff *skb)
205     {
206             skb_release_data(skb);
207             kmem_cache_free(skbuff_head_cache, skb);
208     }
209
```

*[net/core/skbuff.c]*

## 6.16. [include/linux/skbuff.h]

*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). This function is usually called after the skb_reserve() function has been called on a newly allocated skbuff moving the data pointer, to move the tail pointer and prepare the space to copy over the data. It updates the len field of the skbuff header. This operation can be applied only on linear skbuffs.

```
808     #define SKB_LINEAR_ASSERT(skb)  BUG_ON(skb_is_nonlinear(skb))
809
810     /*
811      *    Add data to an sk_buff
812      */
813     static inline unsigned char *__skb_put(struct sk_buff *skb, unsigned int len)
814     {
815          unsigned char *tmp = skb->tail;
816          SKB_LINEAR_ASSERT(skb);
817          skb->tail += len;
818          skb->len  += len;
819          return tmp;
820     }
821
822     /**
```

```
823     *    skb_put - add data to a buffer
824     *    @skb: buffer to use
825     *    @len: amount of data to add
826     *
827     *    This function extends the used data area of the buffer. If this would
828     *    exceed the total buffer size the kernel will panic. A pointer to the
829     *    first byte of the extra data is returned.
830     */
831     static inline unsigned char *skb_put(struct sk_buff *skb, unsigned int len)
832     {
833           unsigned char *tmp = skb->tail;
834           SKB_LINEAR_ASSERT(skb);
835           skb->tail += len;
836           skb->len  += len;
837           if (unlikely(skb->tail>skb->end))
838                 skb_over_panic(skb, len, current_text_addr());
839           return tmp;
840     }
```

*[include/linux/skbuff.h]*

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

Fig. - An sk_buff containing data

## 6.17.  [include/linux/skbuff.h] kfree_skb()

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



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     *   kfree_skb - free an sk_buff
334     *   @skb: buffer to free
335     *
336     *   Drop a reference to the buffer and free it if the usage count has
337     *   hit zero.
338     */
339    static inline void kfree_skb(struct sk_buff *skb)
340    {
341          if (atomic_read(&skb->users) == 1 || atomic_dec_and_test(&skb->users))
342                __kfree_skb(skb);
343    }
344
```

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

### 6.18. __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     *   Free an skbuff by memory without cleaning the state.
203     */
204    void kfree_skbmem(struct sk_buff *skb)
205    {
206          skb_release_data(skb);
```
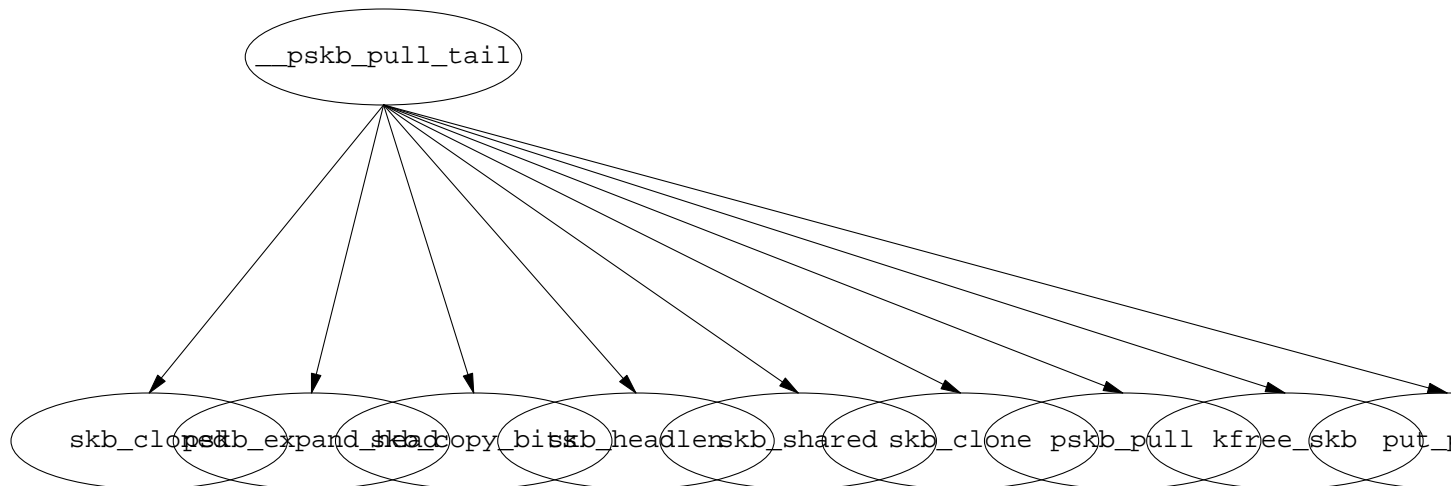
```
207            kmem_cache_free(skbuff_head_cache, skb);
208       }
209
210       /**
211        *    __kfree_skb - private function
212        *    @skb: buffer
213        *
214        *    Free an sk_buff. Release anything attached to the buffer.
215        *    Clean the state. This is an internal helper function. Users should
216        *    always call kfree_skb
217        */
218
219       void __kfree_skb(struct sk_buff *skb)
220       {
221            if (skb->list) {
222                 printk(KERN_WARNING "Warning: kfree_skb passed an skb still "
223                        "on a list (from %p).0, NET_CALLER(skb));
224                 BUG();
225            }
226
227            dst_release(skb->dst);
228       #ifdef CONFIG_XFRM
229            secpath_put(skb->sp);
230       #endif
231            if(skb->destructor) {
232                 if (in_irq())
233                     printk(KERN_WARNING "Warning: kfree_skb on "
234                                "hard IRQ %p0, NET_CALLER(skb));
235                 skb->destructor(skb);
236            }
237       #ifdef CONFIG_NETFILTER
238            nf_conntrack_put(skb->nfct);
239       #if defined(CONFIG_BRIDGE) || defined(CONFIG_BRIDGE_MODULE)
240            nf_bridge_put(skb->nf_bridge);
241       #endif
242       #endif
243            kfree_skbmem(skb);
244       }
245
```

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

## 6.19.  __pskb_pull_tail

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

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, 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. We allocate a new linear part also in the case the skbuff has been cloned since we want to change that data area. Then we copy delta bytes 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 we just have to pull the array of pages, otherwise we have to go through the frag_list. If the array of pages associated with this skbuff is large enough then again we just have to pull the array. Then you go through the frag_list and you eat (kfree_skb) all the complete skbuffs that you can. 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.

**6.20. skb_clone**

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

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

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

327

## 6.21.  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 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 initializes the number of users of the new skb to 1.
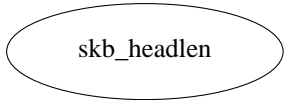
*[net/core/skbuff.c]*

```
328    static void copy_skb_header(struct sk_buff *new, const struct sk_buff *old)
329    {
330         /*
331          *   Shift between the two data areas in bytes
332          */
333         unsigned long offset = new->data - old->data;
334
335         new->list  = NULL;
336         new->sk        = NULL;
337         new->dev   = old->dev;
338         new->real_dev   = old->real_dev;
339         new->priority   = old->priority;
340         new->protocol   = old->protocol;
341         new->dst   = dst_clone(old->dst);
342    #ifdef CONFIG_INET
343         new->sp        = secpath_get(old->sp);
344    #endif
345         new->h.raw = old->h.raw + offset;
346         new->nh.raw     = old->nh.raw + offset;
347         new->mac.raw     = old->mac.raw + offset;
348         memcpy(new->cb, old->cb, sizeof(old->cb));
```

```
349            new->local_df   = old->local_df;
350            new->pkt_type   = old->pkt_type;
351            new->stamp = old->stamp;
352            new->destructor = NULL;
353            new->security   = old->security;
354     #ifdef CONFIG_NETFILTER
355            new->nfmark     = old->nfmark;
356            new->nfcache    = old->nfcache;
357            new->nfct = old->nfct;
358            nf_conntrack_get(old->nfct);
359     #ifdef CONFIG_NETFILTER_DEBUG
360            new->nf_debug   = old->nf_debug;
361     #endif
362     #if defined(CONFIG_BRIDGE) || defined(CONFIG_BRIDGE_MODULE)
363            new->nf_bridge = old->nf_bridge;
364            nf_bridge_get(old->nf_bridge);
365     #endif
366     #endif
367     #ifdef CONFIG_NET_SCHED
368            new->tc_index   = old->tc_index;
369     #endif
370            atomic_set(&new->users, 1);
371     }
```

*[net/core/skbuff.c]*

### 6.22. skb_headlen

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



This function returns the number of data bytes in the first skbuff data area. It should be equal to (skb->tail - skb->data) i think.

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

### 6.23.  skb_copy_bits

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



This function copies bytes of data from an skb to a another area of memory.  The first argument is a pointer
to an skb header from which data area the data should be copied, the 3d arg is a pointer to an area of mem-
ory where the data should be put. 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 The len argument is the number of bytes that should be copied.  This function is
able to treat fragmented skbuff and has code to copy all the fragments in the array of pages and all the even-
tual skbuffs linked together calling iteratively skb_copy_bits for each skbuff in the frag_list.

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

```
820
821     /* Copy some data bits from skb to kernel buffer. */
822
823     int skb_copy_bits(const struct sk_buff *skb, int offset, void *to, int len)
824     {
825          int i, copy;
826          int start = skb_headlen(skb);
827
828          if (offset > (int)skb->len - len)
829               goto fault;
830
831          /* Copy header. */
832          if ((copy = start - offset) > 0) {
833               if (copy > len)
834                    copy = len;
835               memcpy(to, skb->data + offset, copy);
836               if ((len -= copy) == 0)
837                    return 0;
838               offset += copy;
839               to    += copy;
840          }
841
842          for (i = 0; i < skb_shinfo(skb)->nr_frags; i++) {
843               int end;
844
845               BUG_TRAP(start <= offset + len);
846
```

```
847                       end = start + skb_shinfo(skb)->frags[i].size;
848                       if ((copy = end - offset) > 0) {
849                               u8 *vaddr;
850
851                               if (copy > len)
852                                       copy = len;
853
854                               vaddr = kmap_skb_frag(&skb_shinfo(skb)->frags[i]);
855                               memcpy(to,
856                                        vaddr + skb_shinfo(skb)->frags[i].page_offset+
857                                        offset - start, copy);
858                               kunmap_skb_frag(vaddr);
859
860                               if ((len -= copy) == 0)
861                                       return 0;
862                               offset += copy;
863                               to     += copy;
864                       }
865                       start = end;
866               }
867
868               if (skb_shinfo(skb)->frag_list) {
869                       struct sk_buff *list = skb_shinfo(skb)->frag_list;
870
871                       for (; list; list = list->next) {
872                               int end;
873
874                               BUG_TRAP(start <= offset + len);
875
876                               end = start + list->len;
877                               if ((copy = end - offset) > 0) {
878                                       if (copy > len)
879                                               copy = len;
880                                       if (skb_copy_bits(list, offset - start,
881                                                       to, copy))
882                                               goto fault;
883                                       if ((len -= copy) == 0)
884                                               return 0;
885                                       offset += copy;
886                                       to     += copy;
887                               }
888                               start = end;
889                       }
890               }
891               if (!len)
892                       return 0;
893
894       fault:
895           return -EFAULT;
896       }
897
```

### 6.24.  pskb_expand_head

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



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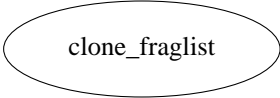
————————————————————————————————————————————————————— *[net/core/skbuff.c]*

```
473    /**
474     *    pskb_expand_head - reallocate header of &sk_buff
475     *    @skb: buffer to reallocate
476     *    @nhead: room to add at head
477     *    @ntail: room to add at tail
478     *    @gfp_mask: allocation priority
479     *
480     *    Expands (or creates identical copy, if &nhead and &ntail are zero)
481     *    header of skb. &sk_buff itself is not changed. &sk_buff MUST have
482     *    reference count of 1. Returns zero in the case of success or error,
483     *    if expansion failed. In the last case, &sk_buff is not changed.
484     *
485     *    All the pointers pointing into skb header may change and must be
486     *    reloaded after call to this function.
487     */
488
489    int pskb_expand_head(struct sk_buff *skb, int nhead, int ntail, int gfp_mask)
490    {
491          int i;
492          u8 *data;
493          int size = nhead + (skb->end - skb->head) + ntail;
494          long off;
495
496          if (skb_shared(skb))
```

```
497                    BUG();
498
499            size = SKB_DATA_ALIGN(size);
500
501            data = kmalloc(size + sizeof(struct skb_shared_info), gfp_mask);
502            if (!data)
503                    goto nodata;
504
505            /* Copy only real data... and, alas, header. This should be
506             * optimized for the cases when header is void. */
507            memcpy(data + nhead, skb->head, skb->tail - skb->head);
508            memcpy(data + size, skb->end, sizeof(struct skb_shared_info));
509
510            for (i = 0; i < skb_shinfo(skb)->nr_frags; i++)
511                    get_page(skb_shinfo(skb)->frags[i].page);
512
513            if (skb_shinfo(skb)->frag_list)
514                    skb_clone_fraglist(skb);
515
516            skb_release_data(skb);
517
518            off = (data + nhead) - skb->head;
519
520            skb->head    = data;
521            skb->end     = data + size;
522            skb->data    += off;
523            skb->tail    += off;
524            skb->mac.raw += off;
525            skb->h.raw   += off;
526            skb->nh.raw  += off;
527            skb->cloned  = 0;
528            atomic_set(&skb_shinfo(skb)->dataref, 1);
529            return 0;
530
531    nodata:
532            return -ENOMEM;
533    }
```

*— [net/core/skbuff.c]*

### 6.25.  [net/core/skbuff.c] clone_fraglist()

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

### 6.26.  [include/linux/skbuff.h] skb_get()

It increments the number of users of the skbuff.

```
314     /**
315      *   skb_get - reference buffer
316      *   @skb: buffer to reference
317      *
318      *   Makes another reference to a socket buffer and returns a pointer
319      *   to the buffer.
320      */
321     static inline struct sk_buff *skb_get(struct sk_buff *skb)
322     {
323           atomic_inc(&skb->users);
324           return skb;
325     }
326
```

**6.27.  [net/core/skbuff.c] skb_copy()**

Makes a complete copy of an skbuff header and its data.  It converts a nonlinear skbuff to a linear one.  The headroom of the old skbuff is computed , and inappropriately called haederlen.

```
                              ┌───────────┐
                              │ skb_copy  │
                              └───────────┘
```



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. skb->len + headroom = skb->end - skb->head + skb->data_len If it is not possible to allocate such an skbuff returns NULL. A headroom equal to the one in the old skb is reserved in the new skb. It sets the tail pointer in the new skb at skb->data+ skb->len. It copies the checksum and the ip_summed flag. Then it calls the skb_copy_bits function to copy everything since the very beginning of the old skbuff (skb->head not skb->data !!!!! ) to the end of the data. This means that it copies also the headroom. Finally it copies the skbuff header. and returns the pointer for the new skbuff.

```
373    /**
374     *   skb_copy    -    create private copy of an sk_buff
375     *   @skb: buffer to copy
376     *   @gfp_mask: allocation priority
377     *
378     *   Make a copy of both an &sk_buff and its data. This is used when the
379     *   caller wishes to modify the data and needs a private copy of the
380     *   data to alter. Returns %NULL on failure or the pointer to the buffer
381     *   on success. The returned buffer has a reference count of 1.
382     *
383     *   As by-product this function converts non-linear &sk_buff to linear
384     *   one, so that &sk_buff becomes completely private and caller is allowed
385     *   to modify all the data of returned buffer. This means that this
386     *   function is not recommended for use in circumstances when only
387     *   header is going to be modified. Use pskb_copy() instead.
388     */
389
390    struct sk_buff *skb_copy(const struct sk_buff *skb, int gfp_mask)
391    {
392         int headerlen = skb->data - skb->head;
393         /*
394          *    Allocate the copy buffer
395          */
396         struct sk_buff *n = alloc_skb(skb->end - skb->head + skb->data_len,
397                                 gfp_mask);
398         if (!n)
399              return NULL;
400
401         /* Set the data pointer */
402         skb_reserve(n, headerlen);
403         /* Set the tail pointer and length */
404         skb_put(n, skb->len);
405         n->csum      = skb->csum;
406         n->ip_summed = skb->ip_summed;
407
408         if (skb_copy_bits(skb, -headerlen, n->head, headerlen + skb->len))
409              BUG();
410
411         copy_skb_header(n, skb);
412         return n;
413    }
414
```

## 6.28. [net/core/skbuff.c] pskb_copy()

It allocates an skb sufficient for headlen + headroom. It reserves the same headroom available in the old skbuff. It copies the headlen bytes of data from the old skbuff to the new one. Copies checksum, ip_summed flag,data_len,len from old to new. Runs through the array of fragments, copy frag descriptors from the old to the new skbuff. And for each page increments the usage count. Then it has a frag_list, it copies the pointer, and it goes through the list of fragments (frag_list) and increases the usage count (skb_clone_fraglist). Finally it copies the skb header (copy_skb_header). And it returns a pointer to the new skbuff head.

```
416    /**
417     *    pskb_copy -    create copy of an sk_buff with private head.
418     *    @skb: buffer to copy
419     *    @gfp_mask: allocation priority
420     *
421     *    Make a copy of both an &sk_buff and part of its data, located
422     *    in header. Fragmented data remain shared. This is used when
423     *    the caller wishes to modify only header of &sk_buff and needs
424     *    private copy of the header to alter. Returns %NULL on failure
425     *    or the pointer to the buffer on success.
426     *    The returned buffer has a reference count of 1.
427     */
428
429    struct sk_buff *pskb_copy(struct sk_buff *skb, int gfp_mask)
430    {
431        /*
432         *    Allocate the copy buffer
433         */
434        struct sk_buff *n = alloc_skb(skb->end - skb->head, gfp_mask);
435
436        if (!n)
437            goto out;
438
439        /* Set the data pointer */
440        skb_reserve(n, skb->data - skb->head);
441        /* Set the tail pointer and length */
442        skb_put(n, skb_headlen(skb));
443        /* Copy the bytes */
444        memcpy(n->data, skb->data, n->len);
445        n->csum       = skb->csum;
446        n->ip_summed = skb->ip_summed;
447
448        n->data_len  = skb->data_len;
449        n->len        = skb->len;
450
451        if (skb_shinfo(skb)->nr_frags) {
452            int i;
453
454            for (i = 0; i < skb_shinfo(skb)->nr_frags; i++) {
455                skb_shinfo(n)->frags[i] = skb_shinfo(skb)->frags[i];
456                get_page(skb_shinfo(n)->frags[i].page);
457            }
458            skb_shinfo(n)->nr_frags = i;
459        }
460        skb_shinfo(n)->tso_size = skb_shinfo(skb)->tso_size;
461        skb_shinfo(n)->tso_segs = skb_shinfo(skb)->tso_segs;
462
463        if (skb_shinfo(skb)->frag_list) {
464            skb_shinfo(n)->frag_list = skb_shinfo(skb)->frag_list;
465            skb_clone_fraglist(n);
466        }
467
468        copy_skb_header(n, skb);
469    out:
```

```
470          return n;
471     }
472
```

# Table of Contents

**References**

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