# Received Packet Processing

In this section we consider the processing of a received packet as it moves from the device driver to the IP layer.   The *device driver* relies principally upon two kernel functions

| | |
|---|---|
| *dev_alloc_skb()* | Allocates an *sk_buff* of the required size prior to transferring the packet to kernel memory.  Two hardware strategies are commonly used.  If packets are received directly into system memory owned by the kernel, the *sk_buff* must be allocated prior to initiating the receive operation.  If packets are first received into NIC buffers and then transferred via DMA to system memory,  an *sk_buff* of the exact size needed may be allocated after the packet has been received but before the DMA transfer is initiated. |
| *netif_rx()* | Used to pass the *sk_buff* to the generic device layer when a receive operation completes. |

The example below is taken from drivers/net/3c59x.c

```
skb = dev_alloc_skb(pkt_len + 5);
skb->protocol = eth_type_trans(skb, dev);
 :
netif_rx(skb);
```

## Allocation and initialization of the *sk_buff*

*dev_alloc_skb()* is defined in include/linux/skbuff.h. It merely calls *__dev_alloc_skb* with the GFP_ATOMIC flag set.   This flag forces *kmalloc()* to return an error code rather than sleeping if no memory is available.   It is necessary because sleeping in an interrupt context is a fatal error.

```
1053 static inline struct sk_buff *dev_alloc_skb(unsigned
                                   int length)
1054 {
1055     return __dev_alloc_skb(length, GFP_ATOMIC);
1056 }
```

__dev_alloc_skb() allocates memory for the sk_buff of specified size. A few additional bytes are always allocated for alignment optimization purposes.

```
1028
1029 static inline struct sk_buff *__dev_alloc_skb(unsigned
                              int length, int gfp_mask)
1031 {
1032     struct sk_buff *skb;
1033
1034     skb = alloc_skb(length+16, gfp_mask);
```

The alloc_skb() function allocates the sk_buff. Recall that an sk_buff consists of a fixed size header of type struct sk_ buff which is allocated from a cache of such objects and a variable size data buffer allocated from one of the general caches and actually holds the packet's headers and user data.

```
164 struct sk_buff *alloc_skb(unsigned int size, int gfp_mask)
165 {
166     struct sk_buff *skb;
167     u8 *data;
168
169     if (in_interrupt() && (gfp_mask & __GFP_WAIT)) {
170         static int count = 0;
171         if (++count < 5) {
172             printk(KERN_ERR "alloc_skb called nonatomically "
173                     "from interrupt %p\n", NET_CALLER(size));
174                     BUG();
175         }
176         gfp_mask &= ~__GFP_WAIT;
177     }
178
179     /* Get the HEAD */
180     skb = skb_head_from_pool();
181     if (skb == NULL) {
182         skb = kmem_cache_alloc(skbuff_head_cache, gfp_mask &
                    ~__GFP_DMA);
183         if (skb == NULL)
184             goto nohead;
185     }
186
187     /* Get the DATA. Size must match skb_add_mtu(). */
188     size = SKB_DATA_ALIGN(size);
189     data = kmalloc(size + sizeof(struct skb_shared_info),
                    gfp_mask);
190     if (data == NULL)
191         goto nodata;
```

When head and data have been successfully allocated, the head is initialized.

```
192
193             /* XXX: does not include slab overhead */
194         skb->truesize = size + sizeof(struct sk_buff);
195
196             /* Load the data pointers. */
197         skb->head = data;
198         skb->data = data;
199         skb->tail = data;
200         skb->end = data + size;
201
202             /* Set up other state */
203         skb->len = 0;
204         skb->cloned = 0;
205         skb->data_len = 0;
206
207         atomic_set(&skb->users, 1);
208         atomic_set(&(skb_shinfo(skb)->dataref), 1);
209         skb_shinfo(skb)->nr_frags = 0;
210         skb_shinfo(skb)->frag_list = NULL;
211         return skb;
212
213 nodata:
214         skb_head_to_pool (skb);
215 nohead:
216         return NULL;
217 }
218
219
```

On return to __dev_alloc_skb() space is reserved for the MAC header.

```
1035        if (skb)
1036            skb_reserve(skb,16);
1037        return skb;
1038 }
```

skb_reserve moves the data and tail pointers to point to first byte after the 16 bytes of headroom.

```
911 static inline void skb_reserve(struct sk_buff *skb,
                                    unsigned int len)
912 {
913     skb->data+=len;
914     skb->tail+=len;
915 }
```

The skb_put() function can be used to update the len and tail values after data has been placed in the sk_buff(). The actual filling of the buffer is most commonly performed by a DMA transfer.

```
786 static inline unsigned char *skb_put(struct sk_buff
                          *skb, unsigned int len)
787 {
788     unsigned char *tmp=skb->tail;
789     SKB_LINEAR_ASSERT(skb);
790     skb->tail += len;
791     skb->len += len;
792     if(skb->tail>skb->end) {
793         skb_over_panic(skb,len,current_text_addr());
794     }
795     return tmp;
796 }
```

Non−linear sk_buffs are those consisting of unmapped page buffers and additional chained struct sk_buffs. A non−zero value of data_len is an indicator of non−linearity. For obvious reasons the simple skb_put() function neither supports nor tolerates non−linearity. SKB_LINEAR_ASSERT checks value of data_len through function skb_is_nonlinear. A non−zero value results in an error message to be logged by BUG.

```
761 #define SKB_LINEAR_ASSERT(skb)
        do { if (skb_is_nonlinear(skb)) BUG(); } while (0)
```

skb_is_nonlinear is defined as below.

```
749 static inline int skb_is_nonlinear(const struct
                          sk_buff *skb)
750 {
751     return skb->data_len;
752 }
```

Queuing the packet with netif_rx()

The netif_rx() function is defined in net/core/dev.c. It typically runs in the context of the hardware interrupt that signalled the completion of the DMA transfer. Its function is to queue the sk_ buff for processing by network layer. The buffer may, however, be dropped during processing for congestion control. After queuing the packet, netif_rx() raises the NET_RX_SOFTIRQ. The bulk of the processing of an input packet is done in the context of this softirq by the net_rx_action() function. The netif_rx() function returns one of the following values which are defined in include/linux/netdevice.h.

```
55 /* Backlog congestion levels */
56 #define NET_RX_SUCCESS  0 /* keep 'em coming, baby */
57 #define NET_RX_DROP     1 /* packet dropped */
58 #define NET_RX_CN_LOW   2 /* storm alert, just in case */
59 #define NET_RX_CN_MOD   3 /* Storm on its way! */
60 #define NET_RX_CN_HIGH  4 /* The storm is here */
61 #define NET_RX_BAD      5  /* packet dropped due to
                               kernel error */
```

Incoming packets are placed on per–cpu queues so that no locking is needed. The softnet_data array, defined in include/linux/netdevice.h.,  consists of a struct softnet_data for each CPU.

```
473 struct softnet_data
474 {
475     int                     throttle;
476     int                     cng_level;
477     int                     avg_blog;
478     struct sk_buff_head     input_pkt_queue;
479     struct net_device       *output_queue;
480     struct sk_buff          *completion_queue;
481 } __attribute__((__aligned__(SMP_CACHE_BYTES)));

484 extern struct softnet_data softnet_data[NR_CPUS];

 97 struct sk_buff_head {
 98     /* These two members must be first. */
 99     struct sk_buff  * next;
100     struct sk_buff  * prev;
101
102     __u32           qlen;
103     spinlock_t      lock;
104 };
```

These are the congestion management parameters.

```
1073 int netdev_max_backlog = 300;
1074 /* These numbers are selected based on intuition and some
1075  * experimentatiom, if you have more scientific way
1076  * please go ahead and fix things.
1077  */
1078 int no_cong_thresh = 10;
1079 int no_cong = 20;
1080 int lo_cong = 100;
1081 int mod_cong = 290;
1082
```

```
1214 int netif_rx(struct sk_buff *skb)
1215 {
1216     int this_cpu = smp_processor_id();
1217     struct softnet_data *queue;
1218     unsigned long flags;
```

If the device driver has not already time stamped the packet, it is done here.

```
1220     if (skb->stamp.tv_sec == 0)
1221         do_gettimeofday(&skb->stamp);
```

The local variable queue is set to point to the struct sofnet_data for this cpu.

```
1223     /*  The code is rearranged so that the path is
            the most short when CPU is congested, but is
            still operating.
1225     */
1226     queue = &softnet_data[this_cpu];
1227
```

Interrupts are disabled on this CPU while the packet is queued.

```
1228     local_irq_save(flags);
1229
1230     netdev_rx_stat[this_cpu].total++;
```

The length of the input packet queue is compared against its maximum backlog. If the queue is full, the sk_buff is discarded. The value of netdev_max_backlog is declared to be 300 packets in net/core/dev.c.

```
1231        if (queue->input_pkt_queue.qlen <= netdev_max_backlog) {
```

The following compound if first tests to see if the input queue is not empty. If the queue is not empty, then the throttle flag is tested to see if the packet should be dropped. The throttle flag indicates the presence (1) or absence (0) of congestion. If congestion is present, the sk_buff is discarded.

```
1232            if (queue->input_pkt_queue.qlen) {
1233                if (queue->throttle)
1234                    goto drop;
```

If the throttle flag is not set, the sk_buff is added to the input packet queue. Since the if statement above found qlen > 0, the queue is guaranteed to be non−empty here.

```
1236 enqueue:
1237                dev_hold(skb->dev);
1238                __skb_queue_tail(&queue->input_pkt_queue,
                            skb);
```

The cpu_raise_softirq() function sets a flag to indicate that the NET_RX_SOFTIRQ software interrupt is now pending. Most of the actual work of handling the packet will take place in the context of the softirq.

```
1239        /*   Runs from irqs or BH's, no need to wake BH */
1240                cpu_raise_softirq(this_cpu, NET_RX_SOFTIRQ);
1241                local_irq_restore(flags);
```

After raising the softirq, netif_rx() returns the congestion level from softnet_data structure. The get_sample_stats() function in net/core/dev.c sets the congestion level.

```
1242 #ifndef OFFLINE_SAMPLE
1243                get_sample_stats(this_cpu);
1244 #endif
1245                return softnet_data[this_cpu].cng_level;
1246            }
```

If we reach this point in netif_rx(), the input packet queue is empty. The throttle flag, if set, is cleared.

```
1248            if (queue->throttle) {
1249                queue->throttle = 0;
1250 #ifdef CONFIG_NET_HW_FLOWCONTROL
1251                if (atomic_dec_and_test
                               (&netdev_dropping))
1252                    netdev_wakeup();
1253 #endif
1254            }
```

Here we jump back to the code that queue the sk_buff and raises the soft irq.

```
1255            goto enqueue;
1256        }
```

The if block that began at line 1231 ends here. If control reaches this point the input packet queue is full of sk_buffs. The throttle flag is set to indicate congestion. Note that once the queue becomes throttled it must drain completely to become ''unthrottled''.

```
1258    if (queue->throttle == 0) {
1259        queue->throttle = 1;
1260        netdev_rx_stat[this_cpu].throttled++;
1261 #ifdef CONFIG_NET_HW_FLOWCONTROL
1262        atomic_inc(&netdev_dropping);
1263 #endif
1264    }
```

The sk_buff is discarded here.

```
1266 drop:
1267    netdev_rx_stat[this_cpu].dropped++;
1268    local_irq_restore(flags);
1269
1270    kfree_skb(skb);
1271    return NET_RX_DROP;
1272 }
```

This is the end of netif_rx. We now turn our attention to the softirq.

Softirqs

In early versions of Linux processing of received packets took place in the context of what was called a bottom half. The softirq mechanism, which was designed to replace the bottom half was introduced in kernel 2.4. The primary advantage of the softirq mechanism is that multiple soft irqs may run concurrently on multiple processes. Bottom halves were permitted to run only on one CPU at a time.

Further handling of our packet is done in the network receive softirq (NET_RX_SOFTIRQ) which is called from kernel/softirq.c:do_softirq().

Recall that the function netdev_init() which runs at boot time registered two softirq handlers.

```
2865      open_softirq(NET_TX_SOFTIRQ, net_tx_action, NULL);
2866      open_softirq(NET_RX_SOFTIRQ, net_rx_action, NULL);
```

Therefore when netif_rx() executed the line of code shown below it indirectly causes the function net_rx_action() to be executed in a softirq context.

```
1240              cpu_raise_softirq(this_cpu, NET_RX_SOFTIRQ);
```

The  do_softirq() function is defined in kernel/softirq.c. It invokes the appropriate action handler for
each softirq raised.

```
61 asmlinkage void do_softirq()
62 {
63     int cpu = smp_processor_id();
64     __u32 pending;
65     long flags;
66     __u32 mask;
67
68     if (in_interrupt())
69         return;
70
71     local_irq_save(flags);
```

The variable pending is a bit mask that indicates which of the possible 32 soft irqs are presently
pending.  It is set to irq_stat[cpu].__softirq_pending and mask  is set to complement  of this value.

```
73     pending = softirq_pending(cpu);
74
75     if (pending) {
76         struct softirq_action *h;
77
78         mask = ~pending;
79         local_bh_disable();
```

local_bh_disable is defined in include/asm−i386/softirq.h. It is replaced by a call to cpu_bh_disable.
local_bh_count  is a macro which translates to  irq_stat[cpu].__local_bh_count.

```
10 #define cpu_bh_disable(cpu) \
11     do { local_bh_count(cpu)++; barrier(); } while (0)
12
13 #define local_bh_disable()
                       cpu_bh_disable(smp_processor_id())
```

Now in do_softirq, clear irq_stat[cpu].__softirq_pending.

```
80 restart:
81    /* Reset the pending bitmask before enabling irqs */
82            softirq_pending(cpu) = 0;
83
84            local_irq_enable();
```

"h" is set to point to first element in softirq_vec array.

```
86            h = softirq_vec;
87
88            do {
```

Softirqs are checked in order of their priority (HI_SOFTIRQ, NET_TX_SOFTIRQ ...)
and the respective function handler is called. In the case of NET_RX_SOFTIRQ, it is net_rx_action.

```
89                    if (pending & 1)
90                            h->action(h);
```

"h" is now set to point to next element in softirq_vec array.

```
91                    h++;
92                    pending >>= 1;
93            } while (pending);
94
95            local_irq_disable();
```

If new softirqs (other than those handled above) have been raised, they are handled as well. Recall
that mask was originally set to the complement of pending. The masking here is presumably in
place to prevent livelock type conditions.

```
97            pending = softirq_pending(cpu);
98            if (pending & mask) {
99                    mask &= ~pending;
100                   goto restart;
101           }
```

__local_bh_enable is defined in include/asm−i386/softirq.h.

```
 8 #define __cpu_bh_enable(cpu) \
 9     do { barrier(); local_bh_count(cpu)--; } while (0)
14 #define __local_bh_enable() \
        __cpu_bh_enable(smp_processor_id())

102             __local_bh_enable();
103
104         if (pending)
105             wakeup_softirqd(cpu);
106     }
107
108     local_irq_restore(flags);
109 }
```

Received packet handling in the softirq.

The net_rx_action() function resides in net/core/dev.c and was previously shown to have been installed as the handler for the NET_RX_SOFTIRQ. As might be expected, its mission is to consume packets from the queue that netif_rx() produces to and then to pass them on to the proper handler.

```
1419 static void net_rx_action(struct softirq_action *h)
1420 {
1421     int this_cpu = smp_processor_id();
```

A unique structure of type struct softnet_data is associated with each CPU for the purpose of managing input and output queues at the interface between the protocols and the device driver. Here queue is initalized to point to softnet_data structure for this CPU.

```
1422     struct softnet_data *queue =
                        &softnet_data[this_cpu];
1423     unsigned long start_time = jiffies;
1424     int bugdet = netdev_max_backlog;
1425
1426     br_read_lock(BR_NETPROTO_LOCK);
1427
```

This loop is executed until either the input queue has been emptied, at least one jiffy of CPU time has been consumed, or the value of bugdet becomes negative.

```
1428        for (;;) {
1429               struct sk_buff *skb;
1430               struct net_device *rx_dev;
```

Attempt to dequeue an sk_buff from this CPU's input packet queue. Because the queue is local to the CPU local disabling of interrupts provides safe serialization.

```
1432               local_irq_disable();
1433               skb = __skb_dequeue(&queue->input_pkt_queue);
1434               local_irq_enable();
```

Return if the input packet queue is empty.

```
1436               if (skb == NULL)
1437                   break;
```

The net_device pointer (which was set by the device driver) is potentially adjusted here.

```
1439               skb_bond(skb);
```

The skb_bond() function is defined in net/core/dev.c. It assigns the sk_buff to the master device for present device if such exists. We really don't understand device groups and master devices!

```
      /*  Reparent skb to master device. This function is
          called only from net_rx_action under
          BR_NETPROTO_LOCK.
          It is misuse of BR_NETPROTO_LOCK, but it is OK for
          now.
       */
1314 static __inline__ void skb_bond(struct sk_buff *skb)
1315 {
1316     struct net_device *dev = skb->dev;
1317
1318     if (dev->master) {
1319         dev_hold(dev->master);
1320         skb->dev = dev->master;
1321         dev_put(dev);
1322     }
1323 }
```

Back in net_rx_action the net_device pointer is copied to a local variable.

```
1441              rx_dev = skb->dev;
1442
```

CONFIG_NET_FASTROUTE is an option to allow direct NIC–to–NIC data transfer on a local network. We do will ignore it for now.

```
1443 #ifdef CONFIG_NET_FASTROUTE
1444          if (skb->pkt_type == PACKET_FASTROUTE)  {
1445              netdev_rx_stat[this_cpu].
                      fastroute_deferred_out++;
1446              dev_queue_xmit(skb);
1447              dev_put(rx_dev);
1448              continue;
1449          }
1450 #endif
```

Link level demultiplexing

Here skb->data points to the start of the data area of the buffer. Therefore h.raw and nh.raw are both being set to point to the MAC header.

```
1451              skb->h.raw = skb->nh.raw = skb->data;
```

Here is where link layer demultiplexing takes place. For DIX framing the key to demultiplexing is the standard packet type field that is carried in the MAC header.

```
39 #define ETH_P_LOOP      0x0060    /* Ethernet Loopback packet    */
40 #define ETH_P_PUP       0x0200    /* Xerox PUP packet            */
41 #define ETH_P_PUPAT     0x0201    /* Xerox PUP Addr Trans packet */
42 #define ETH_P_IP        0x0800    /* Internet Protocol packet    */
43 #define ETH_P_X25       0x0805    /* CCITT X.25                  */
44 #define ETH_P_ARP       0x0806    /* Address Resolution packet   */
45 #define ETH_P_BPQ       0x08FF    /* G8BPQ AX.25 Ethernet Packet */
46 #define ETH_P_IEEEPUP   0x0a00    /* Xerox IEEE802.3 PUP packet */
47 #define ETH_P_IEEEPUPAT 0x0a01    /* Xerox IEEE802.3 PUP Addr Trans pkt*/
```

For 802.3 life is a bit more complicated:

```
73 #define ETH_P_802_3     0x0001    /* Dummy type for 802.3 frames  */
74 #define ETH_P_AX25      0x0002    /* Dummy protocol id for AX.25  */
75 #define ETH_P_ALL       0x0003    /* Every packet (be careful!!!) */
76 #define ETH_P_802_2     0x0004    /* 802.2 frames                 */
77 #define ETH_P_SNAP      0x0005    /* Internal only                */
78 #define ETH_P_DDCMP     0x0006    /* DEC DDCMP: Internal only     */
```

For all types this routine depends upon the device driver to have extracted the appropriate type from the MAC header and stored in in skb->protocol in host byte order.   A convenience function is provided to the device driver as shown in the following extract from 3c59x.c

```
2419                skb->protocol = eth_type_trans(skb, dev);
```

```
152 /*
153 * Determine packet's protocol ID. The rule here is that we
154 * assume 802.3 if type field is short enough to be a length.
155 * This is normal and works for any 'now in use' protocol.
156 */
157
158 unsigned short eth_type_trans(struct sk_buff *skb,
         struct net_device *dev)
159 {
160     struct ethhdr *eth;
161     unsigned char *rawp;
162
```

The call to skb_pull() advances skb->data so that it points to the network layer header (or the IEEE 802.2 LLC header for 802.2/3 framing,  and decrements skb->len  by the length of the MAC header (hard_header_len);

```
163     skb->mac.raw = skb->data;
164     skb_pull(skb, dev->hard_header_len);
165     eth = skb->mac.ethernet;
166
```

If the low order bit of the high order byte of the MAC address is 1, then this packet is a broadcast or a multicast.

```
167      if (*eth->h_dest&1)
168      {
169              if(memcmp(eth->h_dest,dev->broadcast, ETH_ALEN)==0)
170                      skb->pkt_type=PACKET_BROADCAST;
171              else
172                      skb->pkt_type=PACKET_MULTICAST;
173      }
174
175 /*
176 *        This ALLMULTI check should be redundant by 1.4
177 *        so don't forget to remove it.
178 *
179 *        Seems, you forgot to remove it. All silly devices
180 *        seems to set IFF_PROMISC.
181 */
182
183      else if(1 /*dev->flags&IFF_PROMISC*/)
184      {
185              if(memcmp(eth->h_dest,dev->dev_addr, ETH_ALEN))
186                      skb->pkt_type = PACKET_OTHERHOST;
187      }
188
```

The two byte field immediately following the destination MAC address is the packet type for DIX framing but it is the packet length for IEEE 802.2/3 framing.   For IP, ARP, RARP, and IPX the packet type is at least 0x800 which is 2048 and thus larger than the maximum frame size.   It does look like something really ugly could ensue here if jumbo frames were used in conjunction with 802.2/3 framing.

```
189      if (ntohs(eth->h_proto) >= 1536)
190              return eth->h_proto;
191
192      rawp = skb->data;
193
194 /*
195 *This is a hack to spot IPX packets. Older Novell  breaks
196 *the proto and runs IPX over 802.3 without an 802.2 LLC
197 *layer. We look for FFFF which isn't a used 802.2 SSAP/DSAP.
198 *This won't work for fault tol netware but does for the rest.
199 */
200      if (*(unsigned short *)rawp == 0xFFFF)
201              return htons(ETH_P_802_3);
202
```

For "real" 802.2/3 framing, the length field is followed by the 802.2 LLC header containing the DSAP,SSAP, and cntl fields which are normally set to 0xaa, 0xaa, 0x03.   This is followed by the 802.2 SNAP header which contains a 3 byte originator code and finally the 2 byte type field.   This module just returns the code for 802_2 in that case and leaves it to the 802.2 module to eventually perform the demultiplexing.

```
203 /*
204  *        Real  802.2 LLC
205  */
206      return htons(ETH_P_802_2);
207 }
```

Recall that protocol packet handlers register themselves by filling in the packet_type structure and passing it to dev_add_pack() where the structure is placed on an appropriate chain.

```
421 struct packet_type
422 {
423      unsigned short  type;    /* really htons(ether_type).*/
424      struct net_device *dev; /* NULL is wildcarded here  */
425      int    (*func) (struct sk_buff *, struct net_device *,
426                         struct packet_type *);
427      void    *data;         /* Private to the packet type */
428      struct packet_type *next;
429 };
430
```

This is the block in net_rx_action in which the packet processing occurs. Protocols which wish to receive all incoming packets are linked into a list pointed to by ptype_all. These protocols register the have type ETH_P_ALL and are processed before considering the protocols that consume only a specific packet type.

```
1452            {
1453                    struct packet_type *ptype, *pt_prev;
1454                    unsigned short type =  skb->protocol;
```

pt_prev is initialized to NULL. We try to match sk buff against each protocol, registered with ptype_all.

```
1456                    pt_prev = NULL;
1457                    for (ptype = ptype_all; ptype; ptype =
                            ptype->next) {
```

Even though every packet handler in this chain says it wants to see all packets, it can also say that it wants to limit the packets to those received on a specific device.    If ptype->dev is null, then any device is acceptable.   A value of 0 in the data field of the packet_type structure indicates that this is an old protocol that does not understand shared sk_buffs. We don't really grasp what the oddball use of pt_prev is all about, but possibly it is trying to deal with the necessity of cloning an skb that is to be consumed by an old style protocol.   It should be necessary to clone if and only if there is more than one protocol interested.

```
1458                    if (!ptype->dev || ptype->dev==skb->dev){
1459                        if (pt_prev) {
1460                            if  (!pt_prev->data) {
1461                                deliver_to_old_ones
                                        (pt_prev, skb,0);
1462                            } else {
1463                                atomic_inc(&skb->users);
1464                                pt_prev->func(skb,
1465                                    skb->dev,
1466                                    pt_prev);
1467                            }
1468                        }
```

pt_prev is set to matched protocol. The protocol specific function is called when next  match is found.

```
1469                        pt_prev = ptype;
1470                    }
1471            }
```

18

We are also not interested in diverters at the moment.

```
1473 #ifdef CONFIG_NET_DIVERT
1474              if (skb->dev->divert &&
                         skb->dev->divert->divert)
1475              handle_diverter(skb);
1476 #endif /* CONFIG_NET_DIVERT */
```

CONFIG_BRIDGE is an option to enable ethernet bridging that we are also not considering.

```
1479 #if defined(CONFIG_BRIDGE) ||
            defined(CONFIG_BRIDGE_MODULE)
1480              if (skb->dev->br_port != NULL &&
1481                  br_handle_frame_hook != NULL){
1482              handle_bridge(skb, pt_prev);
1483              dev_put(rx_dev);
1484              continue;
1485          }
1486 #endif
```

This is the point at which protocols wanting only specific packet types are processed. Recall that the packet_type structures of these protocols are placed in a hash table with 16 hash lists. The low order 4 bits of the protocol type is the hash key. The loop processes the hash list corresponding to the protocol type of the current packet.

```
1488              for (ptype=ptype_base[ntohs(type)&15];
                      ptype;ptype=ptype->next) {
```

Test to see if the type of the packet matches the type registered in the packet_type structure. If so, and the protocol also registered a specific struct netdevice then it is necessary to see if the input device matches as well.

```
1489                  if (ptype->type == type &&
1490                      (!ptype->dev ||
                          ptype->dev == skb->dev))
                      {
1491                  if (pt_prev) {
1492                      if (!pt_prev->data)
1493                          deliver_to_old_ones
                                  (pt_prev,
                                   skb, 0);
1494                      else {
1495                          atomic_inc
                                  (&skb->users);
```

Depending on the protocol type, the appropriate handler function is called. It is ip_rcv for ETH_P_IP and arp_rcv for ETH_P_ARP.

```
1496                                    pt_prev->func(skb,
1497                                             skb->dev,
1498                                             pt_prev);
1499                        }
1500                   }
1501                   pt_prev = ptype;
1502              }
1503         }
```

On exit from the loop invoke the function handler for the last matched protocol.

```
1505         if (pt_prev) {
1506              if (!pt_prev->data)
1507                   deliver_to_old_ones(pt_prev,
1508                                       skb, 1);
                 else
1509                   pt_prev->func(skb, skb->dev,
                             pt_prev);
```

If no protocol was matched , pt_prev is NULL and the sk_buff is discarded.

```
1510                        } else
1511                            kfree_skb(skb);
1512                }
1513
```

The dev_put() function decrements the value of  rx_dev->refcount.   This was incremented back in netif_rx() and preserved across the scheduling of the softirq.    But what if the skb_bond() call changed it ...  All is well skb_bond also released and reallocated.

```
1514                dev_put(rx_dev);
```

If more than one jiffy has elapsed while consuming sk_buffs  or netdev_max_backlog buffers have been consumed exit the loop.  Otherwise, continue dequeueing sk_buffs.

```
1516                if (bugdet-- < 0 || jiffies - start_time > 1)
1517                    goto softnet_break;
```

This option enables NIC (Network Interface Card) hardware throttling during periods of extremal congestion. At the moment only a couple of device drivers support it.

```
1519 #ifdef CONFIG_NET_HW_FLOWCONTROL
1520                if (queue->throttle && queue->input_pkt_queue.qlen
                            < no_cong_thresh ) {
1521                    if (atomic_dec_and_test
                                (&netdev_dropping)) {
1522                        queue->throttle = 0;
1523                        netdev_wakeup();
1524                        goto softnet_break;
1525                    }
1526                }
1527 #endif
1528
1529        }
1530        br_read_unlock(BR_NETPROTO_LOCK);
```

Reaching this point means that the for loop was exited via the break at line 1437 and the input packet queue has been completely emptied. The throttle flag is cleared, if set.

```
1532        local_irq_disable();
1533        if (queue->throttle) {
1534            queue->throttle = 0;
1535 #ifdef CONFIG_NET_HW_FLOWCONTROL
1536            if (atomic_dec_and_test(&netdev_dropping))
1537                netdev_wakeup();
1538 #endif
1539        }
1540        local_irq_enable();
1541
1542        NET_PROFILE_LEAVE(softnet_process);
1543        return;
```

Reaching this point means that the loop was exited via the goto at line 1517 which is triggered by the jiffy count. The interrupt is raised again since there are sk_buffs remaining to be processed.

```
1545 softnet_break:
1546        br_read_unlock(BR_NETPROTO_LOCK);
1547
1548        local_irq_disable();
1549        netdev_rx_stat[this_cpu].time_squeeze++;
1550        /*   This already runs in BH context, no need to
                 wake up BH's */
1551        cpu_raise_softirq(this_cpu, NET_RX_SOFTIRQ);
1552        local_irq_enable();
```

NET_PROFILE_LEAVE has effect only when network code profiler is configured.

```
1554        NET_PROFILE_LEAVE(softnet_process);
1555        return;
1556 }
```

The deliver_to_old_ones() function is defined in net/core/dev.c. It invokes the handler function for old protocols that do not understand shared sk_buffs.

```
      /*    Deliver skb to an old protocol, which is not
            threaded well or which do not understand shared
            skbs.
      */
1277 static int deliver_to_old_ones(struct packet_type *pt,
                          struct sk_buff *skb, int last)
1278 {
1279      static spinlock_t net_bh_lock = SPIN_LOCK_UNLOCKED;
1280      int ret = NET_RX_DROP;
1281
```

The value of last is one if and only if this is the last protocol to which the packet must be delivered. In this case it is not necessary to create a clone of the sk_buff.

```
1283      if (!last) {
1284          skb = skb_clone(skb, GFP_ATOMIC);
1285          if (skb == NULL)
1286              return ret;
1287      }
```

An sk_buff which contains data in unmapped page sections is made linear by skb_linearize. A linear sb_buff is one which consists of a fixed length header of type struct sk_buff and a single kmalloc'ed data buffer.

```
1288      if (skb_is_nonlinear(skb) &&
                  skb_linearize(skb, GFP_ATOMIC) != 0) {
1289          kfree_skb(skb);
1290          return ret;
1291      }
1292
```

Here several hacks are inserted to provide the expected environment to the old protocols. We do not consider these in detail.

```
1293  /* The assumption (correct one) is that old
            protocols did not depened on BHs different of
            NET_BH and TIMER_BH. */
1296
1297   /* Emulate NET_BH with special spinlock */
1298      spin_lock(&net_bh_lock);
```

All timers are disabled due to above assumption.

```
1300  /* Disable timers and wait for all timers completion */
1301      tasklet_disable(bh_task_vec+TIMER_BH);
```

The protocol specific function is invoked here. Timers are enabled thereafter.

```
1303        ret = pt->func(skb, skb->dev, pt);
1305        tasklet_hi_enable(bh_task_vec+TIMER_BH);
1306        spin_unlock(&net_bh_lock);
1307        return ret;
1308 }
```