

# Navigation of TCP/IP files in Linux

Divye Kapoor  
Pracheer Agarwal  
Swagat Konchada

# Background Information

# Linux Virtual Filesystem (VFS)

- It is the software layer in the kernel that provides a uniform filesystem interface to userspace programs
- It provides an abstraction within the kernel that allows for transparent working with a variety of filesystems.
- Thus it allows many different filesystem implementations to coexist freely
- Each socket is implemented as a “file” mounted on the sockfs filesystem.
  - file->private points to the socket information.

# Inodes and File Structures

- Inodes provide a method to access the actual data blocks allocated to a file. For sockets, they provide buffer space which can be used to hold socket specific data.
  - struct inode
- Every file is represented in the kernel as an object of the *file* structure. It requires an inode provided to it.
  - struct file

# Structure of Function Pointers

Struct operations {

int (\*read)(int, char \*, int);

void (\*destroy\_inode)(inode \*);

void (\*dirty\_inode) (struct inode \*);

int (\*write\_inode) (struct inode \*, int);

void (\*drop\_inode) (struct inode \*);

void (\*delete\_inode) (struct inode \*);

};

Sizeof(operations) = sizeof(function ptr)\*6

Divye Kapoor

# Walkthrough of Sending

## User Space

Socket, bind, listen, connect, send, recv, write, read etc.



## Socket Functions (Kernel)

sys\_socket, sys\_bind, sys\_listen, sys\_connect etc. in socket.c



## TCP/IP Layer Functions

inet\_create, tcp\_v4\_connect, tcp\_sendmsg, tcp\_recvmsg



## Ethernet Device Layer

dev\_hard\_start\_xmit



# Socket(family, type, proto)

Sys\_socket()

Sock\_create()

Allocate a socket object  
(internally an inode  
Associated with a file object)

Locate the family requested and  
call the create function for that  
family

Inet\_create()  
Lower layer initialization

Sock\_map\_fd()

Sock\_alloc\_fd()  
Allocate a file descriptor

Sock\_attach\_fd()

Fd\_install()



# Sys\_connect(fd, sockaddr \*, len)

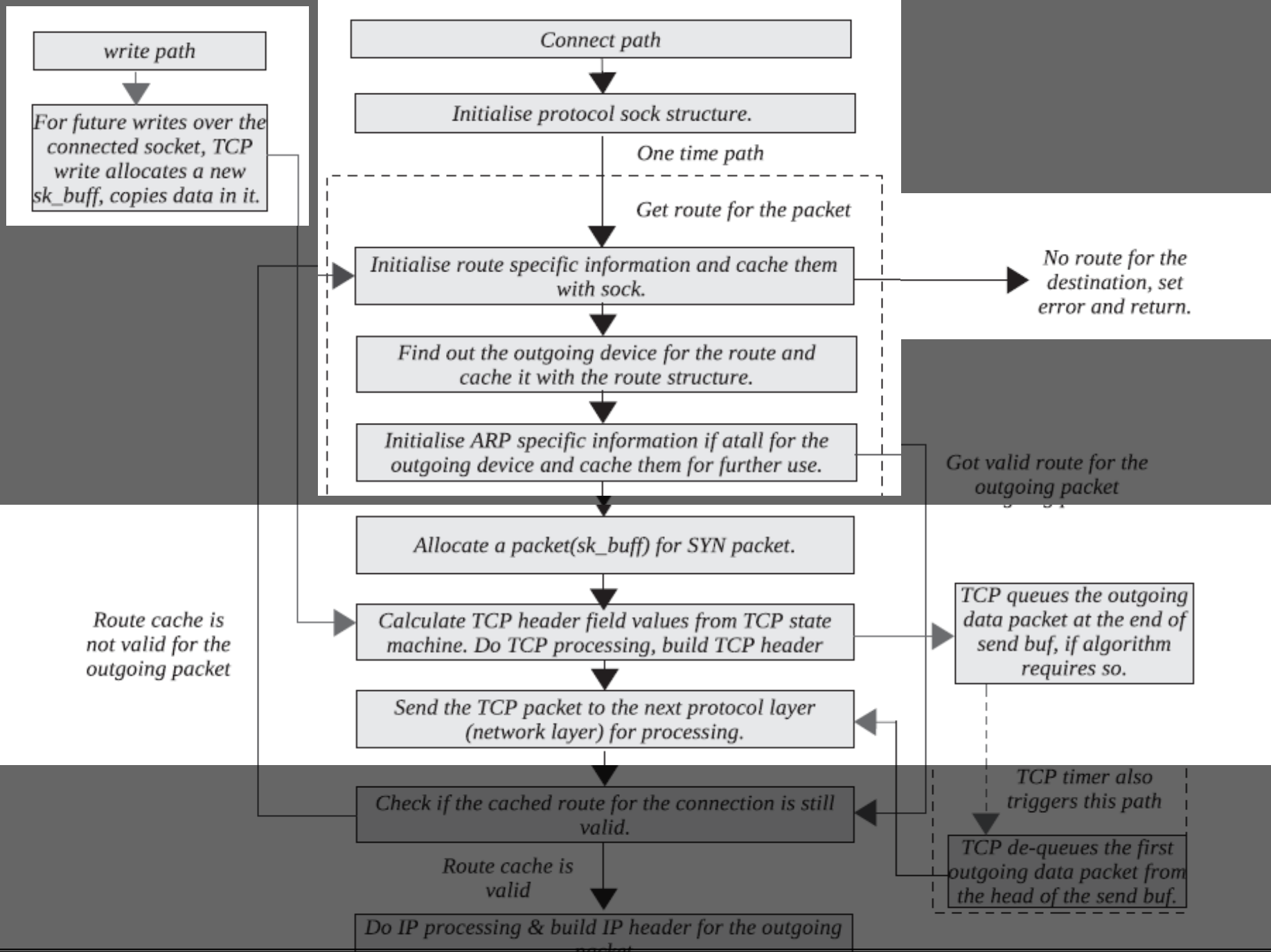
Sys\_connect()

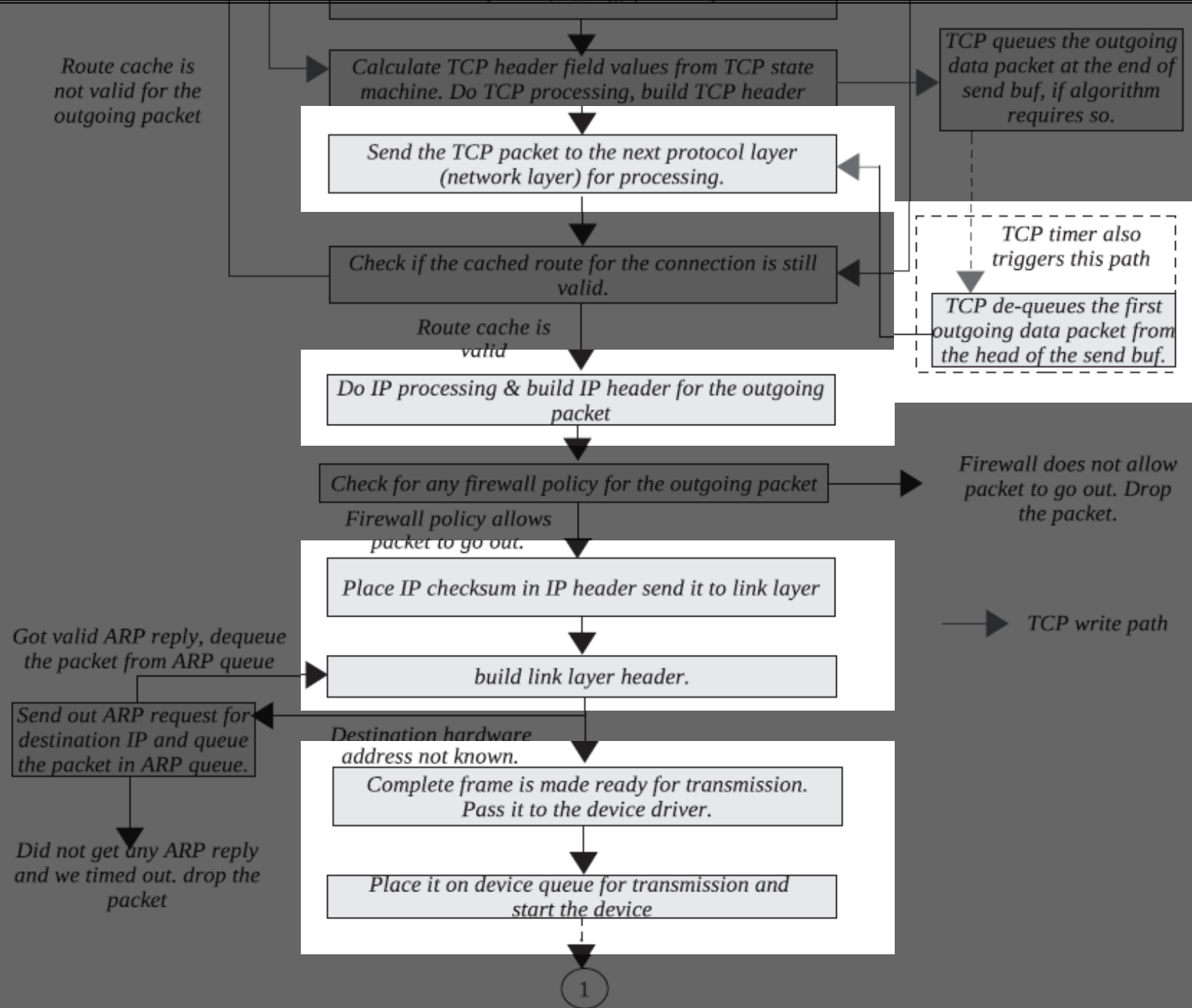
Sockfd\_lookup\_light()  
Returns the socket object  
associated with the given fd

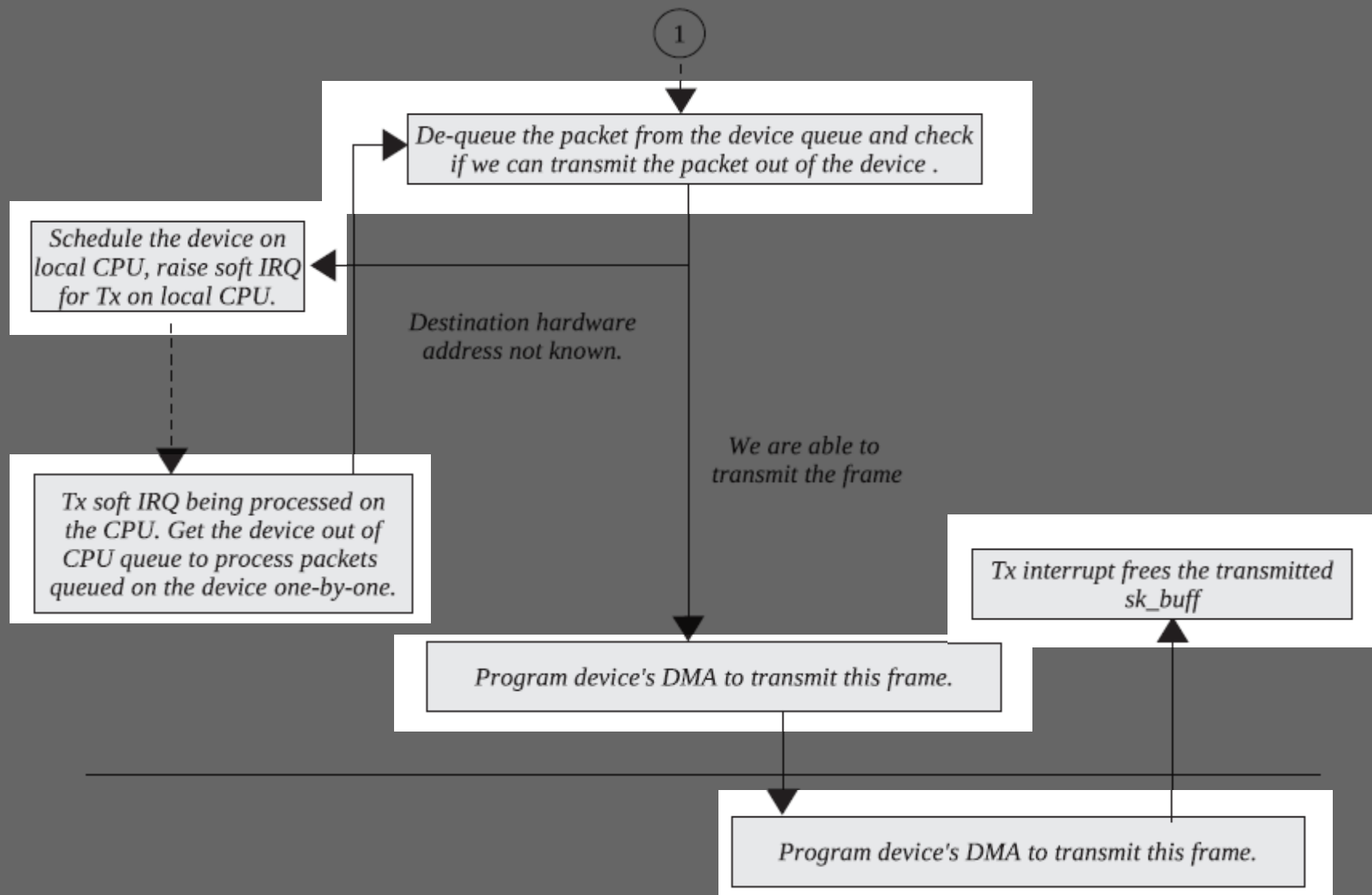
Move\_addr\_to\_kernel()  
For userspace sockaddr \*

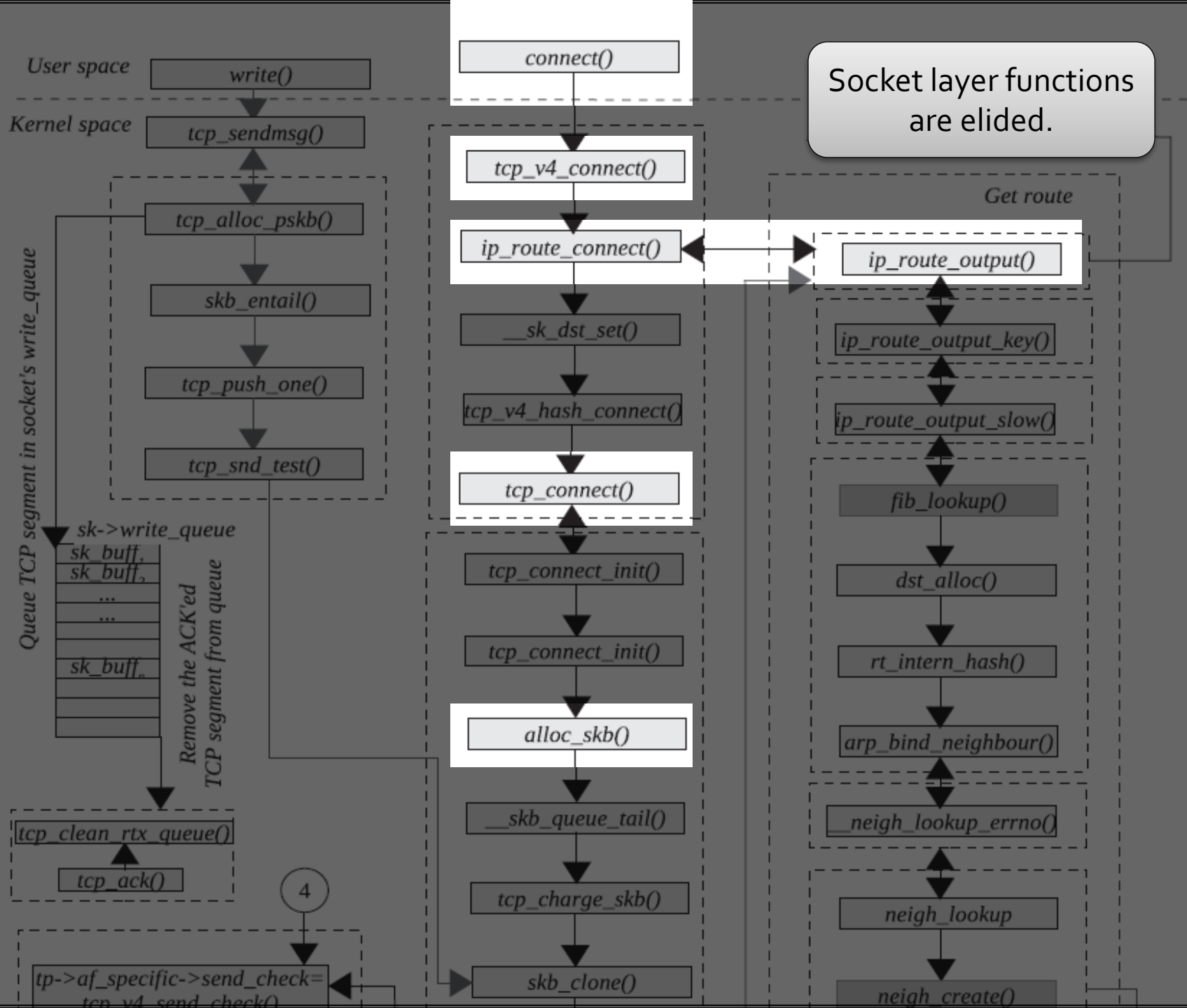
Sock->ops->connect()  
Lower layer call

Tcp\_v4\_connect()









**struct sk\_buff**

# struct sk\_buff

Defined in `<include/linux/skbuff.h>`

- used by every network layer (except the physical layer)
- fields of the structure change as it is passed from one layer to another
- i.e., fields are layer dependent.

# Networking options

```
struct sk_buff {  
    ... ..  
#ifdef CONFIG_NET_SCHED  
    __u32  tc_index;  
#ifdef CONFIG_NET_CLS_ACT  
    __u32  tc_verd;  
    __u32  tc_classid;  
#endif  
#endif  
}
```

sk\_buff is peppered with c preprocessor #ifdef directives.

CONFIG\_NET\_SCHED symbol should be defined at compile time for the structure to have the element tc\_index.

enabled with some version of *make config* by an administrator.



# sk\_buff list

- The kernel maintains all sk\_buff structures in a doubly linked list.

```
struct sk_buff_head { /* only the head of the list */
    /* These two members must be first. */
    struct sk_buff  * next;
    struct sk_buff  * prev;

    __u32    qlen;
    spinlock_t  lock; /* atomicity in accessing a sk_buff list. */
};
```

# Element classification

---

- Layout
- General
- Feature-specific
- Management functions

# Layout

- **struct sock \* sk**

**sock data structure of the socket that owns this buffer**

- unsigned int len

includes both the data in the main buffer (i.e., the one pointed to by head) and the data in the fragments

- unsigned int data\_len

unlike len, data\_len accounts only for the size of the data in the fragments.

- unsigned int truesize

skb->truesize = size + sizeof(struct sk\_buff);

- atomic\_t users

reference count, or the number of entities using this sk\_buff buffer

atomic\_inc and atomic\_dec

# Layout

- `struct sock * sk`

sock data structure of the socket that owns this buffer

- **unsigned int len**

**includes both the data in the main buffer (i.e., the one pointed to by head) and the data in the fragments**

- `unsigned int data_len`

unlike len, data\_len accounts only for the size of the data in the fragments.

- `unsigned int truesize`

`skb->truesize = size + sizeof(struct sk_buff);`

- `atomic_t users`

reference count, or the number of entities using this sk\_buff buffer

`atomic_inc` and `atomic_dec`

# Layout

- `struct sock * sk`  
sock data structure of the socket that owns this buffer
- `unsigned int len`  
includes both the data in the main buffer (i.e., the one pointed to by head) and the data in the fragments
- **`unsigned int data_len`**  
**unlike len, data\_len accounts only for the size of the data in the fragments.**
- `unsigned int truesize`  
`skb->truesize = size + sizeof(struct sk_buff);`
- `atomic_t users`  
reference count, or the number of entities using this `sk_buff` buffer  
`atomic_inc` and `atomic_dec`

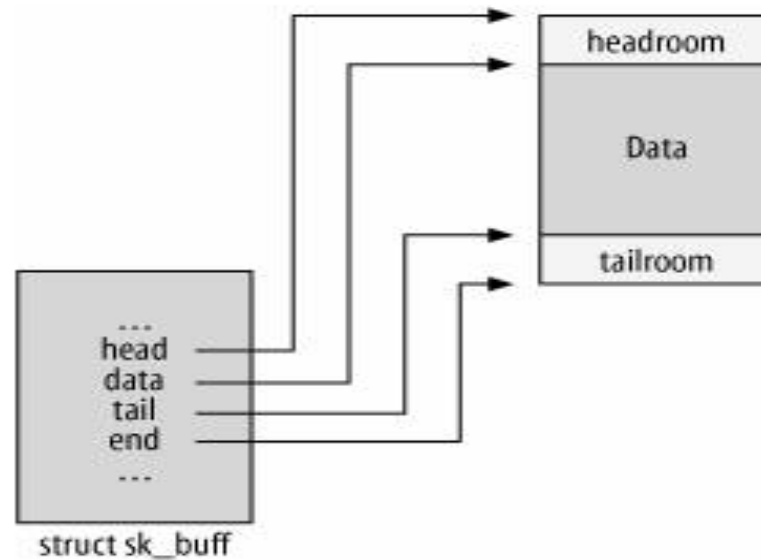
# Layout

- `struct sock * sk`  
sock data structure of the socket that owns this buffer
- `unsigned int len`  
includes both the data in the main buffer (i.e., the one pointed to by `head`) and the data in the fragments
- `unsigned int data_len`  
unlike `len`, `data_len` accounts only for the size of the data in the fragments.
- **`unsigned int truesize`**  
**`skb->truesize = size + sizeof(struct sk_buff);`**
- `atomic_t users`  
reference count, or the number of entities using this `sk_buff` buffer  
`atomic_inc` and `atomic_dec`

# Layout

- `struct sock * sk`  
sock data structure of the socket that owns this buffer
- `unsigned int len`  
includes both the data in the main buffer (i.e., the one pointed to by `head`) and the data in the fragments
- `unsigned int data_len`  
unlike `len`, `data_len` accounts only for the size of the data in the fragments.
- `unsigned int truesize`  
`skb->truesize = size + sizeof(struct sk_buff);`
- **`atomic_t users`**  
reference count, or the number of entities using this `sk_buff` buffer  
`atomic_inc()` and `atomic_dec()`

# position pointers



- `unsigned char *head`
- `sk_buff_data_t end`
- `unsigned char *data`
- `sk_buff_data_t tail`

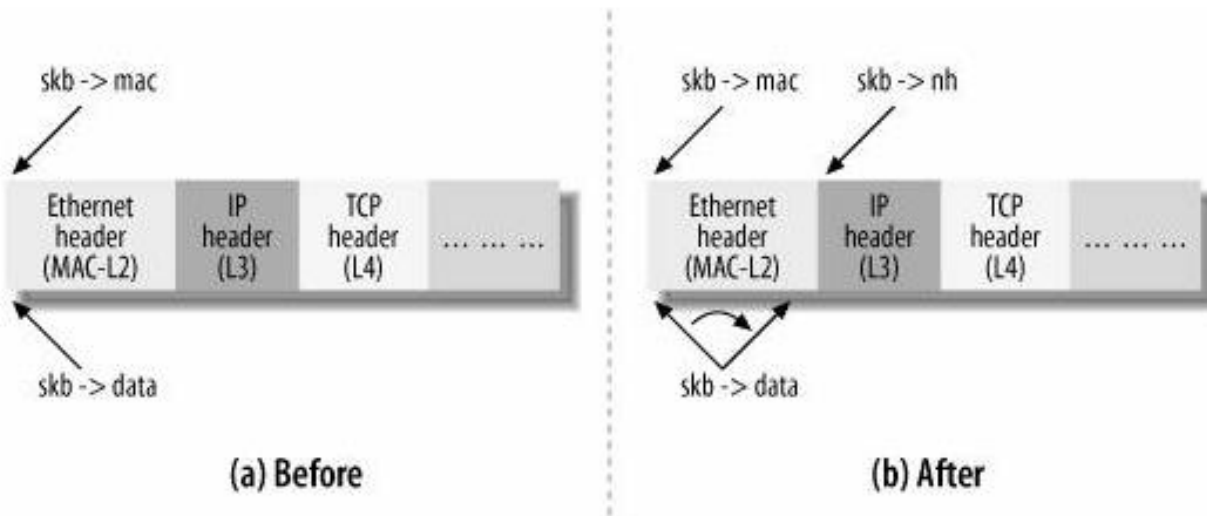


# sk\_buf->dev

struct net\_device \*dev

- represents the receiving interface or the to be transmitted device(or interface) corresponding to the packet.
- usually represents the virtual device's(representation of all devices grouped) net\_device structure.
- Pointers to protocol headers.
- sk\_buff\_data\_t transport\_header;
- sk\_buff\_data\_t network\_header;
- sk\_buff\_data\_t mac\_header;

# pointer modifications



update of data is done using the \*\_header pointers

# Control block

- char cb[40]
- This is a "control buffer," or storage for private information, maintained by each layer for internal use.

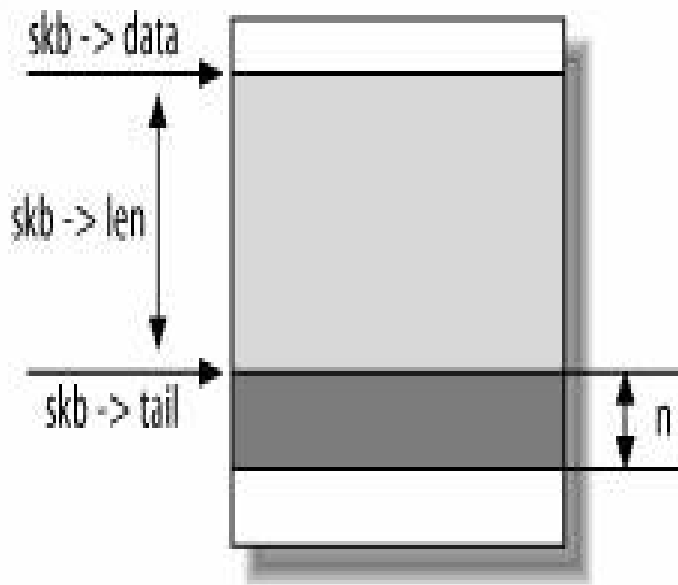
```
struct tcp_skb_cb {  
    ... .. __u32 seq; /* Starting sequence number */  
    __u32 end_seq; /* SEQ + FIN + SYN + datalen*/  
    __u32 when; /* used to compute rtt's */  
    __u8 flags; /* TCP header flags. */  
    ... ..  
};
```

# management functions

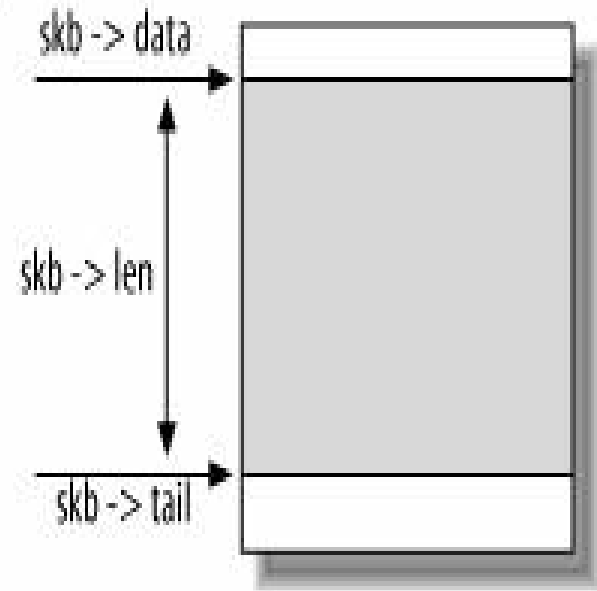
Defined in `<include/linux/skbuff.h>` & `<net/core/skbuff.c>`

`skb_put(struct sk_buff *, unsigned int len)`

(a1)



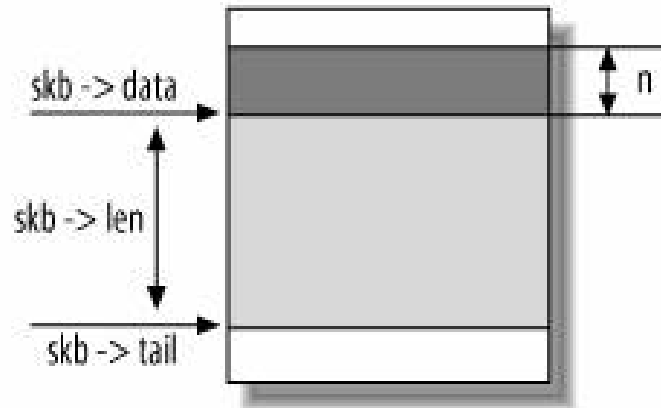
(a2)



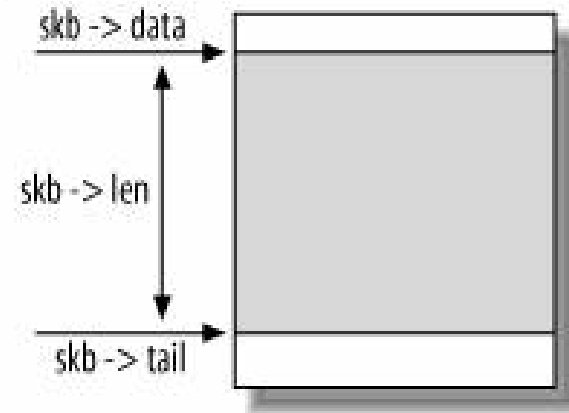
# management functions

`skb_push(struct sk_buff *skb, unsigned int len)`

(b1)

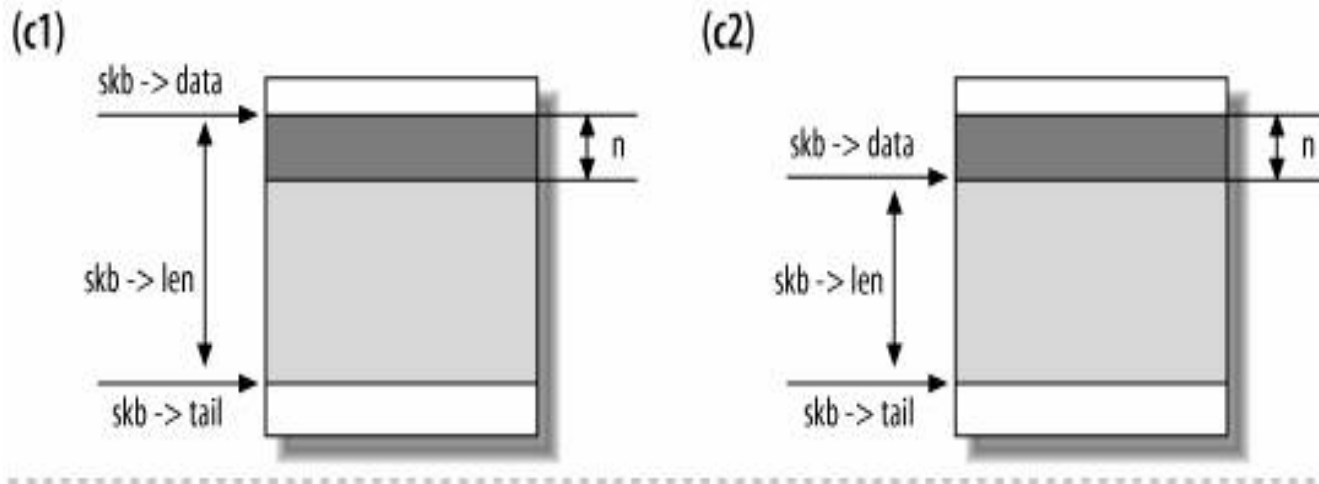


(b2)



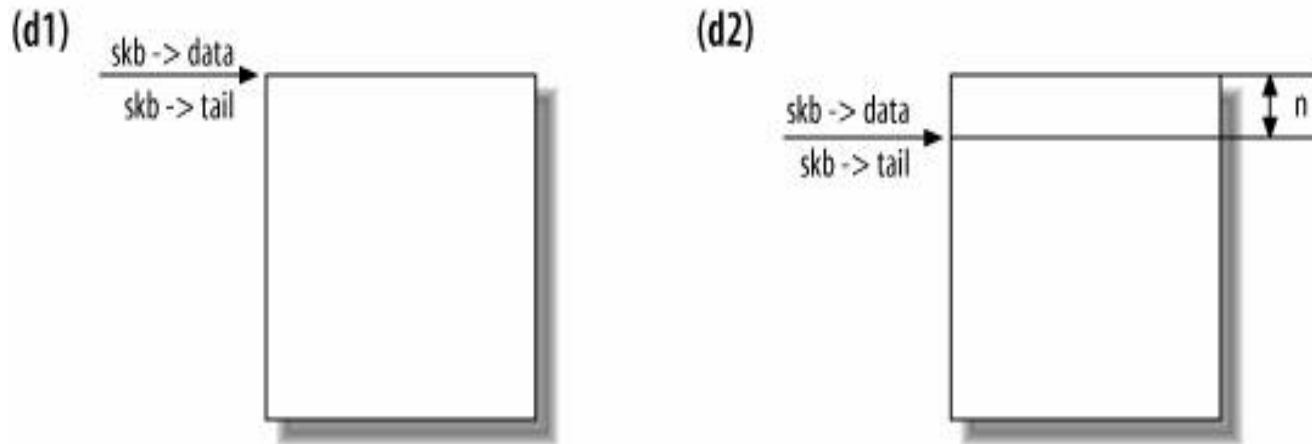
# management functions

`skb_pull(struct sk_buff *skb, unsigned int len)`



# management functions

`skb_reserve(struct sk_buff *skb, int len)`



Each of the above four memory management functions return the data ptr.

# memory allocation

defined in `<net/core/skbuff.c>`

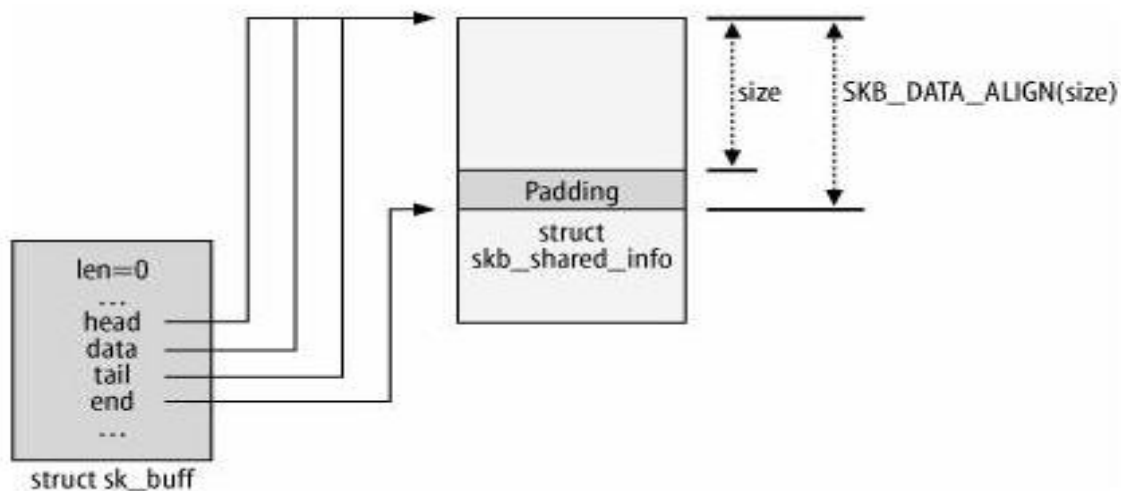
```
struct sk_buff * __alloc_skb(unsigned int size, gfp_t gfp_mask,  
                             int fclone, int node)
```

...

```
size = SKB_DATA_ALIGN(size);
```

```
data = kmalloc(size + sizeof(struct skb_shared_info), gfp_mask);
```

...





# memory allocation

```
struct sk_buff *__netdev_alloc_skb(struct net_device *dev,  
    unsigned int length, gfp_t gfp_mask)
```

The buffer allocation function meant for use by device drivers

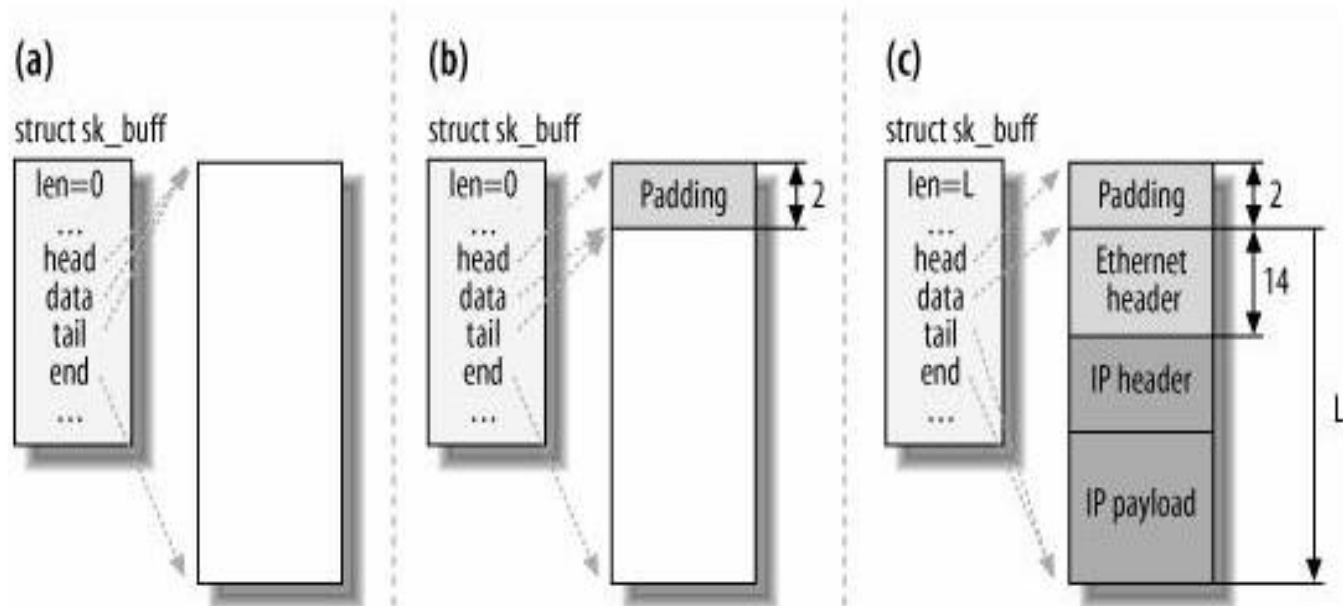
Executed in interrupt mode

**Freeing memory: kfree\_skb and dev\_kfree\_skb**

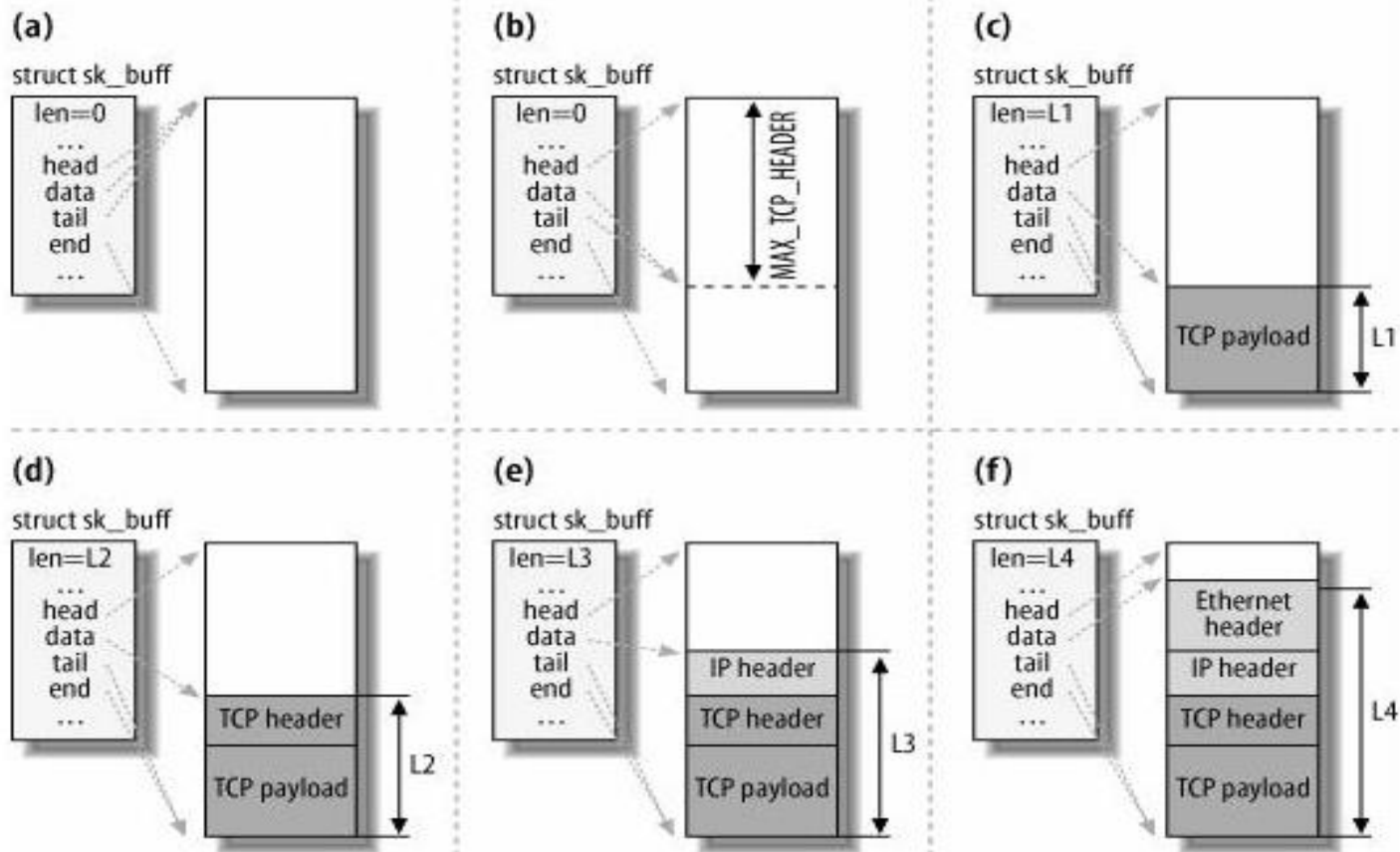
Release buffer back to the buffer-pool.

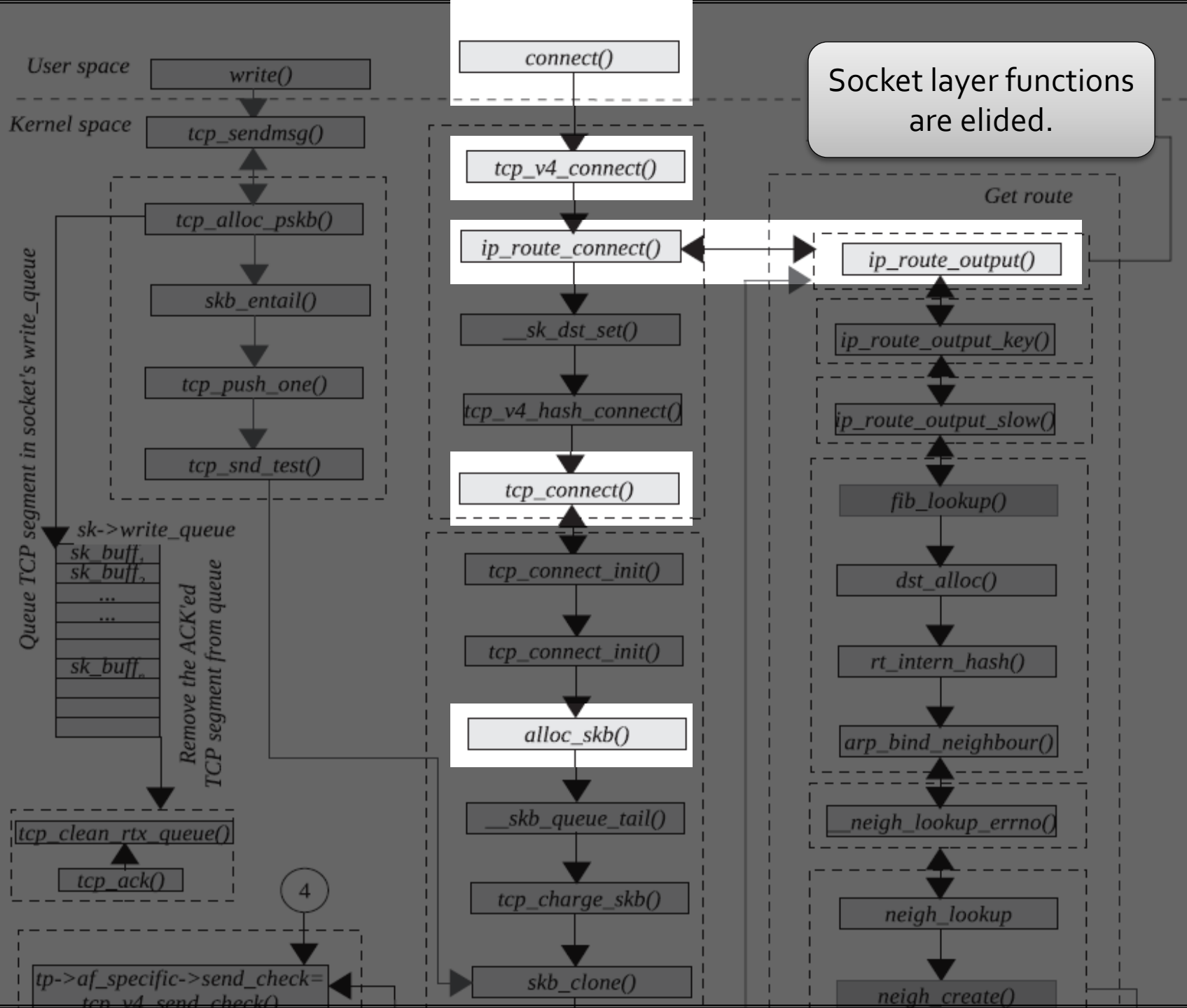
Buffer released only when skb\_users counter is 1. If not, the counter is decremented.

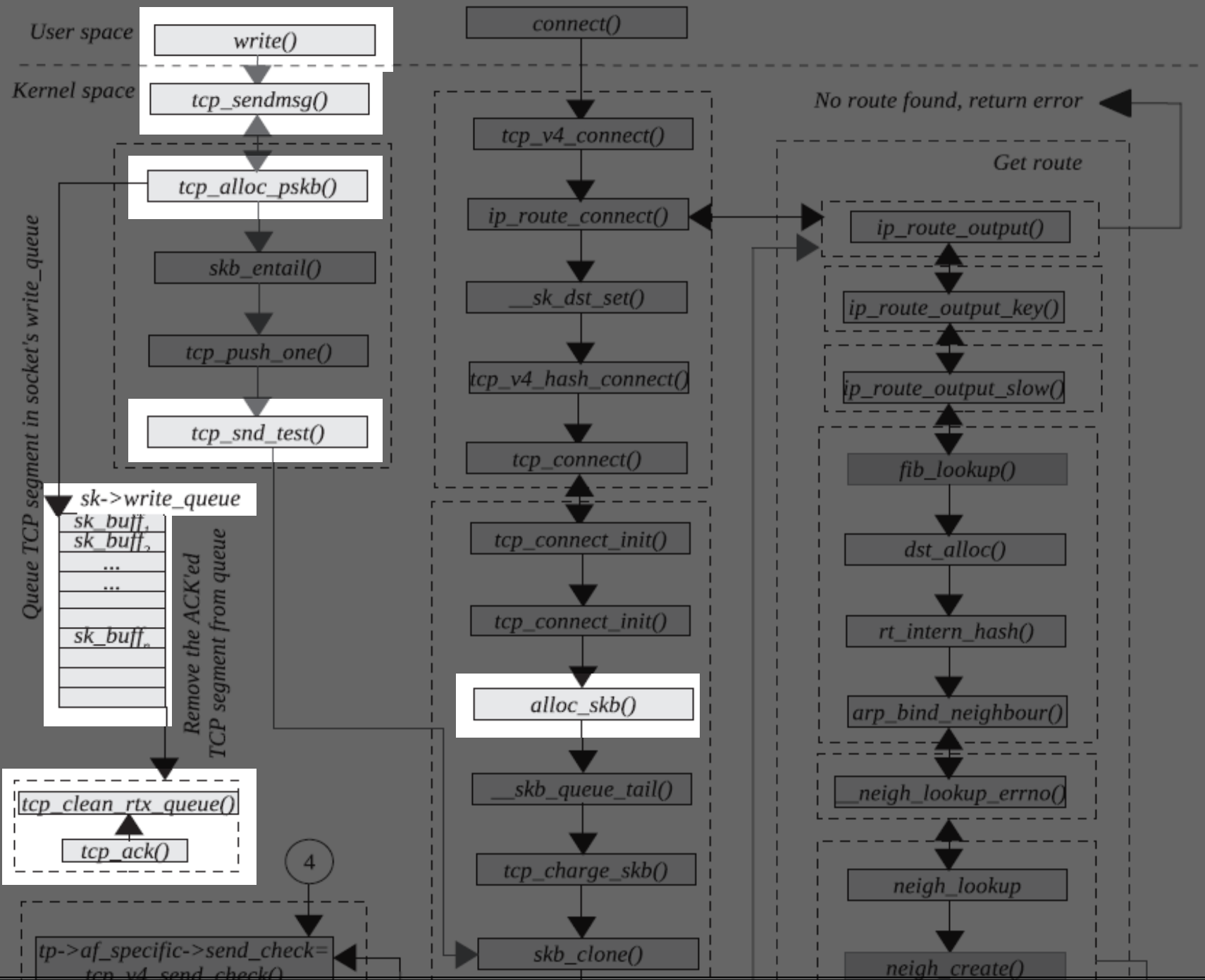
# initializing buffer - reception

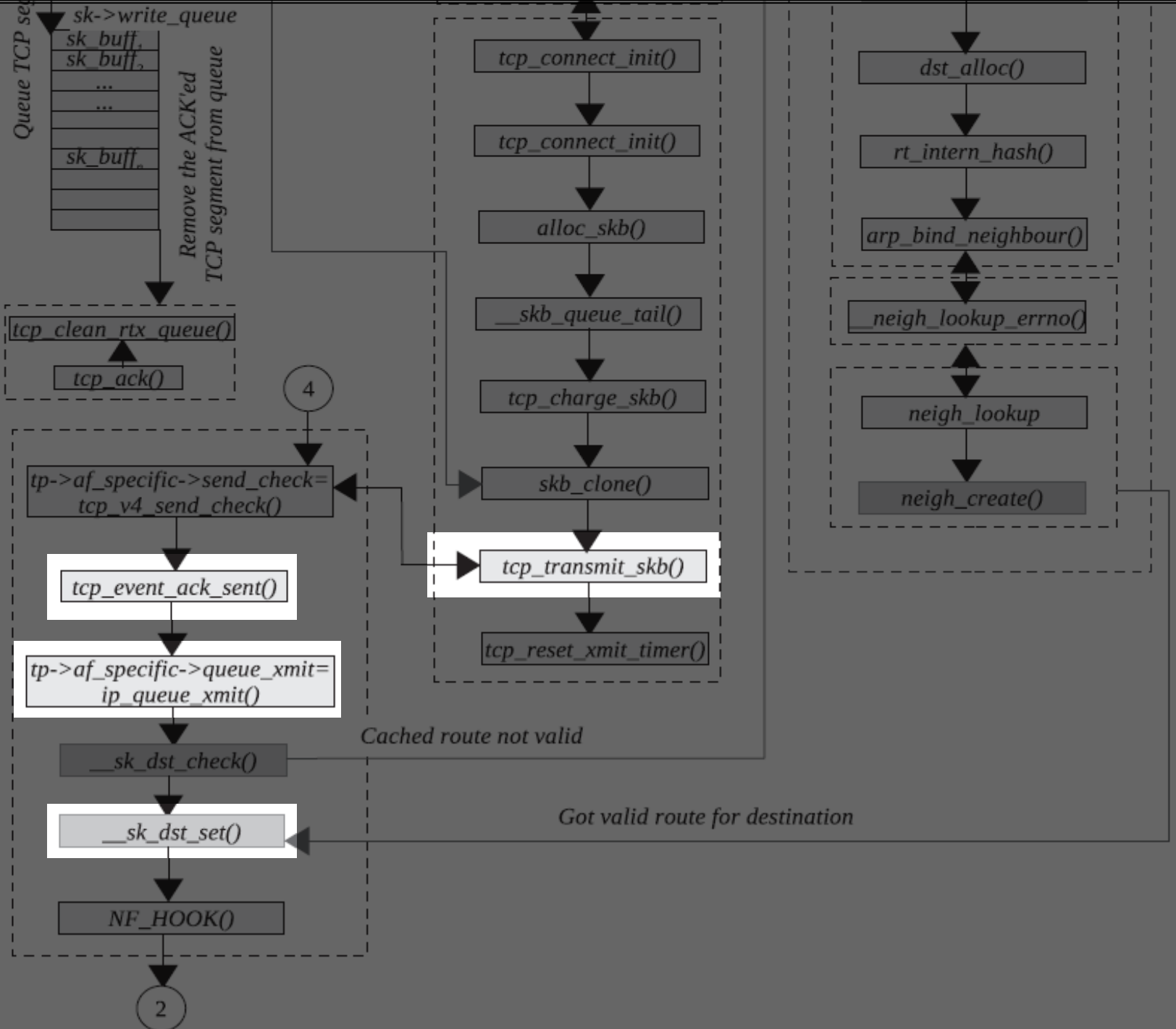


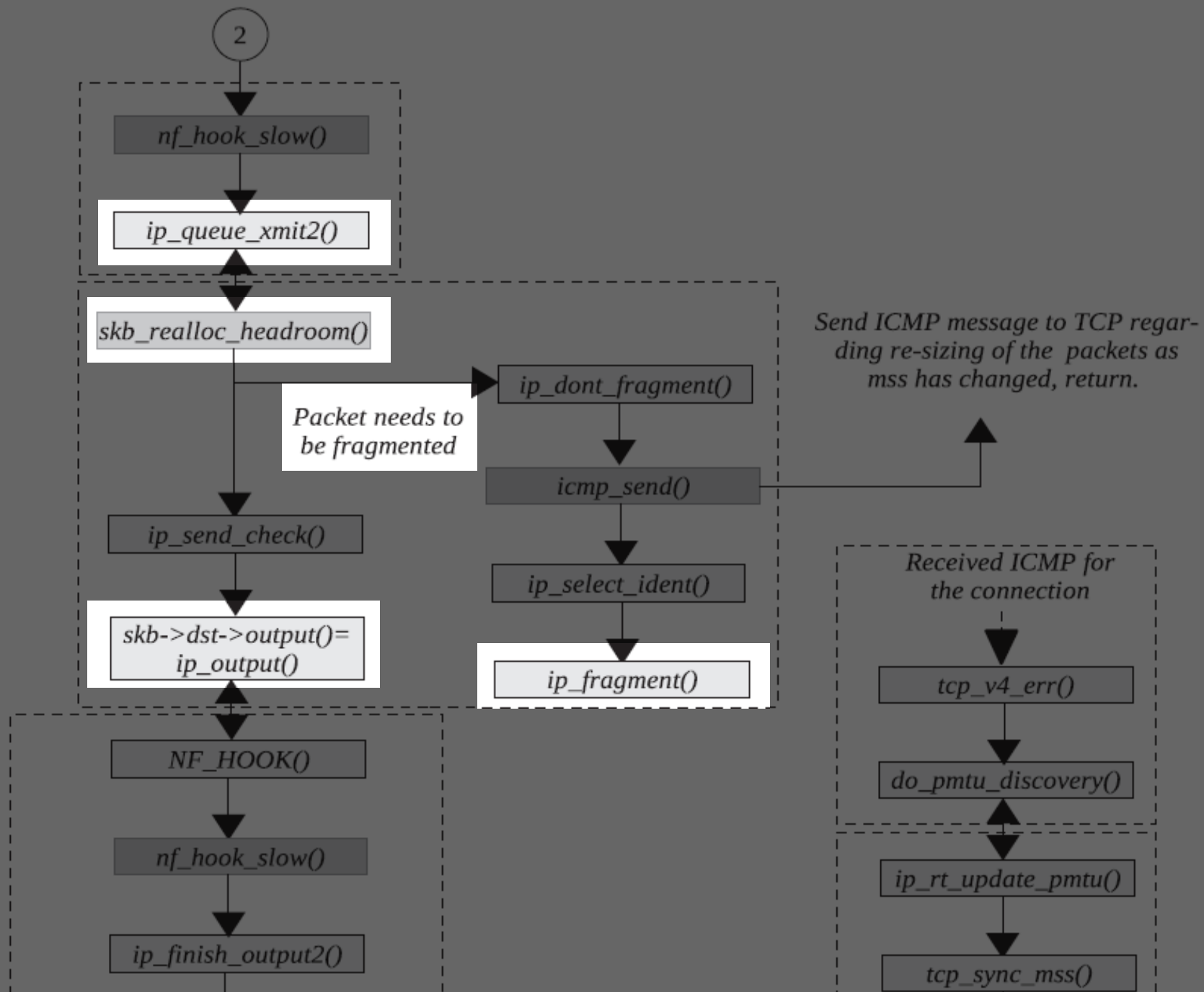
# initializing buffer - transmission

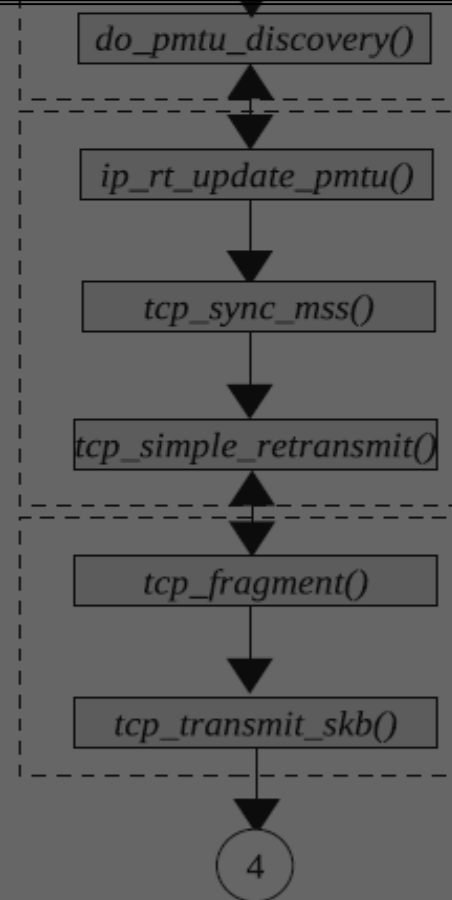
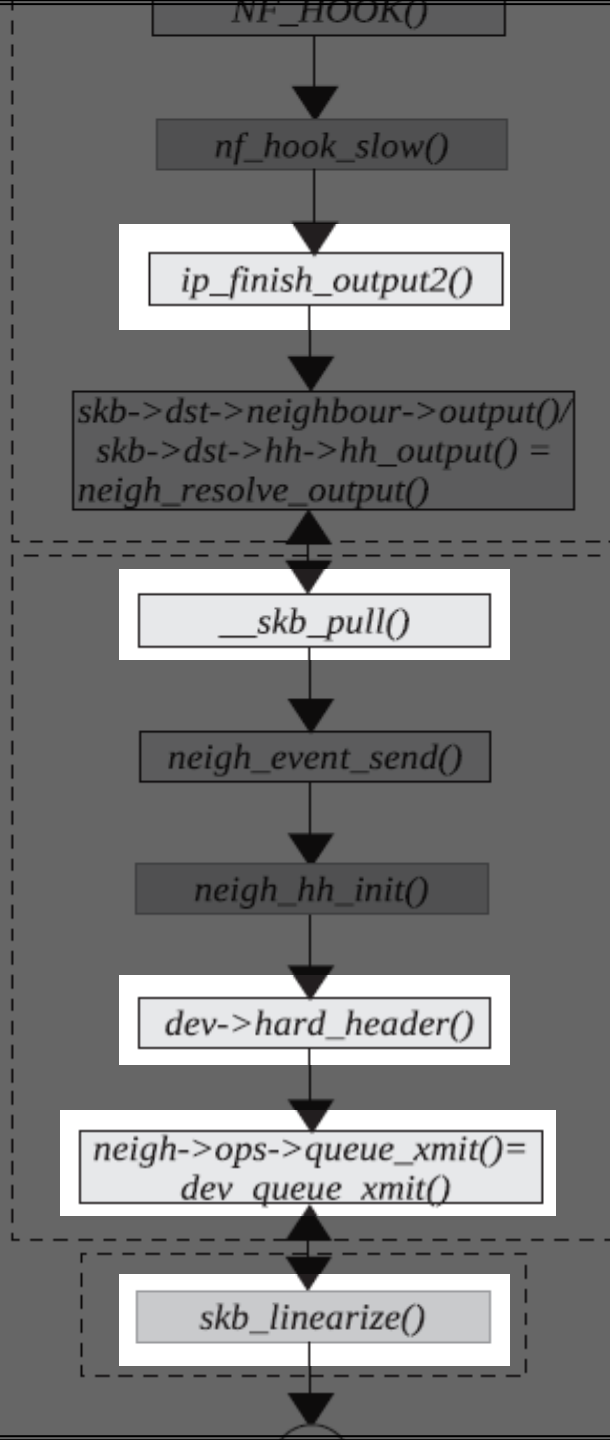




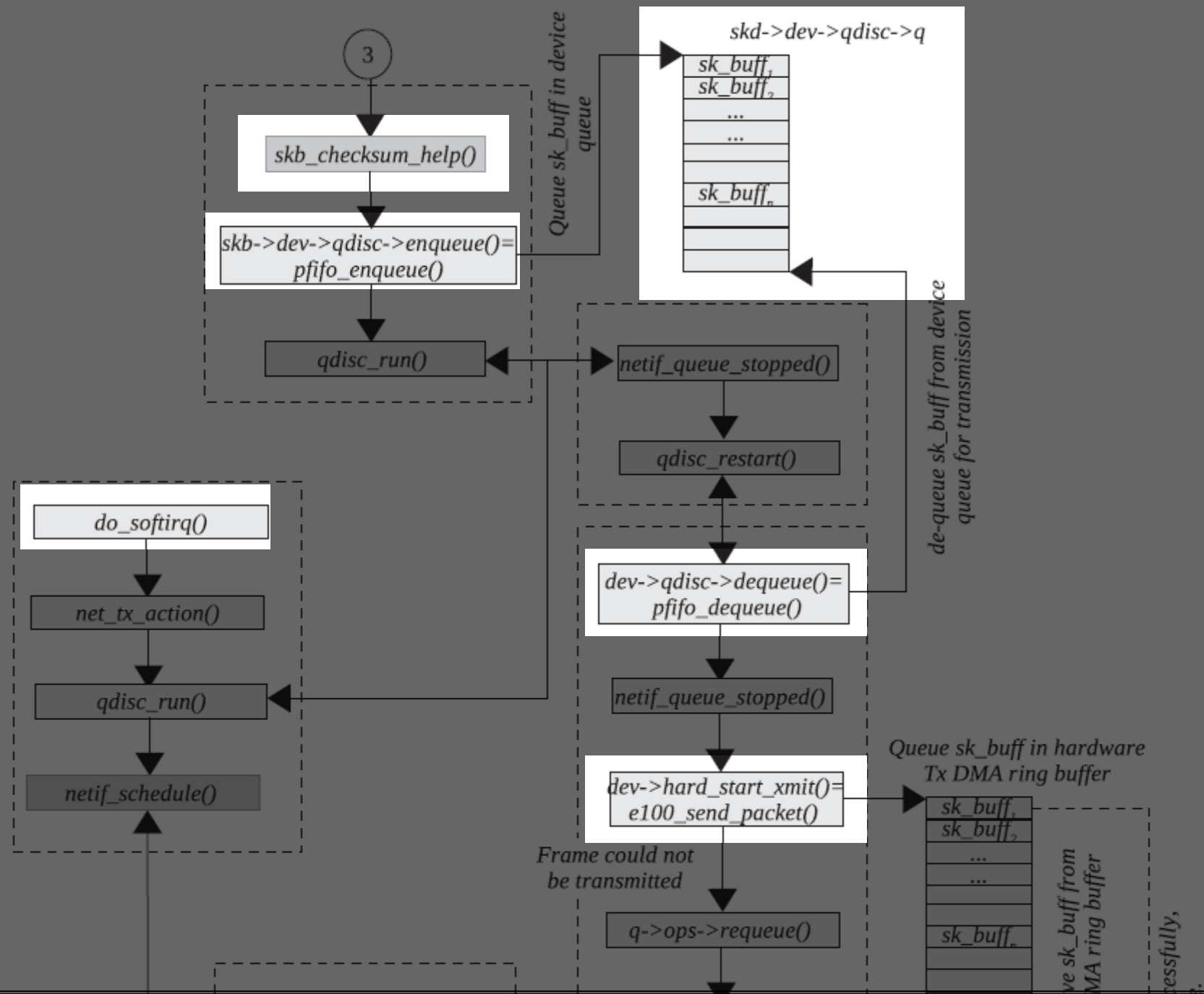


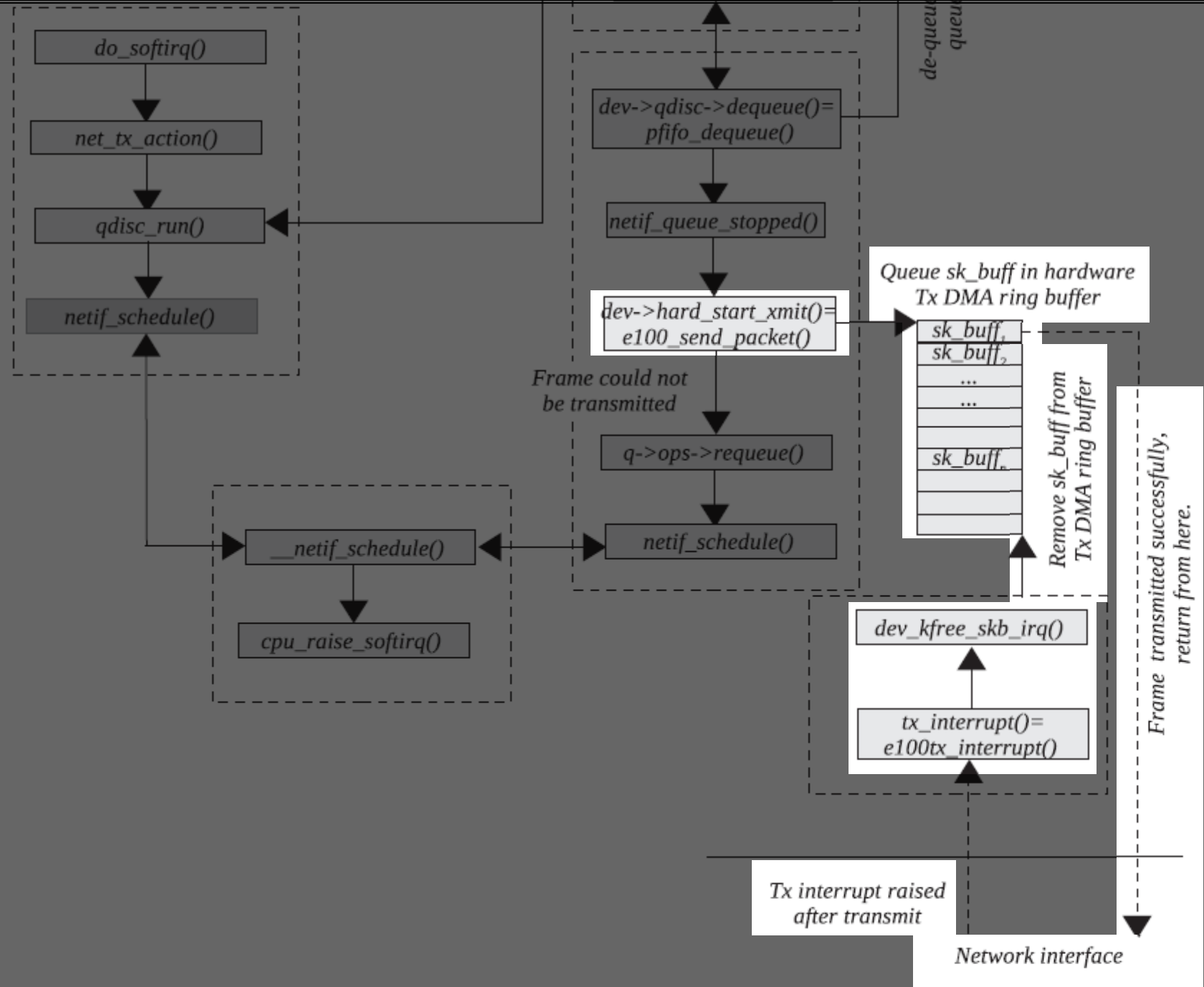




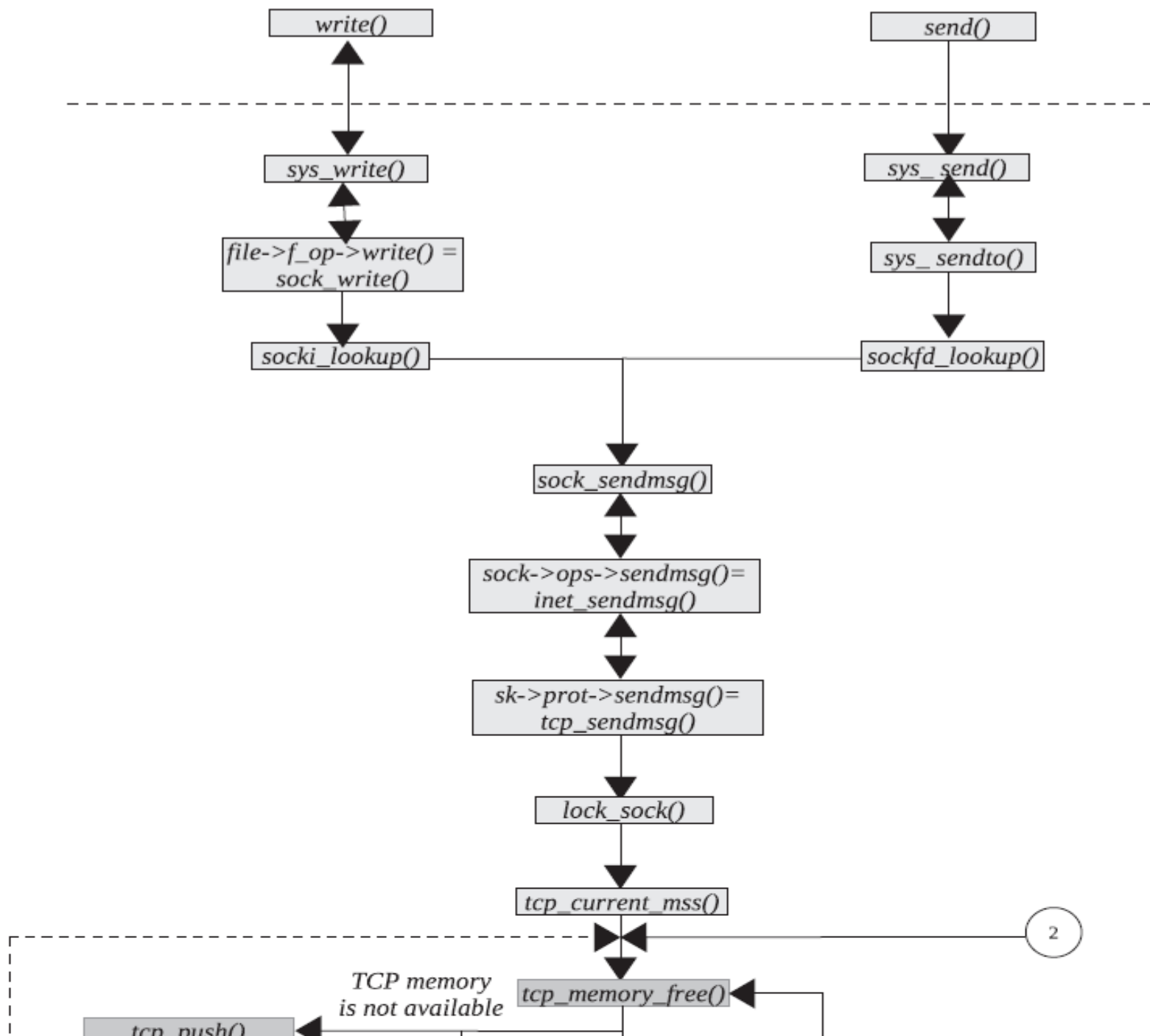


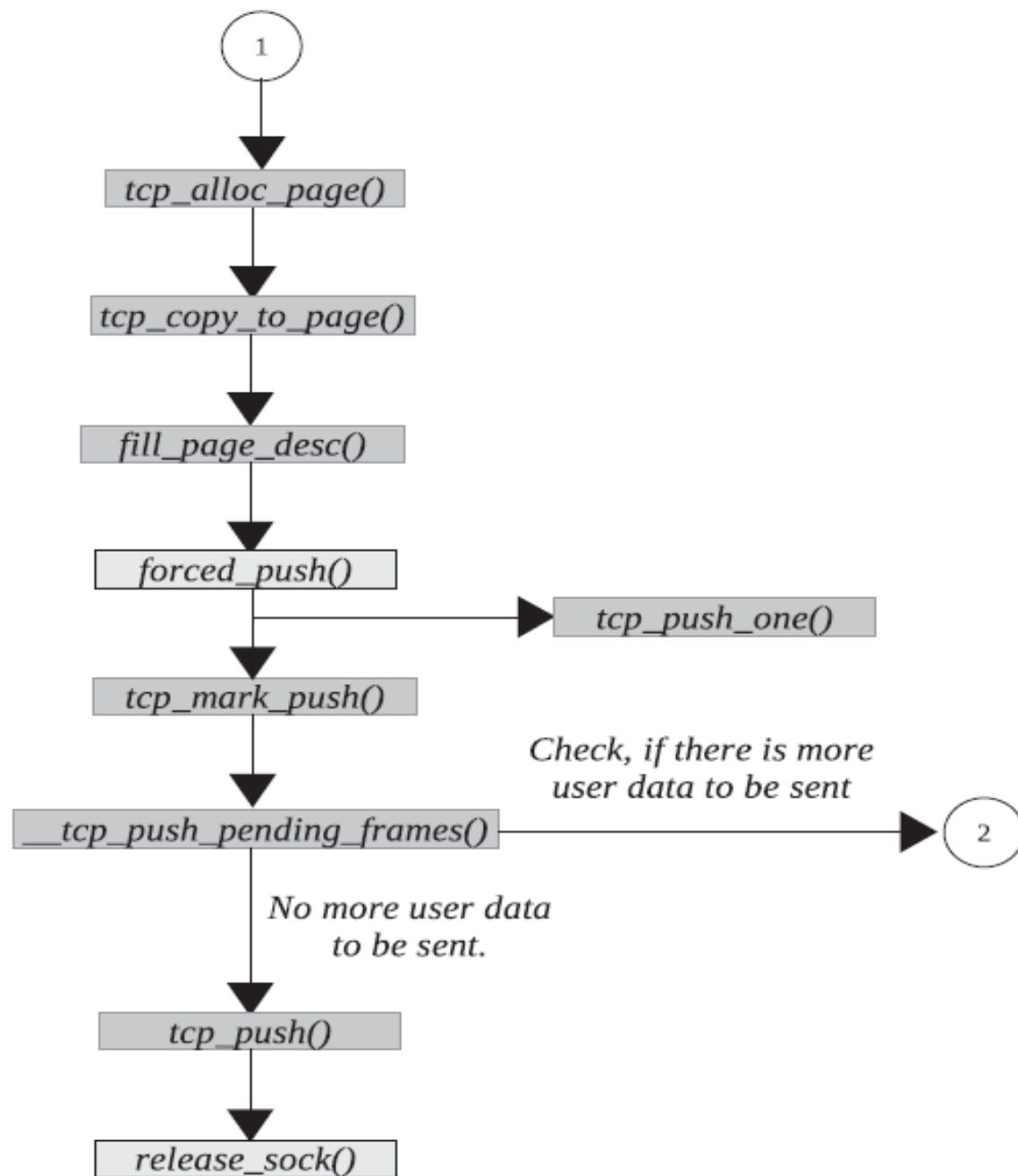






# Overall High Level Functional Overview of Sending





**Figure 7.6b.** Functional flow of TCP send process (*continued*).

# net\_device

- *Defined in <include/linux/netdevice.h>*
- stores all information specifically regarding a network device
- one such structure for each device, both real ones (such as Ethernet NICs) and virtual ones
- Network devices can be classified into types such as *Ethernet cards* and *Token Ring cards*
- Each type may come in several models.
- Model specific parameters are initialized by device driver software.
- Parameters common for different models are initiated by kernel.

# configuration parameters

```
struct net_device{  
    char                name[IFNAMSIZ];  
    int                 ifindex;  
  
    /* device name hash chain, ex: eth0 */  
    struct hlist_node    name_hlist;  
  
    unsigned long        mem_end; /* shared mem end    */  
    unsigned long        mem_start; /* shared mem start */  
    unsigned long        base_addr; /* device I/O address */  
    unsigned int         irq;        /* device IRQ number */  
    unsigned char        if_port;    /* Selectable AUI, TP,.. */  
    unsigned char        dma;        /* DMA channel    */  

```

...

# configuration parameters

```
struct net_device{
    char                name[IFNAMSIZ];
    int                 ifindex;

    /* device name hash chain, ex: eth0 */
    struct hlist_node    name_hlist;

    unsigned long        mem_end;        /* shared mem end */
    unsigned long        mem_start;      /* shared mem start */
    unsigned long        base_addr;      /* device I/O address */
    unsigned int         irq;              /* device IRQ number */
    unsigned char        if_port;          /* Selectable AUI, TP,.. */
    unsigned char        dma;              /* DMA channel */

```

...



# configuration parameters

```
struct net_device{
    char                name[IFNAMSIZ];
    int                 ifindex;

    /* device name hash chain, ex: eth0 */
    struct hlist_node    name_hlist;

    unsigned long        mem_end; /* shared mem end    */
    unsigned long        mem_start; /* shared mem start */
    unsigned long        base_addr; /* device I/O address */
    unsigned int          irq; /* device IRQ number */
    unsigned char        if_port; /* Selectable AUI, TP,.. */
    unsigned char        dma; /* DMA channel
    */
```

# configuration parameters

```
struct net_device{
    char                name[IFNAMSIZ];
    /* device name hash chain, ex: eth0 */
    struct hlist_node    name_hlist;

    unsigned long        mem_end; /* shared mem end    */
    unsigned long        mem_start; /* shared mem start */
    unsigned long        base_addr; /* device I/O address */
    unsigned int         irq; /* device IRQ number */
    unsigned char        if_port; /* Selectable AUI, TP,.. */
    unsigned char        dma; /* DMA channel */
    unsigned short       flags; /* interface flags (a la BSD) */

    ...
}
```

# configuration parameters

```
struct net_device{
    char                name[IFNAMSIZ];
    /* device name hash chain, ex: eth0 */
    struct hlist_node    name_hlist;

    unsigned long        mem_end; /* shared mem end    */
    unsigned long        mem_start; /* shared mem start */
    unsigned long        base_addr; /* device I/O address */
    unsigned int         irq;       /* device IRQ number */
    unsigned char        if_port;   /* Selectable AUI, TP,.. */
    unsigned char        dma;      /* DMA channel    */
    unsigned short       flags;     /* interface flags (a la BSD)*/
    ...
}
```

# configuration parameters

```
struct net_device{
    char                name[IFNAMSIZ];
    /* device name hash chain, ex: eth0 */
    struct hlist_node    name_hlist;

    unsigned long        mem_end; /* shared mem end    */
    unsigned long        mem_start; /* shared mem start */
    unsigned long        base_addr; /* device I/O address */
    unsigned int         irq;      /* device IRQ number */
    unsigned char        if_port; /* Selectable AUI, TP,.. */
    unsigned char        dma;     /* DMA channel        */
    unsigned short      flags; /* interface flags (a la BSD)*/
    /* ex : IFF_UP || IFF_RUNNING || IFF_MULTICAST */
}
```

# configuration parameters

```
struct net_device{
```

```
...
```

```
    unsigned                mtu;                /* interface MTU value */
```

```
    unsigned short          type;                /* interface hardware type */
```

```
    unsigned short          hard_header_len;      /* hardware hdr length */
```

```
    unsigned char           dev_addr[MAX_ADDR_LEN];
```

```
    unsigned char           addr_len;             /* hardware address length */
```

```
    unsigned char           broadcast[MAX_ADDR_LEN];
```

```
    unsigned int            promiscuity;
```

```
...
```

# configuration parameters

```
struct net_device{
```

```
...
```

```
    unsigned          mtu;                /* interface MTU value          */  
    unsigned short    type;              /* interface hardware type*/  
    unsigned short    hard_header_len;    /* hardware hdr length          */
```

```
    unsigned char     dev_addr[MAX_ADDR_LEN];  
    unsigned char     addr_len;           /* hardware address length      */
```

```
    unsigned char     broadcast[MAX_ADDR_LEN];  
    unsigned int       promiscuity;
```

```
...
```

# configuration parameters

```
struct net_device{
```

```
...
```

```
    unsigned          mtu;                /* interface MTU value      */
    unsigned short    type;                /* interface hardware type  */
    unsigned short    hard_header_len; /* hardware hdr length */
```

```
    unsigned char     dev_addr[MAX_ADDR_LEN];
    unsigned char     addr_len;           /* hardware address length  */
```

```
    unsigned char     broadcast[MAX_ADDR_LEN];
    unsigned int       promiscuity;
```

```
...
```

# configuration parameters

```
struct net_device{
```

```
...
```

```
    unsigned          mtu;                /* interface MTU value      */
    unsigned short     type;               /* interface hardware type  */
    unsigned short     hard_header_len;    /* hardware hdr length      */
```

```
    unsigned char      dev_addr[MAX_ADDR_LEN];
    unsigned char      addr_len;          /* hardware address length*/
```

```
    unsigned char      broadcast[MAX_ADDR_LEN];
    unsigned int        promiscuity;
```

```
...
```



# configuration parameters

```
struct net_device{
```

```
...
```

```
    unsigned          mtu;                /* interface MTU value      */
    unsigned short     type;               /* interface hardware type  */
    unsigned short     hard_header_len;    /* hardware hdr length      */
```

```
    unsigned char      dev_addr[MAX_ADDR_LEN];
    unsigned char      addr_len;          /* hardware address length  */
```

```
    unsigned char      broadcast[MAX_ADDR_LEN];
    unsigned int        promiscuity;
```

```
...
```

# configuration parameters

```
struct net_device{
```

```
...
```

```
    unsigned          mtu;                /* interface MTU value      */
    unsigned short     type;               /* interface hardware type  */
    unsigned short     hard_header_len;    /* hardware hdr length      */
```

```
    unsigned char      dev_addr[MAX_ADDR_LEN];
    unsigned char      addr_len;          /* hardware address length  */
```

```
    unsigned char      broadcast[MAX_ADDR_LEN];
    unsigned int         promiscuity;
```

```
...
```

# list management

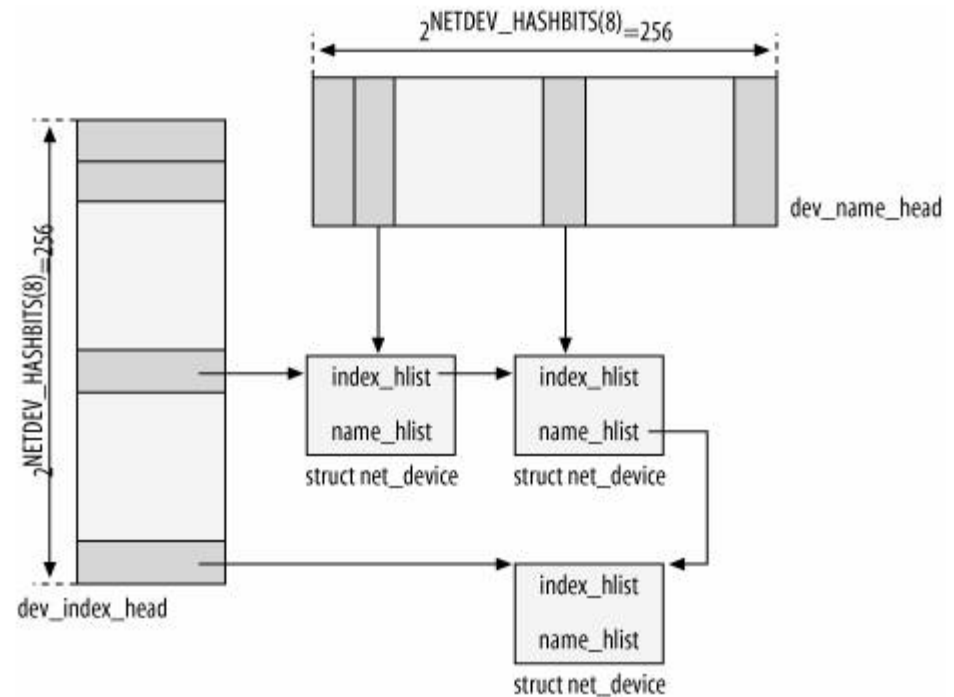
```
struct net_device{
```

```
...
```

```
struct net_device *next;
```

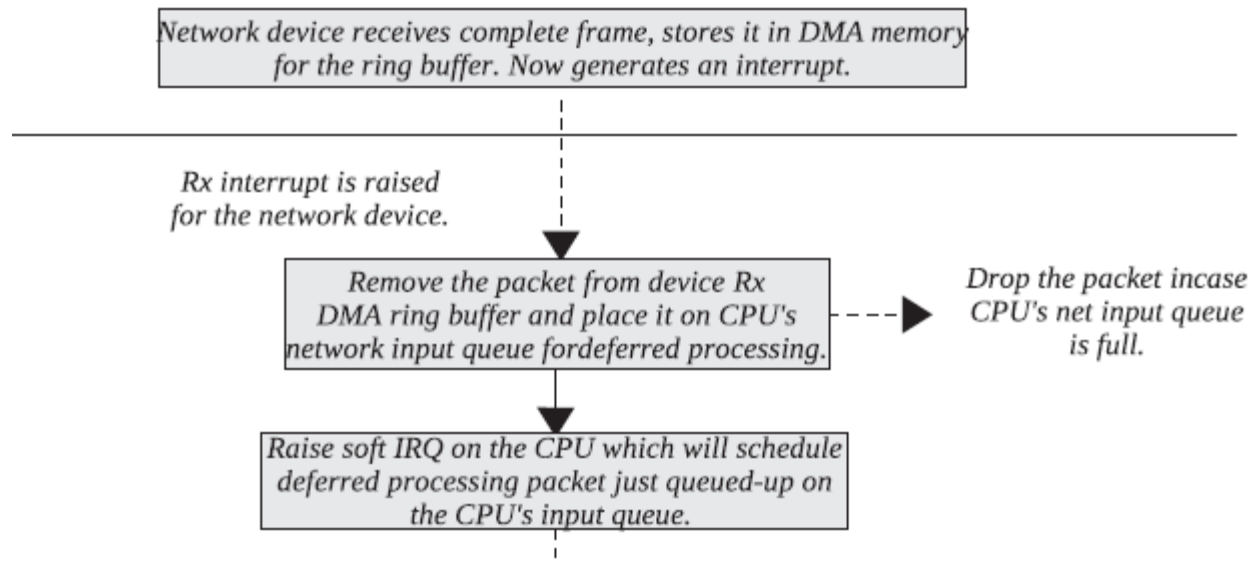
```
struct hlist_node name_hlist;
```

```
struct hlist_node index_hlist
```

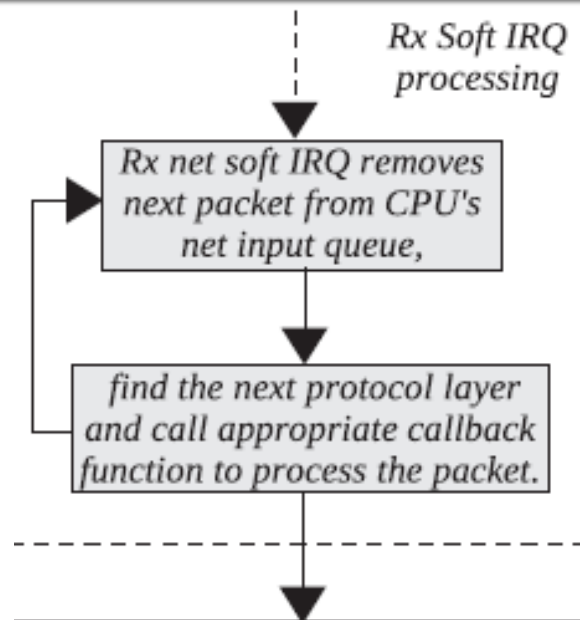


# Walkthrough Reception

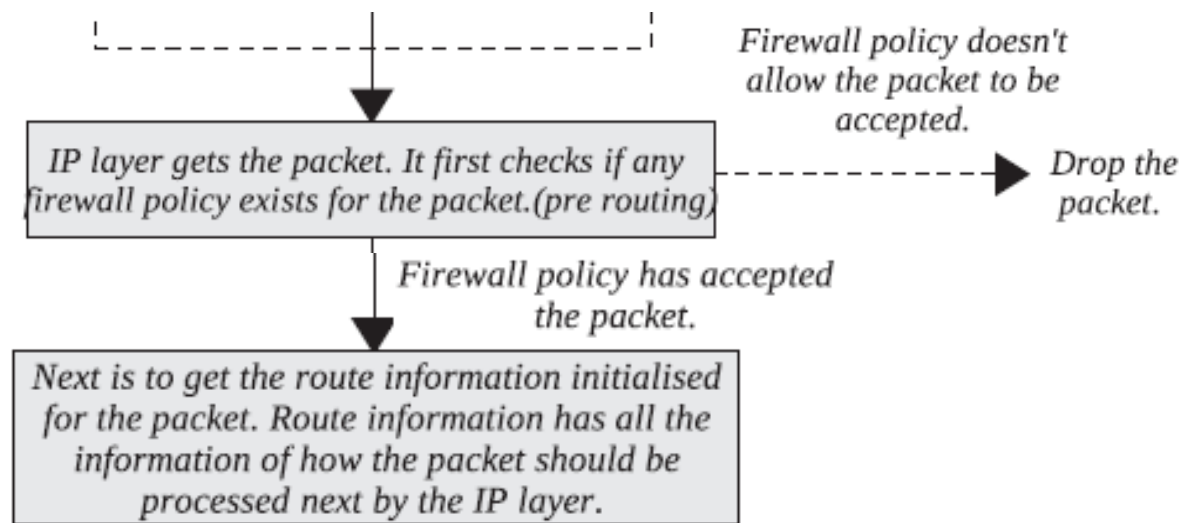
# Walkthrough Reception



- We don't process the packet in the interrupt subroutine.
- `Netif_rx()` – raise the net Rx softIRQ.
- `Net_rx_action()` is called - start processing the packet
- Processing of packet starts with the protocol switching section



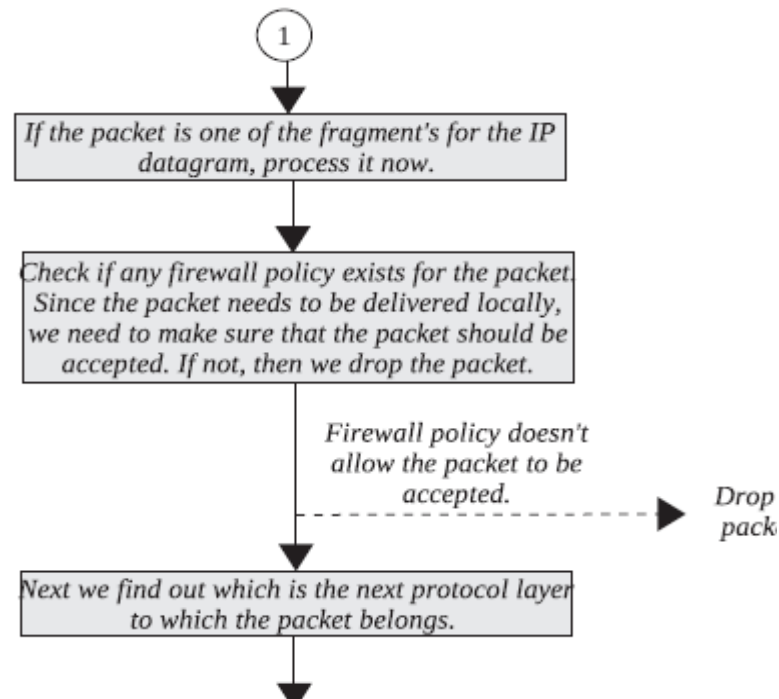
- *Netif\_receive\_skb()* is called to process the packet and find out the next protocol layer.
- Protocol family of the packet is extracted from the link layer header.



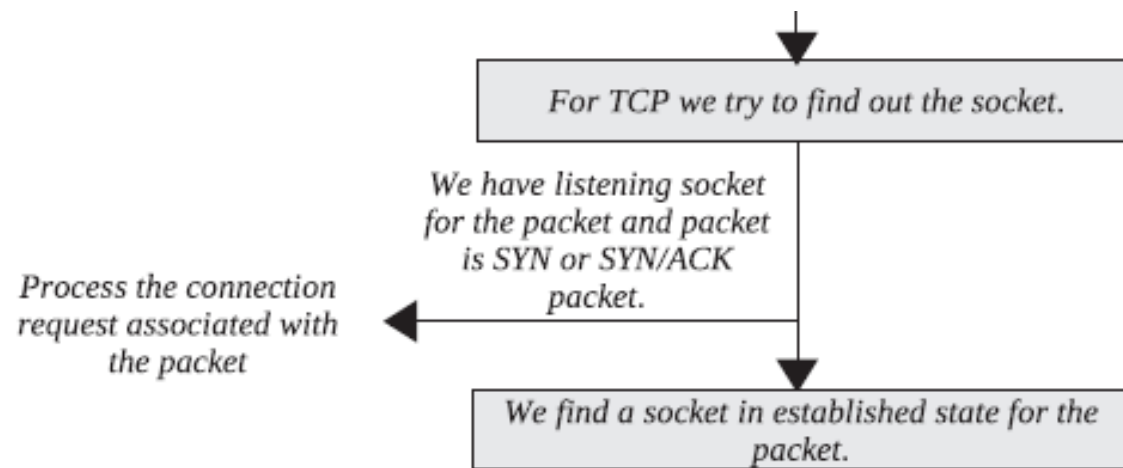
- *ip\_rcv()* is an entry point for IP packets processing.
- Checks if the packet we have is destined for some other host (using `PACKET_OTHERHOST`)
- Check the checksum of the packet by calling `ip_fast_csum()`

- Call `ip_route_input()` , this routine checks kernel routing table `rt_hash_table`.
- If packet needs to be forwarded input routine is `ip_forward()`
- Otherwise `ip_local_deliver()`
- `ip_send()` is called to check if the packet needs to be fragmented
- If yes , fragment the packet by calling `ip_fragment()`
- Packet output path – `ip_finish_output()`
- `ip_local_deliver()` – packets need to delivered locally

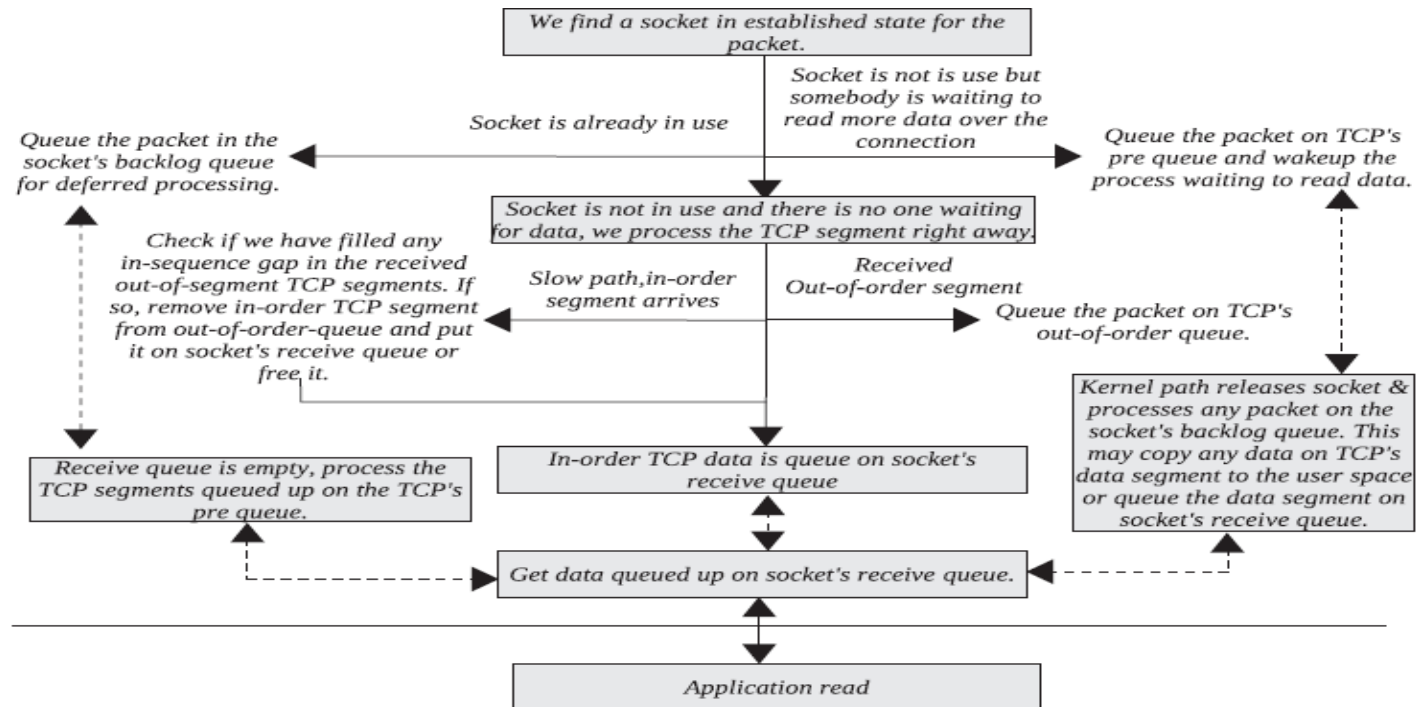




- *ip\_defrag()*
- *Protocol identifier field `skb->np.iph->protocol` (in IP header).*
- *For TCP, we find the receive handler as `tcp_v4_rcv()` (entry point for the TCP layer)*



- `_tcp_v4_lookup()` – find the socket to which the packet belongs
- Established sockets are maintained in the hash table `tcp_ehash`.
- Established socket not found – New connection request for any listening socket
- Search for listening socket – `tcp_v4_lookup_listener()`
- `tcp_rcv_established()`



- Application read the data from the receive queue if it issues `recv()`
- Kernel routine to read data from TCP socket is `tcp_recvmsg()`

**Thank You**  
**Any Questions?**