

---

# **Linux Networking Documentation**

*Release*

**The kernel development community**

**Nov 02, 2017**



<b>1</b>	<b>batman-adv</b>	<b>3</b>
1.1	Configuration . . . . .	3
1.2	Usage . . . . .	4
1.3	Logging/Debugging . . . . .	5
1.4	batctl . . . . .	5
1.5	Contact . . . . .	6
<b>2</b>	<b>Linux Networking and Network Devices APIs</b>	<b>7</b>
2.1	Linux Networking . . . . .	7
2.2	Network device support . . . . .	77
<b>3</b>	<b>Z8530 Programming Guide</b>	<b>125</b>
3.1	Introduction . . . . .	125
3.2	Driver Modes . . . . .	125
3.3	Using the Z85230 driver . . . . .	125
3.4	Attaching Network Interfaces . . . . .	126
3.5	Configuring And Activating The Port . . . . .	126
3.6	Network Layer Functions . . . . .	127
3.7	Porting The Z8530 Driver . . . . .	127
3.8	Known Bugs And Assumptions . . . . .	128
3.9	Public Functions Provided . . . . .	128
3.10	Internal Functions . . . . .	131
	<b>Index</b>	<b>137</b>



Contents:



## **BATMAN-ADV**

Batman advanced is a new approach to wireless networking which does no longer operate on the IP basis. Unlike the batman daemon, which exchanges information using UDP packets and sets routing tables, batman-advanced operates on ISO/OSI Layer 2 only and uses and routes (or better: bridges) Ethernet Frames. It emulates a virtual network switch of all nodes participating. Therefore all nodes appear to be link local, thus all higher operating protocols won't be affected by any changes within the network. You can run almost any protocol above batman advanced, prominent examples are: IPv4, IPv6, DHCP, IPX.

Batman advanced was implemented as a Linux kernel driver to reduce the overhead to a minimum. It does not depend on any (other) network driver, and can be used on wifi as well as ethernet lan, vpn, etc ... (anything with ethernet-style layer 2).

### **Configuration**

Load the batman-adv module into your kernel:

```
$ insmod batman-adv.ko
```

The module is now waiting for activation. You must add some interfaces on which batman can operate. After loading the module batman advanced will scan your systems interfaces to search for compatible interfaces. Once found, it will create subfolders in the /sys directories of each supported interface, e.g.:

```
$ ls /sys/class/net/eth0/batman_adv/  
elp_interval iface_status mesh_iface throughput_override
```

If an interface does not have the batman\_adv subfolder, it probably is not supported. Not supported interfaces are: loopback, non-ethernet and batman's own interfaces.

Note: After the module was loaded it will continuously watch for new interfaces to verify the compatibility. There is no need to reload the module if you plug your USB wifi adapter into your machine after batman advanced was initially loaded.

The batman-adv soft-interface can be created using the iproute2 tool ip:

```
$ ip link add name bat0 type batadv
```

To activate a given interface simply attach it to the bat0 interface:

```
$ ip link set dev eth0 master bat0
```

Repeat this step for all interfaces you wish to add. Now batman starts using/broadcasting on this/these interface(s).

By reading the "iface\_status" file you can check its status:

```
$ cat /sys/class/net/eth0/batman_adv/iface_status  
active
```

To deactivate an interface you have to detach it from the “bat0” interface:

```
$ ip link set dev eth0 nomaster
```

All mesh wide settings can be found in batman’s own interface folder:

```
$ ls /sys/class/net/bat0/mesh/
aggregated_ogms      fragmentation isolation_mark routing_algo
ap_isolation          gw_bandwidth  log_level      vlan0
bonding               gw_mode       multicast_mode
bridge_loop_avoidance gw_sel_class  network_coding
distributed_arp_table hop_penalty   orig_interval
```

There is a special folder for debugging information:

```
$ ls /sys/kernel/debug/batman_adv/bat0/
bla_backbone_table log          neighbors      transtable_local
bla_claim_table    mcast_flags originators
dat_cache          nc          socket
gateways           nc_nodes   transtable_global
```

Some of the files contain all sort of status information regarding the mesh network. For example, you can view the table of originators (mesh participants) with:

```
$ cat /sys/kernel/debug/batman_adv/bat0/originators
```

Other files allow to change batman’s behaviour to better fit your requirements. For instance, you can check the current originator interval (value in milliseconds which determines how often batman sends its broadcast packets):

```
$ cat /sys/class/net/bat0/mesh/orig_interval
1000
```

and also change its value:

```
$ echo 3000 > /sys/class/net/bat0/mesh/orig_interval
```

In very mobile scenarios, you might want to adjust the originator interval to a lower value. This will make the mesh more responsive to topology changes, but will also increase the overhead.

## Usage

To make use of your newly created mesh, batman advanced provides a new interface “bat0” which you should use from this point on. All interfaces added to batman advanced are not relevant any longer because batman handles them for you. Basically, one “hands over” the data by using the batman interface and batman will make sure it reaches its destination.

The “bat0” interface can be used like any other regular interface. It needs an IP address which can be either statically configured or dynamically (by using DHCP or similar services):

```
NodeA: ip link set up dev bat0
NodeA: ip addr add 192.168.0.1/24 dev bat0

NodeB: ip link set up dev bat0
NodeB: ip addr add 192.168.0.2/24 dev bat0
NodeB: ping 192.168.0.1
```

Note: In order to avoid problems remove all IP addresses previously assigned to interfaces now used by batman advanced, e.g.:



```
$ ip addr flush dev eth0
```

## Logging/Debugging

All error messages, warnings and information messages are sent to the kernel log. Depending on your operating system distribution this can be read in one of a number of ways. Try using the commands: `dmesg`, `logread`, or looking in the files `/var/log/kern.log` or `/var/log/syslog`. All batman-adv messages are prefixed with “batman-adv:” So to see just these messages try:

```
$ dmesg | grep batman-adv
```

When investigating problems with your mesh network, it is sometimes necessary to see more detail debug messages. This must be enabled when compiling the batman-adv module. When building batman-adv as part of kernel, use “make menuconfig” and enable the option B.A.T.M.A.N. debugging (`CONFIG_BATMAN_ADV_DEBUG=y`).

Those additional debug messages can be accessed using a special file in debugfs:

```
$ cat /sys/kernel/debug/batman_adv/bat0/log
```

The additional debug output is by default disabled. It can be enabled during run time. Following log\_levels are defined:

0	All debug output disabled
1	Enable messages related to routing / flooding / broadcasting
2	Enable messages related to route added / changed / deleted
4	Enable messages related to translation table operations
8	Enable messages related to bridge loop avoidance
16	Enable messages related to DAT, ARP snooping and parsing
32	Enable messages related to network coding
64	Enable messages related to multicast
128	Enable messages related to throughput meter
255	Enable all messages

The debug output can be changed at runtime using the file `/sys/class/net/bat0/mesh/log_level`. e.g.:

```
$ echo 6 > /sys/class/net/bat0/mesh/log_level
```

will enable debug messages for when routes change.

Counters for different types of packets entering and leaving the batman-adv module are available through `ethtool`:

```
$ ethtool --statistics bat0
```

## batctl

As batman advanced operates on layer 2, all hosts participating in the virtual switch are completely transparent for all protocols above layer 2. Therefore the common diagnosis tools do not work as expected. To overcome these problems, `batctl` was created. At the moment the `batctl` contains ping, traceroute, tcpdump and interfaces to the kernel module settings.

For more information, please see the manpage (`man batctl`).

`batctl` is available on <https://www.open-mesh.org/>

## Contact

Please send us comments, experiences, questions, anything :)

**IRC:** #batman on irc.freenode.org

**Mailing-list:** [b.a.t.m.a.n@open-mesh.org](mailto:b.a.t.m.a.n@open-mesh.org) (optional subscription at <https://lists.open-mesh.org/mm/listinfo/b.a.t.m.a.n>)

You can also contact the Authors:

- Marek Lindner <[mareklindner@neomailbox.ch](mailto:mareklindner@neomailbox.ch)>
- Simon Wunderlich <[sw@simonwunderlich.de](mailto:sw@simonwunderlich.de)>

## LINUX NETWORKING AND NETWORK DEVICES APIS

### Linux Networking

#### Networking Base Types

enum **sock\_type**  
Socket types

##### Constants

**SOCK\_STREAM** stream (connection) socket

**SOCK\_DGRAM** datagram (conn.less) socket

**SOCK\_RAW** raw socket

**SOCK\_RDM** reliably-delivered message

**SOCK\_SEQPACKET** sequential packet socket

**SOCK\_DCCP** Datagram Congestion Control Protocol socket

**SOCK\_PACKET** linux specific way of getting packets at the dev level. For writing rarp and other similar things on the user level.

##### Description

When adding some new socket type please grep ARCH\_HAS\_SOCKET\_TYPE include/asm-\* /socket.h, at least MIPS overrides this enum for binary compat reasons.

struct **socket**  
general BSD socket

##### Definition

```
struct socket {  
    socket_state state;  
    short type;  
    unsigned long flags;  
    struct socket_wq __rcu * wq;  
    struct file * file;  
    struct sock * sk;  
    const struct proto_ops * ops;  
};
```

##### Members

**state** socket state (SS\_CONNECTED, etc)

**type** socket type (SOCK\_STREAM, etc)

**flags** socket flags (SOCK\_NOSPACE, etc)

**wq** wait queue for several uses

**file** File back pointer for gc  
**sk** internal networking protocol agnostic socket representation  
**ops** protocol specific socket operations

## Socket Buffer Functions

**skb\_frag\_foreach\_page**(*f, f\_off, f\_len, p, p\_off, p\_len, copied*)  
loop over pages in a fragment

### Parameters

**f** skb frag to operate on  
**f\_off** offset from start of f->page.p  
**f\_len** length from f\_off to loop over  
**p** (temp var) current page  
**p\_off** (temp var) offset from start of current page, non-zero only on first page.  
**p\_len** (temp var) length in current page, < PAGE\_SIZE only on first and last page.  
**copied** (temp var) length so far, excluding current p\_len.

### Description

A fragment can hold a compound page, in which case per-page operations, notably kmap\_atomic, must be called for each regular page.

struct **skb\_shared\_hwtstamps**  
hardware time stamps

### Definition

```
struct skb_shared_hwtstamps {
    ktime_t hwtstamp;
};
```

### Members

**hwtstamp** hardware time stamp transformed into duration since arbitrary point in time

### Description

Software time stamps generated by ktime\_get\_real() are stored in skb->tstamp.

hwtstamps can only be compared against other hwtstamps from the same device.

This structure is attached to packets as part of the skb\_shared\_info. Use skb\_hwtstamps() to get a pointer.

struct **sk\_buff**  
socket buffer

### Definition

```
struct sk_buff {
    union {unnamed_union};
    __u16 inner_transport_header;
    __u16 inner_network_header;
    __u16 inner_mac_header;
    __be16 protocol;
    __u16 transport_header;
    __u16 network_header;
    __u16 mac_header;
    sk_buff_data_t tail;
};
```

```

sk_buff_data_t end;
unsigned char * head;
unsigned char * data;
unsigned int truesize;
refcount_t users;
};

```

## Members

**{unnamed\_union}** anonymous

**inner\_transport\_header** Inner transport layer header (encapsulation)

**inner\_network\_header** Network layer header (encapsulation)

**inner\_mac\_header** Link layer header (encapsulation)

**protocol** Packet protocol from driver

**transport\_header** Transport layer header

**network\_header** Network layer header

**mac\_header** Link layer header

**tail** Tail pointer

**end** End pointer

**head** Head of buffer

**data** Data head pointer

**truesize** Buffer size

**users** User count - see {datagram,tcp}.c

struct dst\_entry \* **skb\_dst**(const struct *sk\_buff* \* *skb*)  
returns skb dst\_entry

## Parameters

**const struct sk\_buff \* skb** buffer

## Description

Returns skb dst\_entry, regardless of reference taken or not.

void **skb\_dst\_set**(struct *sk\_buff* \* *skb*, struct dst\_entry \* *dst*)  
sets skb dst

## Parameters

**struct sk\_buff \* skb** buffer

**struct dst\_entry \* dst** dst entry

## Description

Sets skb dst, assuming a reference was taken on dst and should be released by *skb\_dst\_drop*()

void **skb\_dst\_set\_noref**(struct *sk\_buff* \* *skb*, struct dst\_entry \* *dst*)  
sets skb dst, hopefully, without taking reference

## Parameters

**struct sk\_buff \* skb** buffer

**struct dst\_entry \* dst** dst entry

### Description

Sets `skb` `dst`, assuming a reference was not taken on `dst`. If `dst` entry is cached, we do not take reference and `dst_release` will be avoided by `refdst_drop`. If `dst` entry is not cached, we take reference, so that last `dst_release` can destroy the `dst` immediately.

bool **skb\_dst\_is\_noref**(const struct *sk\_buff* \* `skb`)  
Test if `skb` `dst` isn't refcounted

### Parameters

const struct *sk\_buff* \* `skb` buffer

bool **skb\_fclone\_busy**(const struct *sock* \* `sk`, const struct *sk\_buff* \* `skb`)  
check if `fclone` is busy

### Parameters

const struct *sock* \* `sk` socket

const struct *sk\_buff* \* `skb` buffer

### Description

Returns true if `skb` is a fast clone, and its clone is not freed. Some drivers call *skb\_orphan()* in their `ndo_start_xmit()`, so we also check that this didnt happen.

int **skb\_pad**(struct *sk\_buff* \* `skb`, int `pad`)  
zero pad the tail of an `skb`

### Parameters

struct *sk\_buff* \* `skb` buffer to pad

int `pad` space to pad

### Description

Ensure that a buffer is followed by a padding area that is zero filled. Used by network drivers which may DMA or transfer data beyond the buffer end onto the wire.

May return error in out of memory cases. The `skb` is freed on error.

int **skb\_queue\_empty**(const struct *sk\_buff\_head* \* `list`)  
check if a queue is empty

### Parameters

const struct *sk\_buff\_head* \* `list` queue head

### Description

Returns true if the queue is empty, false otherwise.

bool **skb\_queue\_is\_last**(const struct *sk\_buff\_head* \* `list`, const struct *sk\_buff* \* `skb`)  
check if `skb` is the last entry in the queue

### Parameters

const struct *sk\_buff\_head* \* `list` queue head

const struct *sk\_buff* \* `skb` buffer

### Description

Returns true if `skb` is the last buffer on the list.

bool **skb\_queue\_is\_first**(const struct *sk\_buff\_head* \* `list`, const struct *sk\_buff* \* `skb`)  
check if `skb` is the first entry in the queue

### Parameters

const struct *sk\_buff\_head* \* `list` queue head

**const struct sk\_buff \* skb** buffer

### Description

Returns true if **skb** is the first buffer on the list.

struct [sk\\_buff](#) \* **skb\_queue\_next**(const struct sk\_buff\_head \* *list*, const struct [sk\\_buff](#) \* *skb*)  
return the next packet in the queue

### Parameters

**const struct sk\_buff\_head \* list** queue head

**const struct sk\_buff \* skb** current buffer

### Description

Return the next packet in **list** after **skb**. It is only valid to call this if [skb\\_queue\\_is\\_last\(\)](#) evaluates to false.

struct [sk\\_buff](#) \* **skb\_queue\_prev**(const struct sk\_buff\_head \* *list*, const struct [sk\\_buff](#) \* *skb*)  
return the prev packet in the queue

### Parameters

**const struct sk\_buff\_head \* list** queue head

**const struct sk\_buff \* skb** current buffer

### Description

Return the prev packet in **list** before **skb**. It is only valid to call this if [skb\\_queue\\_is\\_first\(\)](#) evaluates to false.

struct [sk\\_buff](#) \* **skb\_get**(struct [sk\\_buff](#) \* *skb*)  
reference buffer

### Parameters

**struct sk\_buff \* skb** buffer to reference

### Description

Makes another reference to a socket buffer and returns a pointer to the buffer.

int **skb\_cloned**(const struct [sk\\_buff](#) \* *skb*)  
is the buffer a clone

### Parameters

**const struct sk\_buff \* skb** buffer to check

### Description

Returns true if the buffer was generated with [skb\\_clone\(\)](#) and is one of multiple shared copies of the buffer. Cloned buffers are shared data so must not be written to under normal circumstances.

int **skb\_header\_cloned**(const struct [sk\\_buff](#) \* *skb*)  
is the header a clone

### Parameters

**const struct sk\_buff \* skb** buffer to check

### Description

Returns true if modifying the header part of the buffer requires the data to be copied.

void **skb\_header\_release**(struct [sk\\_buff](#) \* *skb*)  
release reference to header

### Parameters

**struct sk\_buff \* skb** buffer to operate on

### Description

Drop a reference to the header part of the buffer. This is done by acquiring a payload reference. You must not read from the header part of `skb->data` after this.

### Note

Check if you can use `__skb_header_release()` instead.

void **\_\_skb\_header\_release**(struct *sk\_buff* \* *skb*)  
release reference to header

### Parameters

**struct sk\_buff \* skb** buffer to operate on

### Description

Variant of `skb_header_release()` assuming `skb` is private to caller. We can avoid one atomic operation.

int **skb\_shared**(const struct *sk\_buff* \* *skb*)  
is the buffer shared

### Parameters

**const struct sk\_buff \* skb** buffer to check

### Description

Returns true if more than one person has a reference to this buffer.

struct *sk\_buff* \* **skb\_share\_check**(struct *sk\_buff* \* *skb*, gfp\_t *pri*)  
check if buffer is shared and if so clone it

### Parameters

**struct sk\_buff \* skb** buffer to check

**gfp\_t pri** priority for memory allocation

### Description

If the buffer is shared the buffer is cloned and the old copy drops a reference. A new clone with a single reference is returned. If the buffer is not shared the original buffer is returned. When being called from interrupt status or with spinlocks held `pri` must be `GFP_ATOMIC`.

NULL is returned on a memory allocation failure.

struct *sk\_buff* \* **skb\_unshare**(struct *sk\_buff* \* *skb*, gfp\_t *pri*)  
make a copy of a shared buffer

### Parameters

**struct sk\_buff \* skb** buffer to check

**gfp\_t pri** priority for memory allocation

### Description

If the socket buffer is a clone then this function creates a new copy of the data, drops a reference count on the old copy and returns the new copy with the reference count at 1. If the buffer is not a clone the original buffer is returned. When called with a spinlock held or from interrupt state `pri` must be `GFP_ATOMIC`

NULL is returned on a memory allocation failure.

struct *sk\_buff* \* **skb\_peek**(const struct *sk\_buff\_head* \* *list\_*)  
peek at the head of an *sk\_buff\_head*

### Parameters



**const struct sk\_buff\_head \* list\_** list to peek at

### Description

Peek an *sk\_buff*. Unlike most other operations you *MUST* be careful with this one. A peek leaves the buffer on the list and someone else may run off with it. You must hold the appropriate locks or have a private queue to do this.

Returns NULL for an empty list or a pointer to the head element. The reference count is not incremented and the reference is therefore volatile. Use with caution.

struct *sk\_buff* \* **skb\_peek\_next**(struct *sk\_buff* \* *skb*, const struct sk\_buff\_head \* *list\_*)  
peek skb following the given one from a queue

### Parameters

**struct sk\_buff \* skb** skb to start from

**const struct sk\_buff\_head \* list\_** list to peek at

### Description

Returns NULL when the end of the list is met or a pointer to the next element. The reference count is not incremented and the reference is therefore volatile. Use with caution.

struct *sk\_buff* \* **skb\_peek\_tail**(const struct sk\_buff\_head \* *list\_*)  
peek at the tail of an sk\_buff\_head

### Parameters

**const struct sk\_buff\_head \* list\_** list to peek at

### Description

Peek an *sk\_buff*. Unlike most other operations you *MUST* be careful with this one. A peek leaves the buffer on the list and someone else may run off with it. You must hold the appropriate locks or have a private queue to do this.

Returns NULL for an empty list or a pointer to the tail element. The reference count is not incremented and the reference is therefore volatile. Use with caution.

**\_\_u32 skb\_queue\_len**(const struct sk\_buff\_head \* *list\_*)  
get queue length

### Parameters

**const struct sk\_buff\_head \* list\_** list to measure

### Description

Return the length of an *sk\_buff* queue.

void **\_\_skb\_queue\_head\_init**(struct sk\_buff\_head \* *list*)  
initialize non-spinlock portions of sk\_buff\_head

### Parameters

**struct sk\_buff\_head \* list** queue to initialize

### Description

This initializes only the list and queue length aspects of an sk\_buff\_head object. This allows to initialize the list aspects of an sk\_buff\_head without reinitializing things like the spinlock. It can also be used for on-stack sk\_buff\_head objects where the spinlock is known to not be used.

void **skb\_queue\_splice**(const struct sk\_buff\_head \* *list*, struct sk\_buff\_head \* *head*)  
join two skb lists, this is designed for stacks

### Parameters

**const struct sk\_buff\_head \* list** the new list to add

**struct sk\_buff\_head \* head** the place to add it in the first list

void **skb\_queue\_splice\_init**(struct sk\_buff\_head \* *list*, struct sk\_buff\_head \* *head*)  
join two skb lists and reinitialise the emptied list

**Parameters**

**struct sk\_buff\_head \* list** the new list to add  
**struct sk\_buff\_head \* head** the place to add it in the first list

**Description**

The list at **list** is reinitialised

void **skb\_queue\_splice\_tail**(const struct sk\_buff\_head \* *list*, struct sk\_buff\_head \* *head*)  
join two skb lists, each list being a queue

**Parameters**

**const struct sk\_buff\_head \* list** the new list to add  
**struct sk\_buff\_head \* head** the place to add it in the first list  
void **skb\_queue\_splice\_tail\_init**(struct sk\_buff\_head \* *list*, struct sk\_buff\_head \* *head*)  
join two skb lists and reinitialise the emptied list

**Parameters**

**struct sk\_buff\_head \* list** the new list to add  
**struct sk\_buff\_head \* head** the place to add it in the first list

**Description**

Each of the lists is a queue. The list at **list** is reinitialised

void **\_\_skb\_queue\_after**(struct sk\_buff\_head \* *list*, struct *sk\_buff* \* *prev*, struct *sk\_buff* \* *newsk*)  
queue a buffer at the list head

**Parameters**

**struct sk\_buff\_head \* list** list to use  
**struct sk\_buff \* prev** place after this buffer  
**struct sk\_buff \* newsk** buffer to queue

**Description**

Queue a buffer into the middle of a list. This function takes no locks and you must therefore hold required locks before calling it.

A buffer cannot be placed on two lists at the same time.

void **skb\_queue\_head**(struct sk\_buff\_head \* *list*, struct *sk\_buff* \* *newsk*)  
queue a buffer at the list head

**Parameters**

**struct sk\_buff\_head \* list** list to use  
**struct sk\_buff \* newsk** buffer to queue

**Description**

Queue a buffer at the start of a list. This function takes no locks and you must therefore hold required locks before calling it.

A buffer cannot be placed on two lists at the same time.

void **skb\_queue\_tail**(struct sk\_buff\_head \* *list*, struct *sk\_buff* \* *newsk*)  
queue a buffer at the list tail

**Parameters**

**struct sk\_buff\_head \* list** list to use

**struct sk\_buff \* newsk** buffer to queue

### Description

Queue a buffer at the end of a list. This function takes no locks and you must therefore hold required locks before calling it.

A buffer cannot be placed on two lists at the same time.

**struct sk\_buff \* skb\_dequeue**(**struct sk\_buff\_head \* list**)  
remove from the head of the queue

### Parameters

**struct sk\_buff\_head \* list** list to dequeue from

### Description

Remove the head of the list. This function does not take any locks so must be used with appropriate locks held only. The head item is returned or NULL if the list is empty.

**struct sk\_buff \* skb\_dequeue\_tail**(**struct sk\_buff\_head \* list**)  
remove from the tail of the queue

### Parameters

**struct sk\_buff\_head \* list** list to dequeue from

### Description

Remove the tail of the list. This function does not take any locks so must be used with appropriate locks held only. The tail item is returned or NULL if the list is empty.

**void \_\_skb\_fill\_page\_desc**(**struct sk\_buff \* skb**, **int i**, **struct page \* page**, **int off**, **int size**)  
initialise a paged fragment in an skb

### Parameters

**struct sk\_buff \* skb** buffer containing fragment to be initialised

**int i** paged fragment index to initialise

**struct page \* page** the page to use for this fragment

**int off** the offset to the data with **page**

**int size** the length of the data

### Description

Initialises the **i**'th fragment of **skb** to point to **size** bytes at offset **off** within **page**.

Does not take any additional reference on the fragment.

**void skb\_fill\_page\_desc**(**struct sk\_buff \* skb**, **int i**, **struct page \* page**, **int off**, **int size**)  
initialise a paged fragment in an skb

### Parameters

**struct sk\_buff \* skb** buffer containing fragment to be initialised

**int i** paged fragment index to initialise

**struct page \* page** the page to use for this fragment

**int off** the offset to the data with **page**

**int size** the length of the data

### Description

As per **\_\_skb\_fill\_page\_desc()** - initialises the **i**'th fragment of **skb** to point to **size** bytes at offset **off** within **page**. In addition updates **skb** such that **i** is the last fragment.

Does not take any additional reference on the fragment.

unsigned int **skb\_headroom**(const struct *sk\_buff* \* *skb*)  
bytes at buffer head

#### Parameters

const struct *sk\_buff* \* *skb* buffer to check

#### Description

Return the number of bytes of free space at the head of an *sk\_buff*.

int **skb\_tailroom**(const struct *sk\_buff* \* *skb*)  
bytes at buffer end

#### Parameters

const struct *sk\_buff* \* *skb* buffer to check

#### Description

Return the number of bytes of free space at the tail of an *sk\_buff*

int **skb\_availroom**(const struct *sk\_buff* \* *skb*)  
bytes at buffer end

#### Parameters

const struct *sk\_buff* \* *skb* buffer to check

#### Description

Return the number of bytes of free space at the tail of an *sk\_buff* allocated by *sk\_stream\_alloc()*

void **skb\_reserve**(struct *sk\_buff* \* *skb*, int *len*)  
adjust headroom

#### Parameters

struct *sk\_buff* \* *skb* buffer to alter

int *len* bytes to move

#### Description

Increase the headroom of an empty *sk\_buff* by reducing the tail room. This is only allowed for an empty buffer.

void **skb\_tailroom\_reserve**(struct *sk\_buff* \* *skb*, unsigned int *mtu*, unsigned int *needed\_tailroom*)  
adjust reserved\_tailroom

#### Parameters

struct *sk\_buff* \* *skb* buffer to alter

unsigned int *mtu* maximum amount of headlen permitted

unsigned int *needed\_tailroom* minimum amount of reserved\_tailroom

#### Description

Set *reserved\_tailroom* so that *headlen* can be as large as possible but not larger than *mtu* and *tailroom* cannot be smaller than *needed\_tailroom*. The required headroom should already have been reserved before using this function.

void **pskb\_trim\_unique**(struct *sk\_buff* \* *skb*, unsigned int *len*)  
remove end from a paged unique (not cloned) buffer

#### Parameters

struct *sk\_buff* \* *skb* buffer to alter

unsigned int *len* new length

## Description

This is identical to `pskb_trim` except that the caller knows that the `skb` is not cloned so we should never get an error due to out- of-memory.

```
void skb_orphan(struct sk_buff * skb)  
    orphan a buffer
```

## Parameters

**struct sk\_buff \* skb** buffer to orphan

## Description

If a buffer currently has an owner then we call the owner's destructor function and make the **skb** unowned. The buffer continues to exist but is no longer charged to its former owner.

```
int skb_orphan_frags(struct sk_buff * skb, gfp_t gfp_mask)  
    orphan the frags contained in a buffer
```

## Parameters

**struct sk\_buff \* skb** buffer to orphan frags from

**gfp\_t gfp\_mask** allocation mask for replacement pages

## Description

For each frag in the SKB which needs a destructor (i.e. has an owner) create a copy of that frag and release the original page by calling the destructor.

```
void skb_queue_purge(struct sk_buff_head * list)  
    empty a list
```

## Parameters

**struct sk\_buff\_head \* list** list to empty

## Description

Delete all buffers on an *sk\_buff* list. Each buffer is removed from the list and one reference dropped. This function does not take the list lock and the caller must hold the relevant locks to use it.

```
struct sk_buff * netdev_alloc_skb(struct net_device * dev, unsigned int length)  
    allocate an skb for rx on a specific device
```

## Parameters

**struct net\_device \* dev** network device to receive on

**unsigned int length** length to allocate

## Description

Allocate a new *sk\_buff* and assign it a usage count of one. The buffer has unspecified headroom built in. Users should allocate the headroom they think they need without accounting for the built in space. The built in space is used for optimisations.

NULL is returned if there is no free memory. Although this function allocates memory it can be called from an interrupt.

```
struct page * __dev_alloc_pages(gfp_t gfp_mask, unsigned int order)  
    allocate page for network Rx
```

## Parameters

**gfp\_t gfp\_mask** allocation priority. Set `__GFP_NOMEMALLOC` if not for network Rx

**unsigned int order** size of the allocation

### Description

Allocate a new page.

NULL is returned if there is no free memory.

struct page \* **\_\_dev\_alloc\_page**(gfp\_t *gfp\_mask*)  
allocate a page for network Rx

### Parameters

**gfp\_t gfp\_mask** allocation priority. Set `__GFP_NOMEMALLOC` if not for network Rx

### Description

Allocate a new page.

NULL is returned if there is no free memory.

void **skb\_propagate\_pfmemalloc**(struct page \* *page*, struct *sk\_buff* \* *skb*)  
Propagate pfmemalloc if skb is allocated after RX page

### Parameters

**struct page \* page** The page that was allocated from `skb_alloc_page`

**struct sk\_buff \* skb** The skb that may need pfmemalloc set

struct page \* **skb\_frag\_page**(const skb\_frag\_t \* *frag*)  
retrieve the page referred to by a paged fragment

### Parameters

**const skb\_frag\_t \* frag** the paged fragment

### Description

Returns the struct page associated with **frag**.

void **\_\_skb\_frag\_ref**(skb\_frag\_t \* *frag*)  
take an addition reference on a paged fragment.

### Parameters

**skb\_frag\_t \* frag** the paged fragment

### Description

Takes an additional reference on the paged fragment **frag**.

void **skb\_frag\_ref**(struct *sk\_buff* \* *skb*, int *f*)  
take an addition reference on a paged fragment of an skb.

### Parameters

**struct sk\_buff \* skb** the buffer

**int f** the fragment offset.

### Description

Takes an additional reference on the **f**'th paged fragment of **skb**.

void **\_\_skb\_frag\_unref**(skb\_frag\_t \* *frag*)  
release a reference on a paged fragment.

### Parameters

**skb\_frag\_t \* frag** the paged fragment

### Description

Releases a reference on the paged fragment **frag**.

void **skb\_frag\_unref**(struct *sk\_buff* \* *skb*, int *f*)  
release a reference on a paged fragment of an skb.

**Parameters**

**struct sk\_buff \* skb** the buffer

**int f** the fragment offset

**Description**

Releases a reference on the **f**'th paged fragment of **skb**.

void \* **skb\_frag\_address**(const skb\_frag\_t \* *frag*)  
gets the address of the data contained in a paged fragment

**Parameters**

**const skb\_frag\_t \* frag** the paged fragment buffer

**Description**

Returns the address of the data within **frag**. The page must already be mapped.

void \* **skb\_frag\_address\_safe**(const skb\_frag\_t \* *frag*)  
gets the address of the data contained in a paged fragment

**Parameters**

**const skb\_frag\_t \* frag** the paged fragment buffer

**Description**

Returns the address of the data within **frag**. Checks that the page is mapped and returns NULL otherwise.

void **\_\_skb\_frag\_set\_page**(skb\_frag\_t \* *frag*, struct page \* *page*)  
sets the page contained in a paged fragment

**Parameters**

**skb\_frag\_t \* frag** the paged fragment

**struct page \* page** the page to set

**Description**

Sets the fragment **frag** to contain **page**.

void **skb\_frag\_set\_page**(struct *sk\_buff* \* *skb*, int *f*, struct page \* *page*)  
sets the page contained in a paged fragment of an skb

**Parameters**

**struct sk\_buff \* skb** the buffer

**int f** the fragment offset

**struct page \* page** the page to set

**Description**

Sets the **f**'th fragment of **skb** to contain **page**.

dma\_addr\_t **skb\_frag\_dma\_map**(struct device \* *dev*, const skb\_frag\_t \* *frag*, size\_t *offset*, size\_t *size*,  
enum dma\_data\_direction *dir*)  
maps a paged fragment via the DMA API

**Parameters**

**struct device \* dev** the device to map the fragment to

**const skb\_frag\_t \* frag** the paged fragment to map

**size\_t offset** the offset within the fragment (starting at the fragment's own offset)

**size\_t size** the number of bytes to map

**enum dma\_data\_direction dir** the direction of the mapping (PCI\_DMA\_\*)

### Description

Maps the page associated with **frag** to **device**.

int **skb\_clone\_writable**(const struct *sk\_buff* \* *skb*, unsigned int *len*)  
is the header of a clone writable

### Parameters

**const struct sk\_buff \* skb** buffer to check

**unsigned int len** length up to which to write

### Description

Returns true if modifying the header part of the cloned buffer does not requires the data to be copied.

int **skb\_cow**(struct *sk\_buff* \* *skb*, unsigned int *headroom*)  
copy header of skb when it is required

### Parameters

**struct sk\_buff \* skb** buffer to cow

**unsigned int headroom** needed headroom

### Description

If the skb passed lacks sufficient headroom or its data part is shared, data is reallocated. If reallocation fails, an error is returned and original skb is not changed.

The result is skb with writable area *skb->head...skb->tail* and at least **headroom** of space at head.

int **skb\_cow\_head**(struct *sk\_buff* \* *skb*, unsigned int *headroom*)  
skb\_cow but only making the head writable

### Parameters

**struct sk\_buff \* skb** buffer to cow

**unsigned int headroom** needed headroom

### Description

This function is identical to *skb\_cow* except that we replace the *skb\_cloned* check by *skb\_header\_cloned*. It should be used when you only need to push on some header and do not need to modify the data.

int **skb\_padto**(struct *sk\_buff* \* *skb*, unsigned int *len*)  
pad an skbuff up to a minimal size

### Parameters

**struct sk\_buff \* skb** buffer to pad

**unsigned int len** minimal length

### Description

Pads up a buffer to ensure the trailing bytes exist and are blanked. If the buffer already contains sufficient data it is untouched. Otherwise it is extended. Returns zero on success. The skb is freed on error.

int **\_\_skb\_put\_padto**(struct *sk\_buff* \* *skb*, unsigned int *len*, bool *free\_on\_error*)  
increase size and pad an skbuff up to a minimal size

### Parameters



**struct sk\_buff \* skb** buffer to pad  
**unsigned int len** minimal length  
**bool free\_on\_error** free buffer on error

#### Description

Pads up a buffer to ensure the trailing bytes exist and are blanked. If the buffer already contains sufficient data it is untouched. Otherwise it is extended. Returns zero on success. The skb is freed on error if **free\_on\_error** is true.

int **skb\_put\_padto**(struct *sk\_buff* \* *skb*, unsigned int *len*)  
increase size and pad an skb up to a minimal size

#### Parameters

**struct sk\_buff \* skb** buffer to pad  
**unsigned int len** minimal length

#### Description

Pads up a buffer to ensure the trailing bytes exist and are blanked. If the buffer already contains sufficient data it is untouched. Otherwise it is extended. Returns zero on success. The skb is freed on error.

int **skb\_linearize**(struct *sk\_buff* \* *skb*)  
convert paged skb to linear one

#### Parameters

**struct sk\_buff \* skb** buffer to linearize

#### Description

If there is no free memory -ENOMEM is returned, otherwise zero is returned and the old skb data released.

bool **skb\_has\_shared\_frag**(const struct *sk\_buff* \* *skb*)  
can any frag be overwritten

#### Parameters

**const struct sk\_buff \* skb** buffer to test

#### Description

Return true if the skb has at least one frag that might be modified by an external entity (as in vm-splice()/sendfile())

int **skb\_linearize\_cow**(struct *sk\_buff* \* *skb*)  
make sure skb is linear and writable

#### Parameters

**struct sk\_buff \* skb** buffer to process

#### Description

If there is no free memory -ENOMEM is returned, otherwise zero is returned and the old skb data released.

void **skb\_postpull\_rcsum**(struct *sk\_buff* \* *skb*, const void \* *start*, unsigned int *len*)  
update checksum for received skb after pull

#### Parameters

**struct sk\_buff \* skb** buffer to update  
**const void \* start** start of data before pull  
**unsigned int len** length of data pulled

### Description

After doing a pull on a received packet, you need to call this to update the CHECKSUM\_COMPLETE checksum, or set ip\_summed to CHECKSUM\_NONE so that it can be recomputed from scratch.

void **skb\_postpush\_rcsum**(struct *sk\_buff* \* *skb*, const void \* *start*, unsigned int *len*)  
update checksum for received skb after push

### Parameters

**struct sk\_buff \* skb** buffer to update  
**const void \* start** start of data after push  
**unsigned int len** length of data pushed

### Description

After doing a push on a received packet, you need to call this to update the CHECKSUM\_COMPLETE checksum.

void \* **skb\_push\_rcsum**(struct *sk\_buff* \* *skb*, unsigned int *len*)  
push skb and update receive checksum

### Parameters

**struct sk\_buff \* skb** buffer to update  
**unsigned int len** length of data pulled

### Description

This function performs an skb\_push on the packet and updates the CHECKSUM\_COMPLETE checksum. It should be used on receive path processing instead of skb\_push unless you know that the checksum difference is zero (e.g., a valid IP header) or you are setting ip\_summed to CHECKSUM\_NONE.

int **pskb\_trim\_rcsum**(struct *sk\_buff* \* *skb*, unsigned int *len*)  
trim received skb and update checksum

### Parameters

**struct sk\_buff \* skb** buffer to trim  
**unsigned int len** new length

### Description

This is exactly the same as pskb\_trim except that it ensures the checksum of received packets are still valid after the operation.

bool **skb\_needs\_linearize**(struct *sk\_buff* \* *skb*, netdev\_features\_t *features*)  
check if we need to linearize a given skb depending on the given device features.

### Parameters

**struct sk\_buff \* skb** socket buffer to check  
**netdev\_features\_t features** net device features

### Description

Returns true if either: 1. skb has frag\_list and the device doesn't support FRAGLIST, or 2. skb is fragmented and the device does not support SG.

void **skb\_get\_timestamp**(const struct *sk\_buff* \* *skb*, struct timeval \* *stamp*)  
get timestamp from a skb

### Parameters

**const struct sk\_buff \* skb** skb to get stamp from  
**struct timeval \* stamp** pointer to struct timeval to store stamp in

## Description

Timestamps are stored in the skb as offsets to a base timestamp. This function converts the offset back to a struct timeval and stores it in stamp.

```
void skb_complete_tx_timestamp(struct sk_buff * skb, struct skb_shared_hwtstamps * hwt-  
                                stamps)  
    deliver cloned skb with tx timestamps
```

## Parameters

**struct sk\_buff \* skb** clone of the the original outgoing packet

**struct skb\_shared\_hwtstamps \* hwtstamps** hardware time stamps

## Description

PHY drivers may accept clones of transmitted packets for timestamping via their phy\_driver.txtstamp method. These drivers must call this function to return the skb back to the stack with a timestamp.

```
void skb_tstamp_tx(struct sk_buff * orig_skb, struct skb_shared_hwtstamps * hwtstamps)  
    queue clone of skb with send time stamps
```

## Parameters

**struct sk\_buff \* orig\_skb** the original outgoing packet

**struct skb\_shared\_hwtstamps \* hwtstamps** hardware time stamps, may be NULL if not available

## Description

If the skb has a socket associated, then this function clones the skb (thus sharing the actual data and optional structures), stores the optional hardware time stamping information (if non NULL) or generates a software time stamp (otherwise), then queues the clone to the error queue of the socket. Errors are silently ignored.

```
void skb_tx_timestamp(struct sk_buff * skb)  
    Driver hook for transmit timestamping
```

## Parameters

**struct sk\_buff \* skb** A socket buffer.

## Description

Ethernet MAC Drivers should call this function in their hard\_xmit() function immediately before giving the sk\_buff to the MAC hardware.

Specifically, one should make absolutely sure that this function is called before TX completion of this packet can trigger. Otherwise the packet could potentially already be freed.

```
void skb_complete_wifi_ack(struct sk_buff * skb, bool acked)  
    deliver skb with wifi status
```

## Parameters

**struct sk\_buff \* skb** the original outgoing packet

**bool acked** ack status

```
__sum16 skb_checksum_complete(struct sk_buff * skb)  
    Calculate checksum of an entire packet
```

## Parameters

**struct sk\_buff \* skb** packet to process

## Description

This function calculates the checksum over the entire packet plus the value of skb->csum. The latter can be used to supply the checksum of a pseudo header as used by TCP/UDP. It returns the checksum.

For protocols that contain complete checksums such as ICMP/TCP/UDP, this function can be used to verify that checksum on received packets. In that case the function should return zero if the checksum is correct. In particular, this function will return zero if `skb->ip_summed` is `CHECKSUM_UNNECESSARY` which indicates that the hardware has already verified the correctness of the checksum.

void **skb\_checksum\_none\_assert**(const struct *sk\_buff* \* *skb*)  
make sure `skb` `ip_summed` is `CHECKSUM_NONE`

### Parameters

const struct *sk\_buff* \* *skb* *skb* to check

### Description

fresh skbs have their `ip_summed` set to `CHECKSUM_NONE`. Instead of forcing `ip_summed` to `CHECKSUM_NONE`, we can use this helper, to document places where we make this assertion.

bool **skb\_head\_is\_locked**(const struct *sk\_buff* \* *skb*)  
Determine if the `skb->head` is locked down

### Parameters

const struct *sk\_buff* \* *skb* *skb* to check

### Description

The head on skbs build around a head frag can be removed if they are not cloned. This function returns true if the `skb` head is locked down due to either being allocated via `kmalloc`, or by being a clone with multiple references to the head.

unsigned int **skb\_gso\_network\_seglen**(const struct *sk\_buff* \* *skb*)  
Return length of individual segments of a gso packet

### Parameters

const struct *sk\_buff* \* *skb* GSO *skb*

### Description

`skb_gso_network_seglen` is used to determine the real size of the individual segments, including Layer3 (IP, IPv6) and L4 headers (TCP/UDP).

The MAC/L2 header is not accounted for.

struct **sock\_common**  
minimal network layer representation of sockets

### Definition

```
struct sock_common {
    union {unnamed_union};
};
```

### Members

**{unnamed\_union}** anonymous

### Description

This is the minimal network layer representation of sockets, the header for struct `sock` and struct `inet_twait_sock`.

struct **sock**  
network layer representation of sockets

### Definition

```

struct sock {
    struct sock_common __sk_common;
#define sk_node                __sk_common.skc_node
#define sk_nulls_node         __sk_common.skc_nulls_node
#define sk_refcnt              __sk_common.skc_refcnt
#define sk_tx_queue_mapping   __sk_common.skc_tx_queue_mapping
#define sk_dontcopy_begin     __sk_common.skc_dontcopy_begin
#define sk_dontcopy_end       __sk_common.skc_dontcopy_end
#define sk_hash                __sk_common.skc_hash
#define sk_portpair           __sk_common.skc_portpair
#define sk_num                 __sk_common.skc_num
#define sk_dport               __sk_common.skc_dport
#define sk_addrpair           __sk_common.skc_addrpair
#define sk_daddr               __sk_common.skc_daddr
#define sk_rcv_saddr          __sk_common.skc_rcv_saddr
#define sk_family              __sk_common.skc_family
#define sk_state               __sk_common.skc_state
#define sk_reuse               __sk_common.skc_reuse
#define sk_reuseport          __sk_common.skc_reuseport
#define sk_ipv6only            __sk_common.skc_ipv6only
#define sk_net_refcnt          __sk_common.skc_net_refcnt
#define sk_bound_dev_if       __sk_common.skc_bound_dev_if
#define sk_bind_node          __sk_common.skc_bind_node
#define sk_prot                __sk_common.skc_prot
#define sk_net                 __sk_common.skc_net
#define sk_v6_daddr           __sk_common.skc_v6_daddr
#define sk_v6_rcv_saddr       __sk_common.skc_v6_rcv_saddr
#define sk_cookie              __sk_common.skc_cookie
#define sk_incoming_cpu        __sk_common.skc_incoming_cpu
#define sk_flags               __sk_common.skc_flags
#define sk_rxhash              __sk_common.skc_rxhash
    socket_lock_t sk_lock;
    atomic_t sk_drops;
    int sk_rcvlowat;
    struct sk_buff_head sk_error_queue;
    struct sk_buff_head sk_receive_queue;
    struct {unnamed_struct};
#ifdef CONFIG_XFRM
    struct xfrm_policy __rcu * sk_policy;
#endif
    struct dst_entry * sk_rx_dst;
    struct dst_entry __rcu * sk_dst_cache;
    atomic_t sk_omem_alloc;
    int sk_sndbuf;
    int sk_wmem_queued;
    refcount_t sk_wmem_alloc;
    unsigned long sk_tsq_flags;
    struct sk_buff * sk_send_head;
    struct sk_buff_head sk_write_queue;
    __s32 sk_peek_off;
    int sk_write_pending;
    __u32 sk_dst_pending_confirm;
    u32 sk_pacing_status;
    long sk_sndtimeo;
    struct timer_list sk_timer;
    __u32 sk_priority;
    __u32 sk_mark;
    u32 sk_pacing_rate;
    u32 sk_max_pacing_rate;
    struct page_frag sk_frag;
    netdev_features_t sk_route_caps;
    netdev_features_t sk_route_nocaps;
    int sk_gso_type;

```

```
    unsigned int sk_gso_max_size;
    gfp_t sk_allocation;
    __u32 sk_txhash;
    unsigned int __sk_flags_offset;
#ifdef __BIG_ENDIAN_BITFIELD
#define SK_FL_PROTO_SHIFT 16
#define SK_FL_PROTO_MASK 0x00ff0000
#define SK_FL_TYPE_SHIFT 0
#define SK_FL_TYPE_MASK 0x0000ffff
#else
#define SK_FL_PROTO_SHIFT 8
#define SK_FL_PROTO_MASK 0x0000ff00
#define SK_FL_TYPE_SHIFT 16
#define SK_FL_TYPE_MASK 0xffff0000
#endif
    unsigned int sk_padding:1;
    unsigned int sk_kern_sock:1;
    unsigned int sk_no_check_tx:1;
    unsigned int sk_no_check_rx:1;
    unsigned int sk_userlocks:4;
    unsigned int sk_protocol:8;
    unsigned int sk_type:16;
#define SK_PROTOCOL_MAX U8_MAX
    u16 sk_gso_max_segs;
    unsigned long sk_lingertime;
    struct proto * sk_prot_creator;
    rwlock_t sk_callback_lock;
    int sk_err;
    int sk_err_soft;
    u32 sk_ack_backlog;
    u32 sk_max_ack_backlog;
    kuid_t sk_uid;
    struct pid * sk_peer_pid;
    const struct cred * sk_peer_cred;
    long sk_rcvtimeo;
    ktime_t sk_stamp;
    u16 sk_tsflags;
    u8 sk_shutdown;
    u32 sk_tskey;
    atomic_t sk_zckey;
    struct socket * sk_socket;
    void * sk_user_data;
#ifdef CONFIG_SECURITY
    void * sk_security;
#endif
#endif
    struct sock_cgroup_data sk_cgrp_data;
    struct mem_cgroup * sk_memcg;
    void (* sk_state_change) (struct sock *sk);
    void (* sk_data_ready) (struct sock *sk);
    void (* sk_write_space) (struct sock *sk);
    void (* sk_error_report) (struct sock *sk);
    int (* sk_backlog_rcv) (struct sock *sk, struct sk_buff *skb);
    void (* sk_destruct) (struct sock *sk);
    struct sock_reuseport __rcu * sk_reuseport_cb;
    struct rcu_head sk_rcu;
};
```

## Members

**\_\_sk\_common** shared layout with `inet_timewait_sock`

**sk\_lock** synchronizer

**sk\_drops** raw/udp drops counter

**sk\_rcvlowat** SO\_RCVLOWAT setting

**sk\_error\_queue** rarely used

**sk\_receive\_queue** incoming packets

**{unnamed\_struct}** anonymous

**sk\_policy** flow policy

**sk\_rx\_dst** receive input route used by early demux

**sk\_dst\_cache** destination cache

**sk\_omem\_alloc** “o” is “option” or “other”

**sk\_sndbuf** size of send buffer in bytes

**sk\_wmem\_queued** persistent queue size

**sk\_wmem\_alloc** transmit queue bytes committed

**sk\_tsq\_flags** TCP Small Queues flags

**sk\_send\_head** front of stuff to transmit

**sk\_write\_queue** Packet sending queue

**sk\_peek\_off** current peek\_offset value

**sk\_write\_pending** a write to stream socket waits to start

**sk\_dst\_pending\_confirm** need to confirm neighbour

**sk\_pacing\_status** Pacing status (requested, handled by sch\_fq)

**sk\_sndtimeo** SO\_SNDTIMEO setting

**sk\_timer** sock cleanup timer

**sk\_priority** SO\_PRIORITY setting

**sk\_mark** generic packet mark

**sk\_pacing\_rate** Pacing rate (if supported by transport/packet scheduler)

**sk\_max\_pacing\_rate** Maximum pacing rate (SO\_MAX\_PACING\_RATE)

**sk\_frag** cached page frag

**sk\_route\_caps** route capabilities (e.g. NETIF\_F\_TS0)

**sk\_route\_nocaps** forbidden route capabilities (e.g. NETIF\_F\_GSO\_MASK)

**sk\_gso\_type** GSO type (e.g. SKB\_GSO\_TCPV4)

**sk\_gso\_max\_size** Maximum GSO segment size to build

**sk\_allocation** allocation mode

**sk\_txhash** computed flow hash for use on transmit

**\_\_sk\_flags\_offset** empty field used to determine location of bitfield

**sk\_padding** unused element for alignment

**sk\_kern\_sock** True if sock is using kernel lock classes

**sk\_no\_check\_tx** SO\_NO\_CHECK setting, set checksum in TX packets

**sk\_no\_check\_rx** allow zero checksum in RX packets

**sk\_userlocks** SO\_SNDBUF and SO\_RCVBUF settings

**sk\_protocol** which protocol this socket belongs in this network family

**sk\_type** socket type (SOCK\_STREAM, etc)

**sk\_gso\_max\_segs** Maximum number of GSO segments

**sk\_lingertime** SO\_LINGER l\_linger setting

**sk\_prot\_creator** sk\_prot of original sock creator (see `ipv6_setsockopt`, `IPV6_ADDRRFORM` for instance)

**sk\_callback\_lock** used with the callbacks in the end of this struct

**sk\_err** last error

**sk\_err\_soft** errors that don't cause failure but are the cause of a persistent failure not just 'timed out'

**sk\_ack\_backlog** current listen backlog

**sk\_max\_ack\_backlog** listen backlog set in `listen()`

**sk\_uid** user id of owner

**sk\_peer\_pid** struct pid for this socket's peer

**sk\_peer\_cred** SO\_PEERCRED setting

**sk\_rcvtimeo** SO\_RCVTIMEO setting

**sk\_stamp** time stamp of last packet received

**sk\_tsflags** SO\_TIMESTAMPING socket options

**sk\_shutdown** mask of `SEND_SHUTDOWN` and/or `RCV_SHUTDOWN`

**sk\_tskey** counter to disambiguate concurrent timestamp requests

**sk\_zckey** counter to order `MSG_ZEROCOPY` notifications

**sk\_socket** Identd and reporting IO signals

**sk\_user\_data** RPC layer private data

**sk\_security** used by security modules

**sk\_cgrp\_data** cgroup data for this cgroup

**sk\_memcg** this socket's memory cgroup association

**sk\_state\_change** callback to indicate change in the state of the sock

**sk\_data\_ready** callback to indicate there is data to be processed

**sk\_write\_space** callback to indicate there is bf sending space available

**sk\_error\_report** callback to indicate errors (e.g. `MSG_ERRQUEUE`)

**sk\_backlog\_rcv** callback to process the backlog

**sk\_destruct** called at sock freeing time, i.e. when `all_refcnt == 0`

**sk\_reuseport\_cb** reuseport group container

**sk\_rcu** used during RCU grace period

**sk\_for\_each\_entry\_offset\_rcu**(*tpos*, *pos*, *head*, *offset*)  
iterate over a list at a given struct offset

### Parameters

**tpos** the type \* to use as a loop cursor.

**pos** the struct `hlist_node` to use as a loop cursor.

**head** the head for your list.

**offset** offset of `hlist_node` within the struct.

void **unlock\_sock\_fast**(struct [sock](#) \* *sk*, bool *slow*)  
complement of `lock_sock_fast`

### Parameters



```
struct sock * sk socket
```

```
bool slow slow mode
```

### Description

fast unlock socket for user context. If slow mode is on, we call regular `release_sock()`

```
int sk_wmem_alloc_get(const struct sock * sk)
    returns write allocations
```

### Parameters

```
const struct sock * sk socket
```

### Description

Returns `sk_wmem_alloc` minus initial offset of one

```
int sk_rmem_alloc_get(const struct sock * sk)
    returns read allocations
```

### Parameters

```
const struct sock * sk socket
```

### Description

Returns `sk_rmem_alloc`

```
bool sk_has_allocations(const struct sock * sk)
    check if allocations are outstanding
```

### Parameters

```
const struct sock * sk socket
```

### Description

Returns true if socket has write or read allocations

```
bool skwq_has_sleeper(struct socket_wq * wq)
    check if there are any waiting processes
```

### Parameters

```
struct socket_wq * wq struct socket_wq
```

### Description

Returns true if `socket_wq` has waiting processes

The purpose of the `skwq_has_sleeper` and `sock_poll_wait` is to wrap the memory barrier call. They were added due to the race found within the tcp code.

Consider following tcp code paths:

CPU1	CPU2
<code>sys_select</code>	<code>receive packet</code>
<code>...</code>	<code>...</code>
<code>__add_wait_queue</code>	<code>update tp-&gt;rcv_nxt</code>
<code>...</code>	<code>...</code>
<code>tp-&gt;rcv_nxt check</code>	<code>sock_def_readable</code>
<code>...</code>	{
<code>schedule</code>	<code>:c:func:`rcu_read_lock()``;</code>
	<code>wq = rcu_dereference(sk-&gt;sk_wq);</code>
	<code>if (wq &amp;&amp; waitqueue_active(:c:type:`wq-&gt;wait &lt;wq&gt;`))</code>
	<code>wake_up_interruptible(:c:type:`wq-&gt;wait &lt;wq&gt;`)</code>
	<code>...</code>
	}

The race for tcp fires when the `__add_wait_queue` changes done by CPU1 stay in its cache, and so does the `tp->rcv_nxt` update on CPU2 side. The CPU1 could then endup calling `schedule` and sleep forever if there are no more data on the socket.

void **sock\_poll\_wait**(struct file \* *filp*, wait\_queue\_head\_t \* *wait\_address*, poll\_table \* *p*)  
place memory barrier behind the `poll_wait` call.

### Parameters

**struct file \* filp** file

**wait\_queue\_head\_t \* wait\_address** socket wait queue

**poll\_table \* p** poll\_table

### Description

See the comments in the `wq_has_sleeper` function.

struct page\_frag \* **sk\_page\_frag**(struct sock \* *sk*)  
return an appropriate `page_frag`

### Parameters

**struct sock \* sk** socket

### Description

If socket allocation mode allows current thread to sleep, it means its safe to use the per task `page_frag` instead of the per socket one.

void **sock\_tx\_timestamp**(const struct sock \* *sk*, \_\_u16 *tsflags*, \_\_u8 \* *tx\_flags*)  
checks whether the outgoing packet is to be time stamped

### Parameters

**const struct sock \* sk** socket sending this packet

**\_\_u16 tsflags** timestamping flags to use

**\_\_u8 \* tx\_flags** completed with instructions for time stamping

### Note

callers should take care of initial `*tx_flags` value (usually 0)

void **sk\_eat\_skb**(struct sock \* *sk*, struct sk\_buff \* *skb*)  
Release a `skb` if it is no longer needed

### Parameters

**struct sock \* sk** socket to eat this `skb` from

**struct sk\_buff \* skb** socket buffer to eat

### Description

This routine must be called with interrupts disabled or with the socket locked so that the `sk_buff` queue operation is ok.

int **sk\_state\_load**(const struct sock \* *sk*)  
read `sk->sk_state` for lockless contexts

### Parameters

**const struct sock \* sk** socket pointer

### Description

Paired with `sk_state_store()`. Used in places we do not hold socket lock : `tcp_diag_get_info()`, `tcp_get_info()`, `tcp_poll()`, `get_tcp4_sock()` ...

void **sk\_state\_store**(struct sock \* *sk*, int *newstate*)  
update `sk->sk_state`

**Parameters**

**struct sock \* sk** socket pointer

**int newstate** new state

**Description**

Paired with [sk\\_state\\_load\(\)](#). Should be used in contexts where state change might impact lockless readers.

**struct socket \* sockfd\_lookup**(int *fd*, int \* *err*)  
Go from a file number to its socket slot

**Parameters**

**int fd** file handle

**int \* err** pointer to an error code return

**Description**

The file handle passed in is locked and the socket it is bound to is returned. If an error occurs the err pointer is overwritten with a negative errno code and NULL is returned. The function checks for both invalid handles and passing a handle which is not a socket.

On a success the socket object pointer is returned.

**struct socket \* sock\_alloc**(void)  
allocate a socket

**Parameters**

**void** no arguments

**Description**

Allocate a new inode and socket object. The two are bound together and initialised. The socket is then returned. If we are out of inodes NULL is returned.

**void sock\_release**(struct [socket](#) \* *sock*)  
close a socket

**Parameters**

**struct socket \* sock** socket to close

**Description**

The socket is released from the protocol stack if it has a release callback, and the inode is then released if the socket is bound to an inode not a file.

**int kernel\_recvmsg**(struct [socket](#) \* *sock*, struct msghdr \* *msg*, struct kvec \* *vec*, size\_t *num*, size\_t *size*, int *flags*)  
Receive a message from a socket (kernel space)

**Parameters**

**struct socket \* sock** The socket to receive the message from

**struct msghdr \* msg** Received message

**struct kvec \* vec** Input s/g array for message data

**size\_t num** Size of input s/g array

**size\_t size** Number of bytes to read

**int flags** Message flags (MSG\_DONTWAIT, etc...)

**Description**

On return the msg structure contains the scatter/gather array passed in the vec argument. The array is modified so that it consists of the unfilled portion of the original array.

The returned value is the total number of bytes received, or an error.

int **sock\_register**(const struct net\_proto\_family \* *ops*)  
add a socket protocol handler

### Parameters

const struct net\_proto\_family \* **ops** description of protocol

### Description

This function is called by a protocol handler that wants to advertise its address family, and have it linked into the socket interface. The value *ops->family* corresponds to the socket system call protocol family.

void **sock\_unregister**(int *family*)  
remove a protocol handler

### Parameters

int **family** protocol family to remove

### Description

This function is called by a protocol handler that wants to remove its address family, and have it unlinked from the new socket creation.

If protocol handler is a module, then it can use module reference counts to protect against new references. If protocol handler is not a module then it needs to provide its own protection in the *ops->create* routine.

struct *sk\_buff* \* **\_\_alloc\_skb**(unsigned int *size*, gfp\_t *gfp\_mask*, int *flags*, int *node*)  
allocate a network buffer

### Parameters

unsigned int **size** size to allocate

gfp\_t **gfp\_mask** allocation mask

int **flags** If SKB\_ALLOC\_FCLONE is set, allocate from fclone cache instead of head cache and allocate a cloned (child) skb. If SKB\_ALLOC\_RX is set, \_\_GFP\_MEMALLOC will be used for allocations in case the data is required for writeback

int **node** numa node to allocate memory on

### Description

Allocate a new *sk\_buff*. The returned buffer has no headroom and a tail room of at least *size* bytes. The object has a reference count of one. The return is the buffer. On a failure the return is NULL.

Buffers may only be allocated from interrupts using a **gfp\_mask** of GFP\_ATOMIC.

void \* **netdev\_alloc\_frag**(unsigned int *fragsz*)  
allocate a page fragment

### Parameters

unsigned int **fragsz** fragment size

### Description

Allocates a frag from a page for receive buffer. Uses GFP\_ATOMIC allocations.

struct *sk\_buff* \* **\_\_netdev\_alloc\_skb**(struct *net\_device* \* *dev*, unsigned int *len*, gfp\_t *gfp\_mask*)  
allocate an skbuff for rx on a specific device

### Parameters

struct *net\_device* \* **dev** network device to receive on

unsigned int **len** length to allocate

**gfp\_t gfp\_mask** get\_free\_pages mask, passed to alloc\_skb

### Description

Allocate a new *sk\_buff* and assign it a usage count of one. The buffer has NET\_SKB\_PAD headroom built in. Users should allocate the headroom they think they need without accounting for the built in space. The built in space is used for optimisations.

NULL is returned if there is no free memory.

```
struct sk_buff * __napi_alloc_skb(struct napi_struct * napi, unsigned int len, gfp_t gfp_mask)
```

allocate skbuff for rx in a specific NAPI instance

### Parameters

**struct napi\_struct \* napi** napi instance this buffer was allocated for

**unsigned int len** length to allocate

**gfp\_t gfp\_mask** get\_free\_pages mask, passed to alloc\_skb and alloc\_pages

### Description

Allocate a new sk\_buff for use in NAPI receive. This buffer will attempt to allocate the head from a special reserved region used only for NAPI Rx allocation. By doing this we can save several CPU cycles by avoiding having to disable and re-enable IRQs.

NULL is returned if there is no free memory.

```
void __kfree_skb(struct sk_buff * skb)
```

private function

### Parameters

**struct sk\_buff \* skb** buffer

### Description

Free an sk\_buff. Release anything attached to the buffer. Clean the state. This is an internal helper function. Users should always call kfree\_skb

```
void kfree_skb(struct sk_buff * skb)
```

free an sk\_buff

### Parameters

**struct sk\_buff \* skb** buffer to free

### Description

Drop a reference to the buffer and free it if the usage count has hit zero.

```
void skb_tx_error(struct sk_buff * skb)
```

report an sk\_buff xmit error

### Parameters

**struct sk\_buff \* skb** buffer that triggered an error

### Description

Report xmit error if a device callback is tracking this skb. skb must be freed afterwards.

```
void consume_skb(struct sk_buff * skb)
```

free an skbuff

### Parameters

**struct sk\_buff \* skb** buffer to free

### Description

Drop a ref to the buffer and free it if the usage count has hit zero Functions identically to `kfree_skb`, but `kfree_skb` assumes that the frame is being dropped after a failure and notes that

`struct sk_buff * skb_morph(struct sk_buff * dst, struct sk_buff * src)`  
morph one skb into another

#### Parameters

`struct sk_buff * dst` the skb to receive the contents

`struct sk_buff * src` the skb to supply the contents

#### Description

This is identical to `skb_clone` except that the target skb is supplied by the user.

The target skb is returned upon exit.

`int skb_copy_ubufs(struct sk_buff * skb, gfp_t gfp_mask)`  
copy userspace skb frags buffers to kernel

#### Parameters

`struct sk_buff * skb` the skb to modify

`gfp_t gfp_mask` allocation priority

#### Description

This must be called on `SKBTX_DEV_ZEROCOPY` skb. It will copy all frags into kernel and drop the reference to userspace pages.

If this function is called from an interrupt `gfp_mask()` must be `GFP_ATOMIC`.

Returns 0 on success or a negative error code on failure to allocate kernel memory to copy to.

`struct sk_buff * skb_clone(struct sk_buff * skb, gfp_t gfp_mask)`  
duplicate an `sk_buff`

#### Parameters

`struct sk_buff * skb` buffer to clone

`gfp_t gfp_mask` allocation priority

#### Description

Duplicate an `sk_buff`. The new one is not owned by a socket. Both copies share the same packet data but not structure. The new buffer has a reference count of 1. If the allocation fails the function returns NULL otherwise the new buffer is returned.

If this function is called from an interrupt `gfp_mask()` must be `GFP_ATOMIC`.

`struct sk_buff * skb_copy(const struct sk_buff * skb, gfp_t gfp_mask)`  
create private copy of an `sk_buff`

#### Parameters

`const struct sk_buff * skb` buffer to copy

`gfp_t gfp_mask` allocation priority

#### Description

Make a copy of both an `sk_buff` and its data. This is used when the caller wishes to modify the data and needs a private copy of the data to alter. Returns NULL on failure or the pointer to the buffer on success. The returned buffer has a reference count of 1.

As by-product this function converts non-linear `sk_buff` to linear one, so that `sk_buff` becomes completely private and caller is allowed to modify all the data of returned buffer. This means that this function is not recommended for use in circumstances when only header is going to be modified. Use `pskb_copy()` instead.

```
struct sk_buff * __pskb_copy_fclone(struct sk_buff * skb, int headroom, gfp_t gfp_mask,  
                                   bool fclone)  
    create copy of an sk_buff with private head.
```

#### Parameters

**struct *sk\_buff* \* *skb*** buffer to copy

**int *headroom*** headroom of new *skb*

**gfp\_t *gfp\_mask*** allocation priority

**bool *fclone*** if true allocate the copy of the *skb* from the *fclone* cache instead of the head cache; it is recommended to set this to true for the cases where the copy will likely be cloned

#### Description

Make a copy of both an *sk\_buff* and part of its data, located in header. Fragmented data remain shared. This is used when the caller wishes to modify only header of *sk\_buff* and needs private copy of the header to alter. Returns NULL on failure or the pointer to the buffer on success. The returned buffer has a reference count of 1.

```
int pskb_expand_head(struct sk_buff * skb, int nhead, int ntail, gfp_t gfp_mask)  
    reallocate header of sk_buff
```

#### Parameters

**struct *sk\_buff* \* *skb*** buffer to reallocate

**int *nhead*** room to add at head

**int *ntail*** room to add at tail

**gfp\_t *gfp\_mask*** allocation priority

#### Description

Expands (or creates identical copy, if ***nhead*** and ***ntail*** are zero) header of ***skb***. *sk\_buff* itself is not changed. *sk\_buff* MUST have reference count of 1. Returns zero in the case of success or error, if expansion failed. In the last case, *sk\_buff* is not changed.

All the pointers pointing into *skb* header may change and must be reloaded after call to this function.

```
struct sk_buff * skb_copy_expand(const struct sk_buff * skb, int newheadroom, int newtailroom,  
                                gfp_t gfp_mask)  
    copy and expand sk_buff
```

#### Parameters

**const struct *sk\_buff* \* *skb*** buffer to copy

**int *newheadroom*** new free bytes at head

**int *newtailroom*** new free bytes at tail

**gfp\_t *gfp\_mask*** allocation priority

#### Description

Make a copy of both an *sk\_buff* and its data and while doing so allocate additional space.

This is used when the caller wishes to modify the data and needs a private copy of the data to alter as well as more space for new fields. Returns NULL on failure or the pointer to the buffer on success. The returned buffer has a reference count of 1.

You must pass GFP\_ATOMIC as the allocation priority if this function is called from an interrupt.

```
int __skb_pad(struct sk_buff * skb, int pad, bool free_on_error)  
    zero pad the tail of an skb
```

#### Parameters

**struct sk\_buff \* skb** buffer to pad  
**int pad** space to pad  
**bool free\_on\_error** free buffer on error

#### Description

Ensure that a buffer is followed by a padding area that is zero filled. Used by network drivers which may DMA or transfer data beyond the buffer end onto the wire.

May return error in out of memory cases. The skb is freed on error if **free\_on\_error** is true.

void \* **pskb\_put**(struct *sk\_buff* \* *skb*, struct *sk\_buff* \* *tail*, int *len*)  
add data to the tail of a potentially fragmented buffer

#### Parameters

**struct sk\_buff \* skb** start of the buffer to use  
**struct sk\_buff \* tail** tail fragment of the buffer to use  
**int len** amount of data to add

#### Description

This function extends the used data area of the potentially fragmented buffer. **tail** must be the last fragment of **skb** – or **skb** itself. If this would exceed the total buffer size the kernel will panic. A pointer to the first byte of the extra data is returned.

void \* **skb\_put**(struct *sk\_buff* \* *skb*, unsigned int *len*)  
add data to a buffer

#### Parameters

**struct sk\_buff \* skb** buffer to use  
**unsigned int len** amount of data to add

#### Description

This function extends the used data area of the buffer. If this would exceed the total buffer size the kernel will panic. A pointer to the first byte of the extra data is returned.

void \* **skb\_push**(struct *sk\_buff* \* *skb*, unsigned int *len*)  
add data to the start of a buffer

#### Parameters

**struct sk\_buff \* skb** buffer to use  
**unsigned int len** amount of data to add

#### Description

This function extends the used data area of the buffer at the buffer start. If this would exceed the total buffer headroom the kernel will panic. A pointer to the first byte of the extra data is returned.

void \* **skb\_pull**(struct *sk\_buff* \* *skb*, unsigned int *len*)  
remove data from the start of a buffer

#### Parameters

**struct sk\_buff \* skb** buffer to use  
**unsigned int len** amount of data to remove

#### Description

This function removes data from the start of a buffer, returning the memory to the headroom. A pointer to the next data in the buffer is returned. Once the data has been pulled future pushes will overwrite the old data.



void **skb\_trim**(struct *sk\_buff* \* *skb*, unsigned int *len*)  
remove end from a buffer

#### Parameters

**struct sk\_buff \* skb** buffer to alter

**unsigned int len** new length

#### Description

Cut the length of a buffer down by removing data from the tail. If the buffer is already under the length specified it is not modified. The skb must be linear.

void \* **\_\_pskb\_pull\_tail**(struct *sk\_buff* \* *skb*, int *delta*)  
advance tail of skb header

#### Parameters

**struct sk\_buff \* skb** buffer to reallocate

**int delta** number of bytes to advance tail

#### Description

The function makes a sense only on a fragmented *sk\_buff*, it expands header moving its tail forward and copying necessary data from fragmented part.

*sk\_buff* MUST have reference count of 1.

Returns NULL (and *sk\_buff* does not change) if pull failed or value of new tail of skb in the case of success.

All the pointers pointing into skb header may change and must be reloaded after call to this function.

int **skb\_copy\_bits**(const struct *sk\_buff* \* *skb*, int *offset*, void \* *to*, int *len*)  
copy bits from skb to kernel buffer

#### Parameters

**const struct sk\_buff \* skb** source skb

**int offset** offset in source

**void \* to** destination buffer

**int len** number of bytes to copy

#### Description

Copy the specified number of bytes from the source skb to the destination buffer.

**CAUTION ! :** If its prototype is ever changed, check arch/{\*}/net/{\*}.S files, since it is called from BPF assembly code.

int **skb\_store\_bits**(struct *sk\_buff* \* *skb*, int *offset*, const void \* *from*, int *len*)  
store bits from kernel buffer to skb

#### Parameters

**struct sk\_buff \* skb** destination buffer

**int offset** offset in destination

**const void \* from** source buffer

**int len** number of bytes to copy

#### Description

Copy the specified number of bytes from the source buffer to the destination skb. This function handles all the messy bits of traversing fragment lists and such.

int **skb\_zerocopy**(struct *sk\_buff* \* *to*, struct *sk\_buff* \* *from*, int *len*, int *hlen*)  
Zero copy skb to skb

#### Parameters

**struct sk\_buff \* to** destination buffer  
**struct sk\_buff \* from** source buffer  
**int len** number of bytes to copy from source buffer  
**int hlen** size of linear headroom in destination buffer

#### Description

Copies up to *len* bytes from *from* to *to* by creating references to the frags in the source buffer. The *hlen* as calculated by `skb_zerocopy_headlen()` specifies the headroom in the *to* buffer. Return value: 0: everything is OK -ENOMEM: couldn't orphan frags of **from** due to lack of memory -EFAULT: `skb_copy_bits()` found some problem with skb geometry

struct *sk\_buff* \* **skb\_dequeue**(struct sk\_buff\_head \* *list*)  
remove from the head of the queue

#### Parameters

**struct sk\_buff\_head \* list** list to dequeue from

#### Description

Remove the head of the list. The list lock is taken so the function may be used safely with other locking list functions. The head item is returned or NULL if the list is empty.

struct *sk\_buff* \* **skb\_dequeue\_tail**(struct sk\_buff\_head \* *list*)  
remove from the tail of the queue

#### Parameters

**struct sk\_buff\_head \* list** list to dequeue from

#### Description

Remove the tail of the list. The list lock is taken so the function may be used safely with other locking list functions. The tail item is returned or NULL if the list is empty.

void **skb\_queue\_purge**(struct sk\_buff\_head \* *list*)  
empty a list

#### Parameters

**struct sk\_buff\_head \* list** list to empty

#### Description

Delete all buffers on an *sk\_buff* list. Each buffer is removed from the list and one reference dropped. This function takes the list lock and is atomic with respect to other list locking functions.

void **skb\_queue\_head**(struct sk\_buff\_head \* *list*, struct *sk\_buff* \* *newsk*)  
queue a buffer at the list head

#### Parameters

**struct sk\_buff\_head \* list** list to use  
**struct sk\_buff \* newsk** buffer to queue

#### Description

Queue a buffer at the start of the list. This function takes the list lock and can be used safely with other locking *sk\_buff* functions safely.

A buffer cannot be placed on two lists at the same time.

void **skb\_queue\_tail**(struct sk\_buff\_head \* *list*, struct *sk\_buff* \* *newsk*)  
queue a buffer at the list tail

#### Parameters

struct sk\_buff\_head \* **list** list to use

struct sk\_buff \* **newsk** buffer to queue

#### Description

Queue a buffer at the tail of the list. This function takes the list lock and can be used safely with other locking *sk\_buff* functions safely.

A buffer cannot be placed on two lists at the same time.

void **skb\_unlink**(struct *sk\_buff* \* *skb*, struct sk\_buff\_head \* *list*)  
remove a buffer from a list

#### Parameters

struct sk\_buff \* **skb** buffer to remove

struct sk\_buff\_head \* **list** list to use

#### Description

Remove a packet from a list. The list locks are taken and this function is atomic with respect to other list locked calls

You must know what list the SKB is on.

void **skb\_append**(struct *sk\_buff* \* *old*, struct *sk\_buff* \* *newsk*, struct sk\_buff\_head \* *list*)  
append a buffer

#### Parameters

struct sk\_buff \* **old** buffer to insert after

struct sk\_buff \* **newsk** buffer to insert

struct sk\_buff\_head \* **list** list to use

#### Description

Place a packet after a given packet in a list. The list locks are taken and this function is atomic with respect to other list locked calls. A buffer cannot be placed on two lists at the same time.

void **skb\_insert**(struct *sk\_buff* \* *old*, struct *sk\_buff* \* *newsk*, struct sk\_buff\_head \* *list*)  
insert a buffer

#### Parameters

struct sk\_buff \* **old** buffer to insert before

struct sk\_buff \* **newsk** buffer to insert

struct sk\_buff\_head \* **list** list to use

#### Description

Place a packet before a given packet in a list. The list locks are taken and this function is atomic with respect to other list locked calls.

A buffer cannot be placed on two lists at the same time.

void **skb\_split**(struct *sk\_buff* \* *skb*, struct *sk\_buff* \* *skb1*, const u32 *len*)  
Split fragmented skb to two parts at length len.

#### Parameters

struct sk\_buff \* **skb** the buffer to split

struct sk\_buff \* **skb1** the buffer to receive the second part

```
const u32 len new length for skb
```

```
void skb_prepare_seq_read(struct sk_buff *skb, unsigned int from, unsigned int to, struct
                        skb_seq_state *st)
    Prepare a sequential read of skb data
```

## Parameters

```
struct sk_buff * skb the buffer to read
unsigned int from lower offset of data to be read
unsigned int to upper offset of data to be read
struct skb_seq_state * st state variable
```

### Description

Initializes the specified state variable. Must be called before invoking `skb_seq_read()` for the first time.

```
unsigned int skb_seq_read(unsigned int consumed, const u8 ** data, struct skb_seq_state * st)
```

Sequentially read skb data

## Parameters

```

unsigned int consumed number of bytes consumed by the caller so far
const u8 ** data destination pointer for data to be returned
struct skb_seq_state * st state variable

```

### Description

Reads a block of skb data at **consumed** relative to the lower offset specified to `skb_prepare_seq_read()`. Assigns the head of the data block to **data** and returns the length of the block or 0 if the end of the skb data or the upper offset has been reached.

The caller is not required to consume all of the data returned, i.e. **consumed** is typically set to the number of bytes already consumed and the next call to `skb_seq_read()` will return the remaining part of the block.

**Note 1: The size of each block of data returned can be arbitrary,** this limitation is the cost for zero-copy sequential reads of potentially non linear data.

**Note 2: Fragment lists within fragments are not implemented** at the moment, `state->root_skb` could be replaced with a stack for this purpose.

```
void skb_abort_seq_read(struct skb_seq_state * st)
    Abort a sequential read of skb data
```

## Parameters

```
struct skb_seq_state * st state variable
```

### Description

Must be called if `skb_seq_read()` was not called until it returned 0.

```
unsigned int skb_find_text(struct sk_buff * skb, unsigned int from, unsigned int to, struct ts_config
                        * config)
    Find a text pattern in skb data
```

## Parameters

```
struct sk_buff * skb the buffer to look in
unsigned int from search offset
unsigned int to search limit
struct ts_config * config textsearch configuration
```

## Description

Finds a pattern in the skb data according to the specified textsearch configuration. Use `textsearch_next()` to retrieve subsequent occurrences of the pattern. Returns the offset to the first occurrence or `UINT_MAX` if no match was found.

```
int skb_append_datato_frags(struct sock *sk, struct sk_buff *skb, int (*getfrag) (void *from,
                                                                    char *to, int offset, int len, int odd, struct sk_buff *skb, void * from,
                                                                    int length)
    append the user data to a skb
```

## Parameters

**struct sock \* sk** sock structure

**struct sk\_buff \* skb** skb structure to be appended with user data.

**int (\*)(void \*from, char \*to, int offset, int len, int odd, struct sk\_buff \*skb) getfrag** call back function to be used for getting the user data

**void \* from** pointer to user message iov

**int length** length of the iov message

## Description

This procedure append the user data in the fragment part of the skb if any page alloc fails user this procedure returns `-ENOMEM`

```
void * skb_pull_rcsum(struct sk_buff *skb, unsigned int len)
    pull skb and update receive checksum
```

## Parameters

**struct sk\_buff \* skb** buffer to update

**unsigned int len** length of data pulled

## Description

This function performs an `skb_pull` on the packet and updates the `CHECKSUM_COMPLETE` checksum. It should be used on receive path processing instead of `skb_pull` unless you know that the checksum difference is zero (e.g., a valid IP header) or you are setting `ip_summed` to `CHECKSUM_NONE`.

```
struct sk_buff * skb_segment(struct sk_buff * head_skb, netdev_features_t features)
    Perform protocol segmentation on skb.
```

## Parameters

**struct sk\_buff \* head\_skb** buffer to segment

**netdev\_features\_t features** features for the output path (see `dev->features`)

## Description

This function performs segmentation on the given skb. It returns a pointer to the first in a list of new skbs for the segments. In case of error it returns `ERR_PTR(err)`.

```
int skb_to_sgvec(struct sk_buff *skb, struct scatterlist *sg, int offset, int len)
    Fill a scatter-gather list from a socket buffer
```

## Parameters

**struct sk\_buff \* skb** Socket buffer containing the buffers to be mapped

**struct scatterlist \* sg** The scatter-gather list to map into

**int offset** The offset into the buffer's contents to start mapping

**int len** Length of buffer space to be mapped

## Description

Fill the specified scatter-gather list with mappings/pointers into a region of the buffer space attached to a socket buffer. Returns either the number of scatterlist items used, or -EMSGSIZE if the contents could not fit.

int **skb\_cow\_data**(struct *sk\_buff* \* *skb*, int *tailbits*, struct *sk\_buff* \*\* *trailer*)  
Check that a socket buffer's data buffers are writable

### Parameters

**struct sk\_buff \* skb** The socket buffer to check.  
**int tailbits** Amount of trailing space to be added  
**struct sk\_buff \*\* trailer** Returned pointer to the skb where the **tailbits** space begins

### Description

Make sure that the data buffers attached to a socket buffer are writable. If they are not, private copies are made of the data buffers and the socket buffer is set to use these instead.

If **tailbits** is given, make sure that there is space to write **tailbits** bytes of data beyond current end of socket buffer. **trailer** will be set to point to the skb in which this space begins.

The number of scatterlist elements required to completely map the COW'd and extended socket buffer will be returned.

struct *sk\_buff* \* **skb\_clone\_skb**(struct *sk\_buff* \* *skb*)  
create clone of skb, and take reference to socket

### Parameters

**struct sk\_buff \* skb** the skb to clone

### Description

This function creates a clone of a buffer that holds a reference on *sk\_refcnt*. Buffers created via this function are meant to be returned using *sock\_queue\_err\_skb*, or free via *kfree\_skb*.

When passing buffers allocated with this function to *sock\_queue\_err\_skb* it is necessary to wrap the call with *sock\_hold*/*sock\_put* in order to prevent the socket from being released prior to being enqueued on the *sk\_error\_queue*.

bool **skb\_partial\_csum\_set**(struct *sk\_buff* \* *skb*, u16 *start*, u16 *off*)  
set up and verify partial csum values for packet

### Parameters

**struct sk\_buff \* skb** the skb to set  
**u16 start** the number of bytes after *skb->data* to start checksumming.  
**u16 off** the offset from start to place the checksum.

### Description

For untrusted partially-checksummed packets, we need to make sure the values for *skb->csum\_start* and *skb->csum\_offset* are valid so we don't oops.

This function checks and sets those values and *skb->ip\_summed*: if this returns false you should drop the packet.

int **skb\_checksum\_setup**(struct *sk\_buff* \* *skb*, bool *recalculate*)  
set up partial checksum offset

### Parameters

**struct sk\_buff \* skb** the skb to set up  
**bool recalculate** if true the pseudo-header checksum will be recalculated

```
struct sk_buff * skb_checksum_trimmed(struct sk_buff * skb, unsigned int transport_len,  
                                       __sum16(*skb_chkf) (struct sk_buff *skb)  
                                       validate checksum of an skb
```

### Parameters

**struct *sk\_buff* \* *skb*** the skb to check

**unsigned int *transport\_len*** the data length beyond the network header

**\_\_sum16(\*)(*struct sk\_buff* \**skb*) *skb\_chkf*** checksum function to use

### Description

Applies the given checksum function *skb\_chkf* to the provided *skb*. Returns a checked and maybe trimmed *skb*. Returns NULL on error.

If the *skb* has data beyond the given transport length, then a trimmed & cloned *skb* is checked and returned.

Caller needs to set the *skb* transport header and free any returned *skb* if it differs from the provided *skb*.

```
bool skb_try_coalesce(struct sk_buff * to, struct sk_buff * from, bool * fragstolen, int  
                      * delta_truesize)  
    try to merge skb to prior one
```

### Parameters

**struct *sk\_buff* \* *to*** prior buffer

**struct *sk\_buff* \* *from*** buffer to add

**bool \* *fragstolen*** pointer to boolean

**int \* *delta\_truesize*** how much more was allocated than was requested

```
void skb_scrub_packet(struct sk_buff * skb, bool xnet)  
    scrub an skb
```

### Parameters

**struct *sk\_buff* \* *skb*** buffer to clean

**bool *xnet*** packet is crossing netns

### Description

*skb\_scrub\_packet* can be used after encapsulating or decapsulating a packet into/from a tunnel. Some information have to be cleared during these operations. *skb\_scrub\_packet* can also be used to clean a *skb* before injecting it in another namespace (***xnet*** == true). We have to clear all information in the *skb* that could impact namespace isolation.

```
unsigned int skb_gso_transport_seglen(const struct sk_buff * skb)  
    Return length of individual segments of a gso packet
```

### Parameters

**const struct *sk\_buff* \* *skb*** GSO *skb*

### Description

*skb\_gso\_transport\_seglen* is used to determine the real size of the individual segments, including Layer4 headers (TCP/UDP).

The MAC/L2 or network (IP, IPv6) headers are not accounted for.

```
bool skb_gso_validate_mtu(const struct sk_buff * skb, unsigned int mtu)  
    Return in case such skb fits a given MTU
```

### Parameters

**const struct *sk\_buff* \* *skb*** GSO *skb*

**unsigned int mtu** MTU to validate against

### Description

`skb_gso_validate_mtu` validates if a given `skb` will fit a wanted MTU once split.

`struct sk_buff * alloc_skb_with_frags`(unsigned long *header\_len*, unsigned long *data\_len*,  
int *max\_page\_order*, int \* *errcode*, gfp\_t *gfp\_mask*)  
allocate `skb` with page frags

### Parameters

**unsigned long header\_len** size of linear part

**unsigned long data\_len** needed length in frags

**int max\_page\_order** max page order desired.

**int \* errcode** pointer to error code if any

**gfp\_t gfp\_mask** allocation mask

### Description

This can be used to allocate a paged `skb`, given a maximal order for frags.

`bool sk_ns_capable`(const struct *sock* \* *sk*, struct user\_namespace \* *user\_ns*, int *cap*)  
General socket capability test

### Parameters

**const struct sock \* sk** Socket to use a capability on or through

**struct user\_namespace \* user\_ns** The user namespace of the capability to use

**int cap** The capability to use

### Description

Test to see if the opener of the socket had when the socket was created and the current process has the capability **cap** in the user namespace **user\_ns**.

`bool sk_capable`(const struct *sock* \* *sk*, int *cap*)  
Socket global capability test

### Parameters

**const struct sock \* sk** Socket to use a capability on or through

**int cap** The global capability to use

### Description

Test to see if the opener of the socket had when the socket was created and the current process has the capability **cap** in all user namespaces.

`bool sk_net_capable`(const struct *sock* \* *sk*, int *cap*)  
Network namespace socket capability test

### Parameters

**const struct sock \* sk** Socket to use a capability on or through

**int cap** The capability to use

### Description

Test to see if the opener of the socket had when the socket was created and the current process has the capability **cap** over the network namespace the socket is a member of.

`void sk_set_memalloc`(struct *sock* \* *sk*)  
sets `SOCK_MEMALLOC`

### Parameters



**struct sock \* sk** socket to set it on

### Description

Set SOCK\_MEMALLOC on a socket for access to emergency reserves. It's the responsibility of the admin to adjust min\_free\_kbytes to meet the requirements

struct [sock](#) \* **sk\_alloc**(struct net \* *net*, int *family*, gfp\_t *priority*, struct proto \* *prot*, int *kern*)  
All socket objects are allocated here

### Parameters

**struct net \* net** the applicable net namespace

**int family** protocol family

**gfp\_t priority** for allocation (GFP\_KERNEL, GFP\_ATOMIC, etc)

**struct proto \* prot** struct proto associated with this new sock instance

**int kern** is this to be a kernel socket?

struct [sock](#) \* **sk\_clone\_lock**(const struct [sock](#) \* *sk*, const gfp\_t *priority*)  
clone a socket, and lock its clone

### Parameters

**const struct sock \* sk** the socket to clone

**const gfp\_t priority** for allocation (GFP\_KERNEL, GFP\_ATOMIC, etc)

### Description

Caller must unlock socket even in error path (bh\_unlock\_sock(newsk))

bool **skb\_page\_frag\_refill**(unsigned int *sz*, struct page\_frag \* *pfrag*, gfp\_t *gfp*)  
check that a page\_frag contains enough room

### Parameters

**unsigned int sz** minimum size of the fragment we want to get

**struct page\_frag \* pfrag** pointer to page\_frag

**gfp\_t gfp** priority for memory allocation

### Note

While this allocator tries to use high order pages, there is no guarantee that allocations succeed. Therefore, **sz** MUST be less or equal than PAGE\_SIZE.

int **sk\_wait\_data**(struct [sock](#) \* *sk*, long \* *timeo*, const struct [sk\\_buff](#) \* *skb*)  
wait for data to arrive at sk\_receive\_queue

### Parameters

**struct sock \* sk** sock to wait on

**long \* timeo** for how long

**const struct sk\_buff \* skb** last skb seen on sk\_receive\_queue

### Description

Now socket state including sk->sk\_err is changed only under lock, hence we may omit checks after joining wait queue. We check receive queue before schedule() only as optimization; it is very likely that release\_sock() added new data.

int **\_\_sk\_mem\_raise\_allocated**(struct [sock](#) \* *sk*, int *size*, int *amt*, int *kind*)  
increase memory\_allocated

### Parameters

**struct sock \* sk** socket

**int size** memory size to allocate

**int amt** pages to allocate

**int kind** allocation type

### Description

Similar to `__sk_mem_schedule()`, but does not update `sk_forward_alloc`

**int** `__sk_mem_schedule`(struct *sock* \* *sk*, int *size*, int *kind*)  
increase `sk_forward_alloc` and `memory_allocated`

### Parameters

**struct sock** \* *sk* socket

**int size** memory size to allocate

**int kind** allocation type

### Description

If `kind` is `SK_MEM_SEND`, it means `wmem` allocation. Otherwise it means `rmem` allocation. This function assumes that protocols which have `memory_pressure` use `sk_wmem_queued` as write buffer accounting.

**void** `__sk_mem_reduce_allocated`(struct *sock* \* *sk*, int *amount*)  
reclaim `memory_allocated`

### Parameters

**struct sock** \* *sk* socket

**int amount** number of quanta

### Description

Similar to `__sk_mem_reclaim()`, but does not update `sk_forward_alloc`

**void** `__sk_mem_reclaim`(struct *sock* \* *sk*, int *amount*)  
reclaim `sk_forward_alloc` and `memory_allocated`

### Parameters

**struct sock** \* *sk* socket

**int amount** number of bytes (rounded down to a `SK_MEM_QUANTUM` multiple)

**bool** `lock_sock_fast`(struct *sock* \* *sk*)  
fast version of `lock_sock`

### Parameters

**struct sock** \* *sk* socket

### Description

This version should be used for very small section, where process wont block return false if fast path is taken:

`sk_lock.slock` locked, owned = 0, BH disabled

return true if slow path is taken:

`sk_lock.slock` unlocked, owned = 1, BH enabled

**struct sk\_buff** \* `__skb_try_recv_datagram`(struct *sock* \* *sk*, unsigned int *flags*, void (\**destructor*)  
(struct *sock* \**sk*, struct *sk\_buff* \**skb*, int \**peeked*, int  
\* *off*, int \* *err*, struct *sk\_buff* \*\* *last*)

Receive a datagram skbuff

### Parameters

```
struct sock * sk socket
unsigned int flags MSG_flags
void (*)(struct sock *sk, struct sk_buff *skb) destructor invoked under the receive lock on
    successful dequeue
int * peeked returns non-zero if this packet has been seen before
int * off an offset in bytes to peek skb from. Returns an offset within an skb where data actually starts
int * err error code returned
struct sk_buff ** last set to last peeked message to inform the wait function what to look for when
    peeking
```

### Description

Get a datagram skbuff, understands the peeking, nonblocking wakeups and possible races. This replaces identical code in packet, raw and udp, as well as the IPX AX.25 and Appletalk. It also finally fixes the long standing peek and read race for datagram sockets. If you alter this routine remember it must be re-entrant.

This function will lock the socket if a skb is returned, so the caller needs to unlock the socket in that case (usually by calling `skb_free_datagram`). Returns NULL with **err** set to -EAGAIN if no data was available or to some other value if an error was detected.

- It does not lock socket since today. This function is
- free of race conditions. This measure should/can improve
- significantly datagram socket latencies at high loads,
- when data copying to user space takes lots of time.
- (BTW I've just killed the last `cli()` in IP/IPv6/core/netlink/packet
- 8. Great win.)
- -ANK (980729)

The order of the tests when we find no data waiting are specified quite explicitly by POSIX 1003.1g, don't change them without having the standard around please.

```
int skb_kill_datagram(struct sock * sk, struct sk_buff * skb, unsigned int flags)
    Free a datagram skbuff forcibly
```

### Parameters

```
struct sock * sk socket
struct sk_buff * skb datagram skbuff
unsigned int flags MSG_flags
```

### Description

This function frees a datagram skbuff that was received by `skb_recv_datagram`. The flags argument must match the one used for `skb_recv_datagram`.

If the MSG\_PEEK flag is set, and the packet is still on the receive queue of the socket, it will be taken off the queue before it is freed.

This function currently only disables BH when acquiring the `sk_receive_queue` lock. Therefore it must not be used in a context where that lock is acquired in an IRQ context.

It returns 0 if the packet was removed by us.

```
int skb_copy_datagram_iter(const struct sk_buff * skb, int offset, struct iov_iter * to, int len)
    Copy a datagram to an iovec iterator.
```

### Parameters

**const struct sk\_buff \* skb** buffer to copy

**int offset** offset in the buffer to start copying from

**struct iov\_iter \* to** iovec iterator to copy to

**int len** amount of data to copy from buffer to iovec

int **skb\_copy\_datagram\_from\_iter**(struct *sk\_buff* \* *skb*, int *offset*, struct iov\_iter \* *from*, int *len*)  
Copy a datagram from an iov\_iter.

#### Parameters

**struct sk\_buff \* skb** buffer to copy

**int offset** offset in the buffer to start copying to

**struct iov\_iter \* from** the copy source

**int len** amount of data to copy to buffer from iovec

#### Description

Returns 0 or -EFAULT.

int **zerocopy\_sg\_from\_iter**(struct *sk\_buff* \* *skb*, struct iov\_iter \* *from*)  
Build a zerocopy datagram from an iov\_iter

#### Parameters

**struct sk\_buff \* skb** buffer to copy

**struct iov\_iter \* from** the source to copy from

#### Description

The function will first copy up to headlen, and then pin the userspace pages and build frags through them.

Returns 0, -EFAULT or -EMSGSIZE.

int **skb\_copy\_and\_csum\_datagram\_msg**(struct *sk\_buff* \* *skb*, int *hlen*, struct msghdr \* *msg*)  
Copy and checksum skb to user iovec.

#### Parameters

**struct sk\_buff \* skb** skbuff

**int hlen** hardware length

**struct msghdr \* msg** destination

#### Description

Caller \_must\_ check that skb will fit to this iovec.

#### Return

**0 - success.** -EINVAL - checksum failure. -EFAULT - fault during copy.

unsigned int **datagram\_poll**(struct file \* *file*, struct *socket* \* *sock*, poll\_table \* *wait*)  
generic datagram poll

#### Parameters

**struct file \* file** file struct

**struct socket \* sock** socket

**poll\_table \* wait** poll table

#### Description

Datagram poll: Again totally generic. This also handles sequenced packet sockets providing the socket receive queue is only ever holding data ready to receive.

**Note**

when you *don't* use this routine for this protocol, and you use a different write policy from `sock_writable()` then please supply your own `write_space` callback.

int **sk\_stream\_wait\_connect**(struct *sock* \* *sk*, long \* *timeo\_p*)  
Wait for a socket to get into the connected state

**Parameters**

struct *sock* \* *sk* socket to wait on

long \* *timeo\_p* for how long to wait

**Description**

Must be called with the socket locked.

int **sk\_stream\_wait\_memory**(struct *sock* \* *sk*, long \* *timeo\_p*)  
Wait for more memory for a socket

**Parameters**

struct *sock* \* *sk* socket to wait for memory

long \* *timeo\_p* for how long

**Socket Filter**

int **sk\_filter\_trim\_cap**(struct *sock* \* *sk*, struct *sk\_buff* \* *skb*, unsigned int *cap*)  
run a packet through a socket filter

**Parameters**

struct *sock* \* *sk* socket associated with *sk\_buff*

struct *sk\_buff* \* *skb* buffer to filter

unsigned int *cap* limit on how short the eBPF program may trim the packet

**Description**

Run the eBPF program and then cut *skb->data* to correct size returned by the program. If *pkt\_len* is 0 we toss packet. If *skb->len* is smaller than *pkt\_len* we keep whole *skb->data*. This is the socket level wrapper to `BPF_PROG_RUN`. It returns 0 if the packet should be accepted or `-EPERM` if the packet should be tossed.

int **bpf\_prog\_create**(struct *bpf\_prog* \*\* *pfp*, struct *sock\_fprog\_kern* \* *fprog*)  
create an unattached filter

**Parameters**

struct *bpf\_prog* \*\* *pfp* the unattached filter that is created

struct *sock\_fprog\_kern* \* *fprog* the filter program

**Description**

Create a filter independent of any socket. We first run some sanity checks on it to make sure it does not explode on us later. If an error occurs or there is insufficient memory for the filter a negative `errno` code is returned. On success the return is zero.

int **bpf\_prog\_create\_from\_user**(struct *bpf\_prog* \*\* *pfp*, struct *sock\_fprog* \* *fprog*,  
bpf\_aux\_classic\_check\_t *trans*, bool *save\_orig*)  
create an unattached filter from user buffer

**Parameters**

struct *bpf\_prog* \*\* *pfp* the unattached filter that is created

struct *sock\_fprog* \* *fprog* the filter program

**bpf\_aux\_classic\_check\_t trans** post-classic verifier transformation handler

**bool save\_orig** save classic BPF program

### Description

This function effectively does the same as [bpf\\_prog\\_create\(\)](#), only that it builds up its insns buffer from user space provided buffer. It also allows for passing a `bpf_aux_classic_check_t` handler.

int **sk\_attach\_filter**(struct sock\_fprog \* *fprog*, struct [sock](#) \* *sk*)  
attach a socket filter

### Parameters

**struct sock\_fprog \* fprog** the filter program

**struct sock \* sk** the socket to use

### Description

Attach the user's filter code. We first run some sanity checks on it to make sure it does not explode on us later. If an error occurs or there is insufficient memory for the filter a negative errno code is returned. On success the return is zero.

## Generic Network Statistics

struct **gnet\_stats\_basic**  
byte/packet throughput statistics

### Definition

```
struct gnet_stats_basic {
    __u64 bytes;
    __u32 packets;
};
```

### Members

**bytes** number of seen bytes

**packets** number of seen packets

struct **gnet\_stats\_rate\_est**  
rate estimator

### Definition

```
struct gnet_stats_rate_est {
    __u32 bps;
    __u32 pps;
};
```

### Members

**bps** current byte rate

**pps** current packet rate

struct **gnet\_stats\_rate\_est64**  
rate estimator

### Definition

```
struct gnet_stats_rate_est64 {
    __u64 bps;
    __u64 pps;
};
```

**Members****bps** current byte rate**pps** current packet ratestruct **gnet\_stats\_queue**  
queuing statistics**Definition**

```

struct gnet_stats_queue {
    __u32 qlen;
    __u32 backlog;
    __u32 drops;
    __u32 requeues;
    __u32 overlimits;
};

```

**Members****qlen** queue length**backlog** backlog size of queue**drops** number of dropped packets**requeues** number of requeues**overlimits** number of enqueues over the limitstruct **gnet\_estimator**  
rate estimator configuration**Definition**

```

struct gnet_estimator {
    signed char interval;
    unsigned char ewma_log;
};

```

**Members****interval** sampling period**ewma\_log** the log of measurement window weight

int **gnet\_stats\_start\_copy\_compat**(struct *sk\_buff* \**skb*, int *type*, int *tc\_stats\_type*,  
int *xstats\_type*, spinlock\_t \**lock*, struct *gnet\_dump* \**d*,  
int *padattr*)  
start dumping procedure in compatibility mode

**Parameters****struct sk\_buff \* skb** socket buffer to put statistics TLVs into**int type** TLV type for top level statistic TLV**int tc\_stats\_type** TLV type for backward compatibility struct tc\_stats TLV**int xstats\_type** TLV type for backward compatibility xstats TLV**spinlock\_t \* lock** statistics lock**struct gnet\_dump \* d** dumping handle**int padattr** padding attribute**Description**

Initializes the dumping handle, grabs the statistic lock and appends an empty TLV header to the socket buffer for use a container for all other statistic TLVS.

The dumping handle is marked to be in backward compatibility mode telling all `gnet_stats_copy_XXX()` functions to fill a local copy of `struct tc_stats`.

Returns 0 on success or -1 if the room in the socket buffer was not sufficient.

int **gnet\_stats\_start\_copy**(struct *sk\_buff* \* *skb*, int *type*, spinlock\_t \* *lock*, struct gnet\_dump \* *d*,  
int *padattr*)  
start dumping procedure in compatibility mode

### Parameters

**struct sk\_buff \* skb** socket buffer to put statistics TLVs into

**int type** TLV type for top level statistic TLV

**spinlock\_t \* lock** statistics lock

**struct gnet\_dump \* d** dumping handle

**int padattr** padding attribute

### Description

Initializes the dumping handle, grabs the statistic lock and appends an empty TLV header to the socket buffer for use as a container for all other statistic TLVs.

Returns 0 on success or -1 if the room in the socket buffer was not sufficient.

int **gnet\_stats\_copy\_basic**(const seqcount\_t \* *running*, struct gnet\_dump \* *d*, struct  
gnet\_stats\_basic\_cpu \_\_percpu \* *cpu*, struct gnet\_stats\_basic\_packed  
\* *b*)  
copy basic statistics into statistic TLV

### Parameters

**const seqcount\_t \* running** seqcount\_t pointer

**struct gnet\_dump \* d** dumping handle

**struct gnet\_stats\_basic\_cpu \_\_percpu \* cpu** copy statistic per cpu

**struct gnet\_stats\_basic\_packed \* b** basic statistics

### Description

Appends the basic statistics to the top level TLV created by `gnet_stats_start_copy()`.

Returns 0 on success or -1 with the statistic lock released if the room in the socket buffer was not sufficient.

int **gnet\_stats\_copy\_rate\_est**(struct gnet\_dump \* *d*, struct net\_rate\_estimator \_\_rcu \*\* *rate\_est*)  
copy rate estimator statistics into statistics TLV

### Parameters

**struct gnet\_dump \* d** dumping handle

**struct net\_rate\_estimator \_\_rcu \*\* rate\_est** rate estimator

### Description

Appends the rate estimator statistics to the top level TLV created by `gnet_stats_start_copy()`.

Returns 0 on success or -1 with the statistic lock released if the room in the socket buffer was not sufficient.

int **gnet\_stats\_copy\_queue**(struct gnet\_dump \* *d*, struct *gnet\_stats\_queue* \_\_percpu \* *cpu\_q*, struct  
*gnet\_stats\_queue* \* *q*, \_\_u32 *qlen*)  
copy queue statistics into statistics TLV

### Parameters

**struct gnet\_dump \* d** dumping handle



**struct gnet\_stats\_queue \_\_percpu \* cpu\_q** per cpu queue statistics

**struct gnet\_stats\_queue \* q** queue statistics

**\_\_u32 qlen** queue length statistics

### Description

Appends the queue statistics to the top level TLV created by [gnet\\_stats\\_start\\_copy\(\)](#). Using per cpu queue statistics if they are available.

Returns 0 on success or -1 with the statistic lock released if the room in the socket buffer was not sufficient.

**int gnet\_stats\_copy\_app**(struct gnet\_dump \* *d*, void \* *st*, int *len*)  
copy application specific statistics into statistics TLV

### Parameters

**struct gnet\_dump \* d** dumping handle

**void \* st** application specific statistics data

**int len** length of data

### Description

Appends the application specific statistics to the top level TLV created by [gnet\\_stats\\_start\\_copy\(\)](#) and remembers the data for XSTATS if the dumping handle is in backward compatibility mode.

Returns 0 on success or -1 with the statistic lock released if the room in the socket buffer was not sufficient.

**int gnet\_stats\_finish\_copy**(struct gnet\_dump \* *d*)  
finish dumping procedure

### Parameters

**struct gnet\_dump \* d** dumping handle

### Description

Corrects the length of the top level TLV to include all TLVs added by [gnet\\_stats\\_copy\\_XXX\(\)](#) calls. Adds the backward compatibility TLVs if [gnet\\_stats\\_start\\_copy\\_compat\(\)](#) was used and releases the statistics lock.

Returns 0 on success or -1 with the statistic lock released if the room in the socket buffer was not sufficient.

**int gen\_new\_estimator**(struct gnet\_stats\_basic\_packed \* *bstats*, struct gnet\_stats\_basic\_cpu \_\_percpu \* *cpu\_bstats*, struct net\_rate\_estimator \_\_rcu \*\* *rate\_est*, spinlock\_t \* *stats\_lock*, seqcount\_t \* *running*, struct nlattr \* *opt*)  
create a new rate estimator

### Parameters

**struct gnet\_stats\_basic\_packed \* bstats** basic statistics

**struct gnet\_stats\_basic\_cpu \_\_percpu \* cpu\_bstats** bstats per cpu

**struct net\_rate\_estimator \_\_rcu \*\* rate\_est** rate estimator statistics

**spinlock\_t \* stats\_lock** statistics lock

**seqcount\_t \* running** qdisc running seqcount

**struct nlattr \* opt** rate estimator configuration TLV

### Description

Creates a new rate estimator with *bstats* as source and *rate\_est* as destination. A new timer with the interval specified in the configuration TLV is created. Upon each interval, the latest statistics will be read

from `bstats` and the estimated rate will be stored in `rate_est` with the statistics lock grabbed during this period.

Returns 0 on success or a negative error code.

void **gen\_kill\_estimator**(struct net\_rate\_estimator \_\_rcu \*\* *rate\_est*)  
remove a rate estimator

#### Parameters

struct net\_rate\_estimator \_\_rcu \*\* *rate\_est* rate estimator

#### Description

Removes the rate estimator.

int **gen\_replace\_estimator**(struct gnet\_stats\_basic\_packed \* *bstats*, struct gnet\_stats\_basic\_cpu \_\_percpu \* *cpu\_bstats*, struct net\_rate\_estimator \_\_rcu \*\* *rate\_est*, spinlock\_t \* *stats\_lock*, seqcount\_t \* *running*, struct nlattr \* *opt*)  
replace rate estimator configuration

#### Parameters

struct gnet\_stats\_basic\_packed \* *bstats* basic statistics

struct gnet\_stats\_basic\_cpu \_\_percpu \* *cpu\_bstats* bstats per cpu

struct net\_rate\_estimator \_\_rcu \*\* *rate\_est* rate estimator statistics

spinlock\_t \* *stats\_lock* statistics lock

seqcount\_t \* *running* qdisc running seqcount (might be NULL)

struct nlattr \* *opt* rate estimator configuration TLV

#### Description

Replaces the configuration of a rate estimator by calling [gen\\_kill\\_estimator\(\)](#) and [gen\\_new\\_estimator\(\)](#).

Returns 0 on success or a negative error code.

bool **gen\_estimator\_active**(struct net\_rate\_estimator \_\_rcu \*\* *rate\_est*)  
test if estimator is currently in use

#### Parameters

struct net\_rate\_estimator \_\_rcu \*\* *rate\_est* rate estimator

#### Description

Returns true if estimator is active, and false if not.

## SUN RPC subsystem

\_\_be32 \* **xdr\_encode\_opaque\_fixed**(\_\_be32 \* *p*, const void \* *ptr*, unsigned int *nbytes*)  
Encode fixed length opaque data

#### Parameters

\_\_be32 \* *p* pointer to current position in XDR buffer.

const void \* *ptr* pointer to data to encode (or NULL)

unsigned int *nbytes* size of data.

#### Description

Copy the array of data of length *nbytes* at *ptr* to the XDR buffer at position *p*, then align to the next 32-bit boundary by padding with zero bytes (see RFC1832).

#### Note

if `ptr` is `NULL`, only the padding is performed.

Returns the updated current XDR buffer position

`__be32 * xdr_encode_opaque(__be32 * p, const void * ptr, unsigned int nbytes)`  
Encode variable length opaque data

#### Parameters

`__be32 * p` pointer to current position in XDR buffer.

`const void * ptr` pointer to data to encode (or `NULL`)

`unsigned int nbytes` size of data.

#### Description

Returns the updated current XDR buffer position

`void xdr_terminate_string(struct xdr_buf * buf, const u32 len)`  
'0'-terminate a string residing in an `xdr_buf`

#### Parameters

`struct xdr_buf * buf` XDR buffer where string resides

`const u32 len` length of string, in bytes

`void _copy_from_pages(char * p, struct page ** pages, size_t pgbase, size_t len)`

#### Parameters

`char * p` pointer to destination

`struct page ** pages` array of pages

`size_t pgbase` offset of source data

`size_t len` length

#### Description

Copies data into an arbitrary memory location from an array of pages The copy is assumed to be non-overlapping.

`unsigned int xdr_stream_pos(const struct xdr_stream * xdr)`  
Return the current offset from the start of the `xdr_stream`

#### Parameters

`const struct xdr_stream * xdr` pointer to struct `xdr_stream`

`void xdr_init_encode(struct xdr_stream * xdr, struct xdr_buf * buf, __be32 * p)`  
Initialize a struct `xdr_stream` for sending data.

#### Parameters

`struct xdr_stream * xdr` pointer to `xdr_stream` struct

`struct xdr_buf * buf` pointer to XDR buffer in which to encode data

`__be32 * p` current pointer inside XDR buffer

#### Note

**at the moment the RPC client only passes the length of our** scratch buffer in the `xdr_buf`'s header `kvec`. Previously this meant we needed to call `xdr_adjust_iovec()` after encoding the data. With the new scheme, the `xdr_stream` manages the details of the buffer length, and takes care of adjusting the `kvec` length for us.

`void xdr_commit_encode(struct xdr_stream * xdr)`  
Ensure all data is written to buffer

#### Parameters

**struct xdr\_stream \* xdr** pointer to xdr\_stream

### Description

We handle encoding across page boundaries by giving the caller a temporary location to write to, then later copying the data into place; xdr\_commit\_encode does that copying.

Normally the caller doesn't need to call this directly, as the following xdr\_reserve\_space will do it. But an explicit call may be required at the end of encoding, or any other time when the xdr\_buf data might be read.

**\_\_be32 \* xdr\_reserve\_space**(struct xdr\_stream \* *xdr*, size\_t *nbytes*)  
Reserve buffer space for sending

### Parameters

**struct xdr\_stream \* xdr** pointer to xdr\_stream

**size\_t nbytes** number of bytes to reserve

### Description

Checks that we have enough buffer space to encode 'nbytes' more bytes of data. If so, update the total xdr\_buf length, and adjust the length of the current kvec.

**void xdr\_truncate\_encode**(struct xdr\_stream \* *xdr*, size\_t *len*)  
truncate an encode buffer

### Parameters

**struct xdr\_stream \* xdr** pointer to xdr\_stream

**size\_t len** new length of buffer

### Description

Truncates the xdr stream, so that xdr->buf->len == len, and xdr->p points at offset len from the start of the buffer, and head, tail, and page lengths are adjusted to correspond.

If this means moving xdr->p to a different buffer, we assume that that the end pointer should be set to the end of the current page, except in the case of the head buffer when we assume the head buffer's current length represents the end of the available buffer.

This is *not* safe to use on a buffer that already has inlined page cache pages (as in a zero-copy server read reply), except for the simple case of truncating from one position in the tail to another.

**int xdr\_restrict\_buflen**(struct xdr\_stream \* *xdr*, int *newbuflen*)  
decrease available buffer space

### Parameters

**struct xdr\_stream \* xdr** pointer to xdr\_stream

**int newbuflen** new maximum number of bytes available

### Description

Adjust our idea of how much space is available in the buffer. If we've already used too much space in the buffer, returns -1. If the available space is already smaller than newbuflen, returns 0 and does nothing. Otherwise, adjusts xdr->buf->buflen to newbuflen and ensures xdr->end is set at most offset newbuflen from the start of the buffer.

**void xdr\_write\_pages**(struct xdr\_stream \* *xdr*, struct page \*\* *pages*, unsigned int *base*, unsigned int *len*)  
Insert a list of pages into an XDR buffer for sending

### Parameters

**struct xdr\_stream \* xdr** pointer to xdr\_stream

**struct page \*\* pages** list of pages

**unsigned int base** offset of first byte

**unsigned int len** length of data in bytes

**void xdr\_init\_decode**(struct xdr\_stream \* *xdr*, struct xdr\_buf \* *buf*, \_\_be32 \* *p*)  
Initialize an xdr\_stream for decoding data.

#### Parameters

**struct xdr\_stream \* xdr** pointer to xdr\_stream struct

**struct xdr\_buf \* buf** pointer to XDR buffer from which to decode data

**\_\_be32 \* p** current pointer inside XDR buffer

**void xdr\_init\_decode\_pages**(struct xdr\_stream \* *xdr*, struct xdr\_buf \* *buf*, struct page \*\* *pages*,  
unsigned int *len*)  
Initialize an xdr\_stream for decoding into pages

#### Parameters

**struct xdr\_stream \* xdr** pointer to xdr\_stream struct

**struct xdr\_buf \* buf** pointer to XDR buffer from which to decode data

**struct page \*\* pages** list of pages to decode into

**unsigned int len** length in bytes of buffer in pages

**void xdr\_set\_scratch\_buffer**(struct xdr\_stream \* *xdr*, void \* *buf*, size\_t *buflen*)  
Attach a scratch buffer for decoding data.

#### Parameters

**struct xdr\_stream \* xdr** pointer to xdr\_stream struct

**void \* buf** pointer to an empty buffer

**size\_t buflen** size of 'buf'

#### Description

The scratch buffer is used when decoding from an array of pages. If an [\*xdr\\_inline\\_decode\(\)\*](#) call spans across page boundaries, then we copy the data into the scratch buffer in order to allow linear access.

**\_\_be32 \* xdr\_inline\_decode**(struct xdr\_stream \* *xdr*, size\_t *nbytes*)  
Retrieve XDR data to decode

#### Parameters

**struct xdr\_stream \* xdr** pointer to xdr\_stream struct

**size\_t nbytes** number of bytes of data to decode

#### Description

Check if the input buffer is long enough to enable us to decode 'nbytes' more bytes of data starting at the current position. If so return the current pointer, then update the current pointer position.

**unsigned int xdr\_read\_pages**(struct xdr\_stream \* *xdr*, unsigned int *len*)  
Ensure page-based XDR data to decode is aligned at current pointer position

#### Parameters

**struct xdr\_stream \* xdr** pointer to xdr\_stream struct

**unsigned int len** number of bytes of page data

#### Description

Moves data beyond the current pointer position from the XDR head[] buffer into the page list. Any data that lies beyond current position + "len" bytes is moved into the XDR tail[].

Returns the number of XDR encoded bytes now contained in the pages

void **xdr\_enter\_page**(struct xdr\_stream \* *xdr*, unsigned int *len*)  
decode data from the XDR page

#### Parameters

**struct xdr\_stream \* xdr** pointer to xdr\_stream struct

**unsigned int len** number of bytes of page data

#### Description

Moves data beyond the current pointer position from the XDR head[] buffer into the page list. Any data that lies beyond current position + “len” bytes is moved into the XDR tail[]. The current pointer is then repositioned at the beginning of the first XDR page.

int **xdr\_buf\_subsegment**(struct xdr\_buf \* *buf*, struct xdr\_buf \* *subbuf*, unsigned int *base*, unsigned int *len*)  
set subbuf to a portion of buf

#### Parameters

**struct xdr\_buf \* buf** an xdr buffer

**struct xdr\_buf \* subbuf** the result buffer

**unsigned int base** beginning of range in bytes

**unsigned int len** length of range in bytes

#### Description

sets **subbuf** to an xdr buffer representing the portion of **buf** of length **len** starting at offset **base**.

**buf** and **subbuf** may be pointers to the same struct xdr\_buf.

Returns -1 if base of length are out of bounds.

void **xdr\_buf\_trim**(struct xdr\_buf \* *buf*, unsigned int *len*)  
lop at most “len” bytes off the end of “buf”

#### Parameters

**struct xdr\_buf \* buf** buf to be trimmed

**unsigned int len** number of bytes to reduce “buf” by

#### Description

Trim an xdr\_buf by the given number of bytes by fixing up the lengths. Note that it’s possible that we’ll trim less than that amount if the xdr\_buf is too small, or if (for instance) it’s all in the head and the parser has already read too far into it.

ssize\_t **xdr\_stream\_decode\_string\_dup**(struct xdr\_stream \* *xdr*, char \*\* *str*, size\_t *maxlen*, gfp\_t *gfp\_flags*)  
Decode and duplicate variable length string

#### Parameters

**struct xdr\_stream \* xdr** pointer to xdr\_stream

**char \*\* str** location to store pointer to string

**size\_t maxlen** maximum acceptable string length

**gfp\_t gfp\_flags** GFP mask to use

#### Description

**Return values:** On success, returns length of NUL-terminated string stored in **\*ptr** -EBADMSG on XDR buffer overflow -EMSGSIZE if the size of the string would exceed **maxlen** -ENOMEM on memory allocation failure

`char * svc_print_addr(struct svc_rqst * rqstp, char * buf, size_t len)`  
Format `rq_addr` field for printing

#### Parameters

`struct svc_rqst * rqstp` `svc_rqst` struct containing address to print

`char * buf` target buffer for formatted address

`size_t len` length of target buffer

`void svc_reserve(struct svc_rqst * rqstp, int space)`  
change the space reserved for the reply to a request.

#### Parameters

`struct svc_rqst * rqstp` The request in question

`int space` new max space to reserve

#### Description

Each request reserves some space on the output queue of the transport to make sure the reply fits. This function reduces that reserved space to be the amount of space used already, plus **space**.

`struct svc_xprt * svc_find_xprt(struct svc_serv * serv, const char * xcl_name, struct net * net,  
const sa_family_t af, const unsigned short port)`  
find an RPC transport instance

#### Parameters

`struct svc_serv * serv` pointer to `svc_serv` to search

`const char * xcl_name` C string containing transport's class name

`struct net * net` owner net pointer

`const sa_family_t af` Address family of transport's local address

`const unsigned short port` transport's IP port number

#### Description

Return the transport instance pointer for the endpoint accepting connections/peer traffic from the specified transport class, address family and port.

Specifying 0 for the address family or port is effectively a wild-card, and will result in matching the first transport in the service's list that has a matching class name.

`int svc_xprt_names(struct svc_serv * serv, char * buf, const int buflen)`  
format a buffer with a list of transport names

#### Parameters

`struct svc_serv * serv` pointer to an RPC service

`char * buf` pointer to a buffer to be filled in

`const int buflen` length of buffer to be filled in

#### Description

Fills in **buf** with a string containing a list of transport names, each name terminated with 'n'.

Returns positive length of the filled-in string on success; otherwise a negative `errno` value is returned if an error occurs.

`int xprt_register_transport(struct xprt_class * transport)`  
register a transport implementation

#### Parameters

`struct xprt_class * transport` transport to register

## Description

If a transport implementation is loaded as a kernel module, it can call this interface to make itself known to the RPC client.

## Return

0: transport successfully registered -EEXIST: transport already registered -EINVAL: transport module being unloaded

int **xprt\_unregister\_transport**(struct xprt\_class \* *transport*)  
unregister a transport implementation

## Parameters

**struct xprt\_class \* transport** transport to unregister

## Return

0: transport successfully unregistered -ENOENT: transport never registered

int **xprt\_load\_transport**(const char \* *transport\_name*)  
load a transport implementation

## Parameters

**const char \* transport\_name** transport to load

## Return

0: transport successfully loaded -ENOENT: transport module not available

int **xprt\_reserve\_xprt**(struct rpc\_xprt \* *xprt*, struct rpc\_task \* *task*)  
serialize write access to transports

## Parameters

**struct rpc\_xprt \* xprt** pointer to the target transport

**struct rpc\_task \* task** task that is requesting access to the transport

## Description

This prevents mixing the payload of separate requests, and prevents transport connects from colliding with writes. No congestion control is provided.

void **xprt\_release\_xprt**(struct rpc\_xprt \* *xprt*, struct rpc\_task \* *task*)  
allow other requests to use a transport

## Parameters

**struct rpc\_xprt \* xprt** transport with other tasks potentially waiting

**struct rpc\_task \* task** task that is releasing access to the transport

## Description

Note that “task” can be NULL. No congestion control is provided.

void **xprt\_release\_xprt\_cong**(struct rpc\_xprt \* *xprt*, struct rpc\_task \* *task*)  
allow other requests to use a transport

## Parameters

**struct rpc\_xprt \* xprt** transport with other tasks potentially waiting

**struct rpc\_task \* task** task that is releasing access to the transport

## Description

Note that “task” can be NULL. Another task is awoken to use the transport if the transport’s congestion window allows it.



void **xprt\_release\_rqst\_cong**(struct rpc\_task \* *task*)  
housekeeping when request is complete

#### Parameters

**struct rpc\_task \* task** RPC request that recently completed

#### Description

Useful for transports that require congestion control.

void **xprt\_adjust\_cwnd**(struct rpc\_xprt \* *xprt*, struct rpc\_task \* *task*, int *result*)  
adjust transport congestion window

#### Parameters

**struct rpc\_xprt \* xprt** pointer to xprt

**struct rpc\_task \* task** recently completed RPC request used to adjust window

**int result** result code of completed RPC request

#### Description

The transport code maintains an estimate on the maximum number of out-standing RPC requests, using a smoothed version of the congestion avoidance implemented in 44BSD. This is basically the Van Jacobson congestion algorithm: If a retransmit occurs, the congestion window is halved; otherwise, it is incremented by 1/cwnd when

- a reply is received and
- a full number of requests are outstanding and
- the congestion window hasn't been updated recently.

void **xprt\_wake\_pending\_tasks**(struct rpc\_xprt \* *xprt*, int *status*)  
wake all tasks on a transport's pending queue

#### Parameters

**struct rpc\_xprt \* xprt** transport with waiting tasks

**int status** result code to plant in each task before waking it

void **xprt\_wait\_for\_buffer\_space**(struct rpc\_task \* *task*, rpc\_action *action*)  
wait for transport output buffer to clear

#### Parameters

**struct rpc\_task \* task** task to be put to sleep

**rpc\_action action** function pointer to be executed after wait

#### Description

Note that we only set the timer for the case of `RPC_IS_SOFT()`, since we don't in general want to force a socket disconnection due to an incomplete RPC call transmission.

void **xprt\_write\_space**(struct rpc\_xprt \* *xprt*)  
wake the task waiting for transport output buffer space

#### Parameters

**struct rpc\_xprt \* xprt** transport with waiting tasks

#### Description

Can be called in a soft IRQ context, so `xprt_write_space` never sleeps.

void **xprt\_set\_retrans\_timeout\_def**(struct rpc\_task \* *task*)  
set a request's retransmit timeout

#### Parameters

**struct rpc\_task \* task** task whose timeout is to be set

### Description

Set a request's retransmit timeout based on the transport's default timeout parameters. Used by transports that don't adjust the retransmit timeout based on round-trip time estimation.

void **xprt\_set\_retrans\_timeout\_rtt**(struct rpc\_task \* *task*)  
set a request's retransmit timeout

### Parameters

**struct rpc\_task \* task** task whose timeout is to be set

### Description

Set a request's retransmit timeout using the RTT estimator.

void **xprt\_disconnect\_done**(struct rpc\_xprt \* *xprt*)  
mark a transport as disconnected

### Parameters

**struct rpc\_xprt \* xprt** transport to flag for disconnect

void **xprt\_force\_disconnect**(struct rpc\_xprt \* *xprt*)  
force a transport to disconnect

### Parameters

**struct rpc\_xprt \* xprt** transport to disconnect

struct rpc\_rqst \* **xprt\_lookup\_rqst**(struct rpc\_xprt \* *xprt*, \_\_be32 *xid*)  
find an RPC request corresponding to an XID

### Parameters

**struct rpc\_xprt \* xprt** transport on which the original request was transmitted

**\_\_be32 xid** RPC XID of incoming reply

void **xprt\_pin\_rqst**(struct rpc\_rqst \* *req*)  
Pin a request on the transport receive list

### Parameters

**struct rpc\_rqst \* req** Request to pin

### Description

Caller must ensure this is atomic with the call to [\*xprt\\_lookup\\_rqst\(\)\*](#) so should be holding the xprt transport lock.

void **xprt\_unpin\_rqst**(struct rpc\_rqst \* *req*)  
Unpin a request on the transport receive list

### Parameters

**struct rpc\_rqst \* req** Request to pin

### Description

Caller should be holding the xprt transport lock.

void **xprt\_complete\_rqst**(struct rpc\_task \* *task*, int *copied*)  
called when reply processing is complete

### Parameters

**struct rpc\_task \* task** RPC request that recently completed

**int copied** actual number of bytes received from the transport

**Description**

Caller holds transport lock.

**struct rpc\_xprt \* *xprt\_get***(struct rpc\_xprt \* *xprt*)  
return a reference to an RPC transport.

**Parameters**

**struct rpc\_xprt \* *xprt*** pointer to the transport

**void *xprt\_put***(struct rpc\_xprt \* *xprt*)  
release a reference to an RPC transport.

**Parameters**

**struct rpc\_xprt \* *xprt*** pointer to the transport

**void *rpc\_wake\_up***(struct rpc\_wait\_queue \* *queue*)  
wake up all *rpc\_tasks*

**Parameters**

**struct rpc\_wait\_queue \* *queue*** *rpc\_wait\_queue* on which the tasks are sleeping

**Description**

Grabs *queue->lock*

**void *rpc\_wake\_up\_status***(struct rpc\_wait\_queue \* *queue*, int *status*)  
wake up all *rpc\_tasks* and set their status value.

**Parameters**

**struct rpc\_wait\_queue \* *queue*** *rpc\_wait\_queue* on which the tasks are sleeping

**int *status*** status value to set

**Description**

Grabs *queue->lock*

**int *rpc\_malloc***(struct rpc\_task \* *task*)  
allocate RPC buffer resources

**Parameters**

**struct rpc\_task \* *task*** RPC task

**Description**

A single memory region is allocated, which is split between the RPC call and RPC reply that this task is being used for. When this RPC is retired, the memory is released by calling *rpc\_free*.

To prevent *rpciod* from hanging, this allocator never sleeps, returning *-ENOMEM* and suppressing warning if the request cannot be serviced immediately. The caller can arrange to sleep in a way that is safe for *rpciod*.

Most requests are 'small' (under 2KiB) and can be serviced from a mempool, ensuring that NFS reads and writes can always proceed, and that there is good locality of reference for these buffers.

In order to avoid memory starvation triggering more writebacks of NFS requests, we avoid using *GFP\_KERNEL*.

**void *rpc\_free***(struct rpc\_task \* *task*)  
free RPC buffer resources allocated via *rpc\_malloc*

**Parameters**

**struct rpc\_task \* *task*** RPC task

**size\_t *xdr\_skb\_read\_bits***(struct xdr\_skb\_reader \* *desc*, void \* *to*, size\_t *len*)  
copy some data bits from *skb* to internal buffer

### Parameters

**struct xdr\_skb\_reader \* desc** sk\_buff copy helper

**void \* to** copy destination

**size\_t len** number of bytes to copy

### Description

Possibly called several times to iterate over an sk\_buff and copy data out of it.

**ssize\_t xdr\_partial\_copy\_from\_skb**(struct xdr\_buf \* *xdr*, unsigned int *base*, struct xdr\_skb\_reader \* *desc*, xdr\_skb\_read\_actor *copy\_actor*)  
copy data out of an skb

### Parameters

**struct xdr\_buf \* xdr** target XDR buffer

**unsigned int base** starting offset

**struct xdr\_skb\_reader \* desc** sk\_buff copy helper

**xdr\_skb\_read\_actor copy\_actor** virtual method for copying data

**int csum\_partial\_copy\_to\_xdr**(struct xdr\_buf \* *xdr*, struct *sk\_buff* \* *skb*)  
checksum and copy data

### Parameters

**struct xdr\_buf \* xdr** target XDR buffer

**struct sk\_buff \* skb** source skb

### Description

We have set things up such that we perform the checksum of the UDP packet in parallel with the copies into the RPC client iovec. -DaveM

**struct rpc\_iostats \* rpc\_alloc\_iostats**(struct rpc\_clnt \* *clnt*)  
allocate an rpc\_iostats structure

### Parameters

**struct rpc\_clnt \* clnt** RPC program, version, and xprt

**void rpc\_free\_iostats**(struct rpc\_iostats \* *stats*)  
release an rpc\_iostats structure

### Parameters

**struct rpc\_iostats \* stats** doomed rpc\_iostats structure

**void rpc\_count\_iostats\_metrics**(const struct rpc\_task \* *task*, struct rpc\_iostats \* *op\_metrics*)  
tally up per-task stats

### Parameters

**const struct rpc\_task \* task** completed rpc\_task

**struct rpc\_iostats \* op\_metrics** stat structure for OP that will accumulate stats from **task**

**void rpc\_count\_iostats**(const struct rpc\_task \* *task*, struct rpc\_iostats \* *stats*)  
tally up per-task stats

### Parameters

**const struct rpc\_task \* task** completed rpc\_task

**struct rpc\_iostats \* stats** array of stat structures

### Description

Uses the statidx from **task**

int **rpc\_queue\_upcall**(struct rpc\_pipe \* *pipe*, struct rpc\_pipe\_msg \* *msg*)  
queue an upcall message to userspace

#### Parameters

**struct rpc\_pipe \* pipe** upcall pipe on which to queue given message

**struct rpc\_pipe\_msg \* msg** message to queue

#### Description

Call with an **inode** created by `rpc_mkpipe()` to queue an upcall. A userspace process may then later read the upcall by performing a read on an open file for this inode. It is up to the caller to initialize the fields of **msg** (other than **msg->list**) appropriately.

**struct dentry \* rpc\_mkpipe\_dentry**(struct dentry \* *parent*, const char \* *name*, void \* *private*, struct  
rpc\_pipe \* *pipe*)  
make an rpc\_pipefs file for kernel<->userspace communication

#### Parameters

**struct dentry \* parent** dentry of directory to create new “pipe” in

**const char \* name** name of pipe

**void \* private** private data to associate with the pipe, for the caller’s use

**struct rpc\_pipe \* pipe** rpc\_pipe containing input parameters

#### Description

Data is made available for userspace to read by calls to `rpc_queue_upcall()`. The actual reads will result in calls to **ops->upcall**, which will be called with the file pointer, message, and userspace buffer to copy to.

Writes can come at any time, and do not necessarily have to be responses to upcalls. They will result in calls to **msg->downcall**.

The **private** argument passed here will be available to all these methods from the file pointer, via `RPC_I(file_inode(file))->private`.

int **rpc\_unlink**(struct dentry \* *dentry*)  
remove a pipe

#### Parameters

**struct dentry \* dentry** dentry for the pipe, as returned from `rpc_mkpipe`

#### Description

After this call, lookups will no longer find the pipe, and any attempts to read or write using preexisting opens of the pipe will return -EPIPE.

void **rpc\_init\_pipe\_dir\_head**(struct rpc\_pipe\_dir\_head \* *pdh*)  
initialise a struct rpc\_pipe\_dir\_head

#### Parameters

**struct rpc\_pipe\_dir\_head \* pdh** pointer to struct rpc\_pipe\_dir\_head

void **rpc\_init\_pipe\_dir\_object**(struct rpc\_pipe\_dir\_object \* *pdo*, const struct  
rpc\_pipe\_dir\_object\_ops \* *pdo\_ops*, void \* *pdo\_data*)  
initialise a struct rpc\_pipe\_dir\_object

#### Parameters

**struct rpc\_pipe\_dir\_object \* pdo** pointer to struct rpc\_pipe\_dir\_object

**const struct rpc\_pipe\_dir\_object\_ops \* pdo\_ops** pointer to const struct rpc\_pipe\_dir\_object\_ops

**void \* pdo\_data** pointer to caller-defined data

```
int rpc_add_pipe_dir_object(struct net *net, struct rpc_pipe_dir_head *pdh, struct
                           rpc_pipe_dir_object *pdo)
    associate a rpc_pipe_dir_object to a directory
```

#### Parameters

**struct net \* net** pointer to struct net

**struct rpc\_pipe\_dir\_head \* pdh** pointer to struct rpc\_pipe\_dir\_head

**struct rpc\_pipe\_dir\_object \* pdo** pointer to struct rpc\_pipe\_dir\_object

```
void rpc_remove_pipe_dir_object(struct net *net, struct rpc_pipe_dir_head *pdh, struct
                               rpc_pipe_dir_object *pdo)
    remove a rpc_pipe_dir_object from a directory
```

#### Parameters

**struct net \* net** pointer to struct net

**struct rpc\_pipe\_dir\_head \* pdh** pointer to struct rpc\_pipe\_dir\_head

**struct rpc\_pipe\_dir\_object \* pdo** pointer to struct rpc\_pipe\_dir\_object

```
struct rpc_pipe_dir_object * rpc_find_or_alloc_pipe_dir_object(struct net *net, struct
                                                                rpc_pipe_dir_head *pdh,
                                                                int (*match) (struct
                                                                rpc_pipe_dir_object *, void *,
                                                                struct rpc_pipe_dir_object
                                                                *(*alloc) (void *, void * data)
```

#### Parameters

**struct net \* net** pointer to struct net

**struct rpc\_pipe\_dir\_head \* pdh** pointer to struct rpc\_pipe\_dir\_head

**int (\*)(struct rpc\_pipe\_dir\_object \*, void \*) match** match struct rpc\_pipe\_dir\_object to data

**struct rpc\_pipe\_dir\_object (\*)(void \*) alloc** allocate a new struct rpc\_pipe\_dir\_object

**void \* data** user defined data for match() and alloc()

```
void rpcb_getport_async(struct rpc_task *task)
    obtain the port for a given RPC service on a given host
```

#### Parameters

**struct rpc\_task \* task** task that is waiting for portmapper request

#### Description

This one can be called for an ongoing RPC request, and can be used in an async (rpciod) context.

```
struct rpc_clnt * rpc_create(struct rpc_create_args *args)
    create an RPC client and transport with one call
```

#### Parameters

**struct rpc\_create\_args \* args** rpc\_clnt create argument structure

#### Description

Creates and initializes an RPC transport and an RPC client.

It can ping the server in order to determine if it is up, and to see if it supports this program and version. RPC\_CLNT\_CREATE\_NOPING disables this behavior so asynchronous tasks can also use rpc\_create.

```
struct rpc_clnt * rpc_clone_client(struct rpc_clnt *clnt)
    Clone an RPC client structure
```

#### Parameters

**struct rpc\_clnt \* clnt** RPC client whose parameters are copied

### Description

Returns a fresh RPC client or an ERR\_PTR.

**struct rpc\_clnt \* rpc\_clone\_client\_set\_auth**(**struct rpc\_clnt \* clnt**, **rpc\_authflavor\_t flavor**)  
Clone an RPC client structure and set its auth

### Parameters

**struct rpc\_clnt \* clnt** RPC client whose parameters are copied

**rpc\_authflavor\_t flavor** security flavor for new client

### Description

Returns a fresh RPC client or an ERR\_PTR.

**int rpc\_switch\_client\_transport**(**struct rpc\_clnt \* clnt**, **struct xprt\_create \* args**, **const struct rpc\_timeout \* timeout**)

### Parameters

**struct rpc\_clnt \* clnt** pointer to a struct rpc\_clnt

**struct xprt\_create \* args** pointer to the new transport arguments

**const struct rpc\_timeout \* timeout** pointer to the new timeout parameters

### Description

This function allows the caller to switch the RPC transport for the **rpc\_clnt** structure 'clnt' to allow it to connect to a mirrored NFS server, for instance. It assumes that the caller has ensured that there are no active RPC tasks by using some form of locking.

Returns zero if "clnt" is now using the new xprt. Otherwise a negative errno is returned, and "clnt" continues to use the old xprt.

**int rpc\_clnt\_iterate\_for\_each\_xprt**(**struct rpc\_clnt \* clnt**, **int (\*fn)** (**struct rpc\_clnt \***, **struct rpc\_xprt \***, **void \***, **void \* data**)

Apply a function to all transports

### Parameters

**struct rpc\_clnt \* clnt** pointer to client

**int (\*)(struct rpc\_clnt \*, struct rpc\_xprt \*, void \*) fn** function to apply

**void \* data** void pointer to function data

### Description

Iterates through the list of RPC transports currently attached to the client and applies the function **fn**(**clnt**, **xprt**, **data**).

On error, the iteration stops, and the function returns the error value.

**struct rpc\_clnt \* rpc\_bind\_new\_program**(**struct rpc\_clnt \* old**, **const struct rpc\_program \* program**, **u32 vers**)

bind a new RPC program to an existing client

### Parameters

**struct rpc\_clnt \* old** old rpc\_client

**const struct rpc\_program \* program** rpc program to set

**u32 vers** rpc program version

### Description

Clones the rpc client and sets up a new RPC program. This is mainly of use for enabling different RPC programs to share the same transport. The Sun NFSv2/v3 ACL protocol can do this.

**struct rpc\_task \* rpc\_run\_task**(const struct rpc\_task\_setup \* *task\_setup\_data*)  
Allocate a new RPC task, then run `rpc_execute` against it

#### Parameters

**const struct rpc\_task\_setup \* task\_setup\_data** pointer to task initialisation data

**int rpc\_call\_sync**(struct rpc\_clnt \* *clnt*, const struct rpc\_message \* *msg*, int *flags*)  
Perform a synchronous RPC call

#### Parameters

**struct rpc\_clnt \* clnt** pointer to RPC client

**const struct rpc\_message \* msg** RPC call parameters

**int flags** RPC call flags

**int rpc\_call\_async**(struct rpc\_clnt \* *clnt*, const struct rpc\_message \* *msg*, int *flags*, const struct  
rpc\_call\_ops \* *tk\_ops*, void \* *data*)  
Perform an asynchronous RPC call

#### Parameters

**struct rpc\_clnt \* clnt** pointer to RPC client

**const struct rpc\_message \* msg** RPC call parameters

**int flags** RPC call flags

**const struct rpc\_call\_ops \* tk\_ops** RPC call ops

**void \* data** user call data

**size\_t rpc\_peeraddr**(struct rpc\_clnt \* *clnt*, struct sockaddr \* *buf*, size\_t *bufsize*)  
extract remote peer address from *clnt*'s *xprt*

#### Parameters

**struct rpc\_clnt \* clnt** RPC client structure

**struct sockaddr \* buf** target buffer

**size\_t bufsize** length of target buffer

#### Description

Returns the number of bytes that are actually in the stored address.

**const char \* rpc\_peeraddr2str**(struct rpc\_clnt \* *clnt*, enum rpc\_display\_format\_t *format*)  
return remote peer address in printable format

#### Parameters

**struct rpc\_clnt \* clnt** RPC client structure

**enum rpc\_display\_format\_t format** address format

#### Description

NB: the lifetime of the memory referenced by the returned pointer is the same as the *rpc\_xprt* itself. As long as the caller uses this pointer, it must hold the RCU read lock.

**int rpc\_localaddr**(struct rpc\_clnt \* *clnt*, struct sockaddr \* *buf*, size\_t *buflen*)  
discover local endpoint address for an RPC client

#### Parameters

**struct rpc\_clnt \* clnt** RPC client structure

**struct sockaddr \* buf** target buffer

**size\_t buflen** size of target buffer, in bytes



**Description**

Returns zero and fills in “buf” and “buflen” if successful; otherwise, a negative errno is returned.

This works even if the underlying transport is not currently connected, or if the upper layer never previously provided a source address.

The result of this function call is transient: multiple calls in succession may give different results, depending on how local networking configuration changes over time.

int **rpc\_protocol**(struct rpc\_clnt \* *clnt*)  
Get transport protocol number for an RPC client

**Parameters**

struct rpc\_clnt \* *clnt* RPC client to query  
struct net \* **rpc\_net\_ns**(struct rpc\_clnt \* *clnt*)  
Get the network namespace for this RPC client

**Parameters**

struct rpc\_clnt \* *clnt* RPC client to query  
size\_t **rpc\_max\_payload**(struct rpc\_clnt \* *clnt*)  
Get maximum payload size for a transport, in bytes

**Parameters**

struct rpc\_clnt \* *clnt* RPC client to query

**Description**

For stream transports, this is one RPC record fragment (see RFC 1831), as we don't support multi-record requests yet. For datagram transports, this is the size of an IP packet minus the IP, UDP, and RPC header sizes.

size\_t **rpc\_max\_bc\_payload**(struct rpc\_clnt \* *clnt*)  
Get maximum backchannel payload size, in bytes

**Parameters**

struct rpc\_clnt \* *clnt* RPC client to query  
void **rpc\_force\_rebind**(struct rpc\_clnt \* *clnt*)  
force transport to check that remote port is unchanged

**Parameters**

struct rpc\_clnt \* *clnt* client to rebind  
int **rpc\_clnt\_test\_and\_add\_xprt**(struct rpc\_clnt \* *clnt*, struct rpc\_xprt\_switch \* *xps*, struct rpc\_xprt \* *xprt*, void \* *dummy*)  
Test and add a new transport to a rpc\_clnt

**Parameters**

struct rpc\_clnt \* *clnt* pointer to struct rpc\_clnt  
struct rpc\_xprt\_switch \* *xps* pointer to struct rpc\_xprt\_switch,  
struct rpc\_xprt \* *xprt* pointer struct rpc\_xprt  
void \* *dummy* unused  
int **rpc\_clnt\_setup\_test\_and\_add\_xprt**(struct rpc\_clnt \* *clnt*, struct rpc\_xprt\_switch \* *xps*, struct rpc\_xprt \* *xprt*, void \* *data*)

**Parameters**

struct rpc\_clnt \* *clnt* struct rpc\_clnt to get the new transport  
struct rpc\_xprt\_switch \* *xps* the rpc\_xprt\_switch to hold the new transport

**struct rpc\_xprt \* xprt** the `rpc_xprt` to test

**void \* data** a struct `rpc_add_xprt_test` pointer that holds the test function and test function call data

### Description

**This is an `rpc_clnt_add_xprt_setup()` function which returns 1 so:** 1) caller of the test function must dereference the `rpc_xprt_switch` and the `rpc_xprt`. 2) test function must call `rpc_xprt_switch_add_xprt`, usually in the `rpc_call_done` routine.

Upon success (return of 1), the test function adds the new transport to the `rpc_clnt` `xprt` switch

**int `rpc_clnt_add_xprt`**(struct `rpc_clnt` \* *clnt*, struct `xprt_create` \* *xprtargs*, int (\**setup*) (struct `rpc_clnt` \*, struct `rpc_xprt_switch` \*, struct `rpc_xprt` \*, void \*, void \* *data*)  
Add a new transport to a `rpc_clnt`

### Parameters

**struct `rpc_clnt` \* *clnt*** pointer to struct `rpc_clnt`

**struct `xprt_create` \* *xprtargs*** pointer to struct `xprt_create`

**int (\*) (struct `rpc_clnt` \*, struct `rpc_xprt_switch` \*, struct `rpc_xprt` \*, void \*) *setup***  
callback to test and/or set up the connection

**void \* *data*** pointer to setup function data

### Description

Creates a new transport using the parameters set in *args* and adds it to *clnt*. If *ping* is set, then test that connectivity succeeds before adding the new transport.

## WiMAX

**struct `sk_buff` \* `wimax_msg_alloc`**(struct `wimax_dev` \* *wimax\_dev*, const char \* *pipe\_name*, const void \* *msg*, `size_t` *size*, `gfp_t` *gfp\_flags*)  
Create a new `skb` for sending a message to userspace

### Parameters

**struct `wimax_dev` \* *wimax\_dev*** WiMAX device descriptor

**const char \* *pipe\_name*** “named pipe” the message will be sent to

**const void \* *msg*** pointer to the message data to send

**size\_t *size*** size of the message to send (in bytes), including the header.

**gfp\_t *gfp\_flags*** flags for memory allocation.

### Return

0 if ok, negative `errno` code on error

### Description

Allocates an `skb` that will contain the message to send to user space over the messaging pipe and initializes it, copying the payload.

Once this call is done, you can deliver it with `wimax_msg_send()`.

IMPORTANT:

Don't use `skb_push()/skb_pull()/skb_reserve()` on the `skb`, as `wimax_msg_send()` depends on `skb->data` being placed at the beginning of the user message.

Unlike other WiMAX stack calls, this call can be used way early, even before `wimax_dev_add()` is called, as long as the `wimax_dev->net_dev` pointer is set to point to a proper `net_dev`. This is so that drivers can use it early in case they need to send stuff around or communicate with user space.

```
const void * wimax_msg_data_len(struct sk_buff * msg, size_t * size)
```

Return a pointer and size of a message's payload

**Parameters**

**struct sk\_buff \* msg** Pointer to a message created with *wimax\_msg\_alloc()*

**size\_t \* size** Pointer to where to store the message's size

**Description**

Returns the pointer to the message data.

```
const void * wimax_msg_data(struct sk_buff * msg)
```

Return a pointer to a message's payload

**Parameters**

**struct sk\_buff \* msg** Pointer to a message created with *wimax\_msg\_alloc()*

```
ssize_t wimax_msg_len(struct sk_buff * msg)
```

Return a message's payload length

**Parameters**

**struct sk\_buff \* msg** Pointer to a message created with *wimax\_msg\_alloc()*

```
int wimax_msg_send(struct wimax_dev * wimax_dev, struct sk_buff * skb)
```

Send a pre-allocated message to user space

**Parameters**

**struct wimax\_dev \* wimax\_dev** WiMAX device descriptor

**struct sk\_buff \* skb** *struct sk\_buff* returned by *wimax\_msg\_alloc()*. Note the ownership of **skb** is transferred to this function.

**Return**

0 if ok, < 0 errno code on error

**Description**

Sends a free-form message that was preallocated with *wimax\_msg\_alloc()* and filled up.

Assumes that once you pass an skb to this function for sending, it owns it and will release it when done (on success).

IMPORTANT:

Don't use *skb\_push()/skb\_pull()/skb\_reserve()* on the skb, as *wimax\_msg\_send()* depends on *skb->data* being placed at the beginning of the user message.

Unlike other WiMAX stack calls, this call can be used way early, even before *wimax\_dev\_add()* is called, as long as the *wimax\_dev->net\_dev* pointer is set to point to a proper *net\_dev*. This is so that drivers can use it early in case they need to send stuff around or communicate with user space.

```
int wimax_msg(struct wimax_dev * wimax_dev, const char * pipe_name, const void * buf, size_t size,  
               gfp_t gfp_flags)
```

Send a message to user space

**Parameters**

**struct wimax\_dev \* wimax\_dev** WiMAX device descriptor (properly referenced)

**const char \* pipe\_name** "named pipe" the message will be sent to

**const void \* buf** pointer to the message to send.

**size\_t size** size of the buffer pointed to by **buf** (in bytes).

**gfp\_t gfp\_flags** flags for memory allocation.

### Return

0 if ok, negative errno code on error.

### Description

Sends a free-form message to user space on the device **wimax\_dev**.

### NOTES

Once the **skb** is given to this function, who will own it and will release it when done (unless it returns error).

```
int wimax_reset(struct wimax_dev * wimax_dev)
    Reset a WiMAX device
```

### Parameters

**struct wimax\_dev \* wimax\_dev** WiMAX device descriptor

### Return

0 if ok and a warm reset was done (the device still exists in the system).

-ENODEV if a cold/bus reset had to be done (device has disconnected and reconnected, so current handle is not valid any more).

-EINVAL if the device is not even registered.

Any other negative error code shall be considered as non-recoverable.

### Description

Called when wanting to reset the device for any reason. Device is taken back to power on status.

This call blocks; on successful return, the device has completed the reset process and is ready to operate.

```
void wimax_report_rfkill_hw(struct wimax_dev * wimax_dev, enum wimax_rf_state state)
    Reports changes in the hardware RF switch
```

### Parameters

**struct wimax\_dev \* wimax\_dev** WiMAX device descriptor

**enum wimax\_rf\_state state** New state of the RF Kill switch. WIMAX\_RF\_ON radio on, WIMAX\_RF\_OFF radio off.

### Description

When the device detects a change in the state of the hardware RF switch, it must call this function to let the WiMAX kernel stack know that the state has changed so it can be properly propagated.

The WiMAX stack caches the state (the driver doesn't need to). As well, as the change is propagated it will come back as a request to change the software state to mirror the hardware state.

If the device doesn't have a hardware kill switch, just report it on initialization as always on (WIMAX\_RF\_ON, radio on).

```
void wimax_report_rfkill_sw(struct wimax_dev * wimax_dev, enum wimax_rf_state state)
    Reports changes in the software RF switch
```

### Parameters

**struct wimax\_dev \* wimax\_dev** WiMAX device descriptor

**enum wimax\_rf\_state state** New state of the RF kill switch. WIMAX\_RF\_ON radio on, WIMAX\_RF\_OFF radio off.

### Description

Reports changes in the software RF switch state to the WiMAX stack.

The main use is during initialization, so the driver can query the device for its current software radio kill switch state and feed it to the system.

On the side, the device does not change the software state by itself. In practice, this can happen, as the device might decide to switch (in software) the radio off for different reasons.

int **wimax\_rfkill**(struct *wimax\_dev* \* *wimax\_dev*, enum wimax\_rf\_state *state*)  
Set the software RF switch state for a WiMAX device

#### Parameters

**struct wimax\_dev \* wimax\_dev** WiMAX device descriptor

**enum wimax\_rf\_state state** New RF state.

#### Return

>= 0 toggle state if ok, < 0 errno code on error. The toggle state is returned as a bitmap, bit 0 being the hardware RF state, bit 1 the software RF state.

0 means disabled (WIMAX\_RF\_ON, radio on), 1 means enabled radio off (WIMAX\_RF\_OFF).

#### Description

Called by the user when he wants to request the WiMAX radio to be switched on (WIMAX\_RF\_ON) or off (WIMAX\_RF\_OFF). With WIMAX\_RF\_QUERY, just the current state is returned.

#### NOTE

This call will block until the operation is complete.

void **wimax\_state\_change**(struct *wimax\_dev* \* *wimax\_dev*, enum *wimax\_st* *new\_state*)  
Set the current state of a WiMAX device

#### Parameters

**struct wimax\_dev \* wimax\_dev** WiMAX device descriptor (properly referenced)

**enum wimax\_st new\_state** New state to switch to

#### Description

This implements the state changes for the wimax devices. It will

- verify that the state transition is legal (for now it'll just print a warning if not) according to the table in linux/wimax.h's documentation for 'enum wimax\_st'.
- perform the actions needed for leaving the current state and whichever are needed for entering the new state.
- issue a report to user space indicating the new state (and an optional payload with information about the new state).

#### NOTE

**wimax\_dev** must be locked

enum *wimax\_st* **wimax\_state\_get**(struct *wimax\_dev* \* *wimax\_dev*)  
Return the current state of a WiMAX device

#### Parameters

**struct wimax\_dev \* wimax\_dev** WiMAX device descriptor

#### Return

Current state of the device according to its driver.

void **wimax\_dev\_init**(struct *wimax\_dev* \* *wimax\_dev*)  
initialize a newly allocated instance

#### Parameters

**struct wimax\_dev \* wimax\_dev** WiMAX device descriptor to initialize.

## Description

Initializes fields of a freshly allocated **wimax\_dev** instance. This function assumes that after allocation, the memory occupied by **wimax\_dev** was zeroed.

```
int wimax_dev_add(struct wimax_dev * wimax_dev, struct net_device * net_dev)
    Register a new WiMAX device
```

## Parameters

**struct wimax\_dev \* wimax\_dev** WiMAX device descriptor (as embedded in your **net\_dev**'s priv data). You must have called `wimax_dev_init()` on it before.

**struct net\_device \* net\_dev** net device the **wimax\_dev** is associated with. The function expects `SET_NETDEV_DEV()` and `register_netdev()` were already called on it.

## Description

Registers the new WiMAX device, sets up the user-kernel control interface (generic netlink) and common WiMAX infrastructure.

Note that the parts that will allow interaction with user space are setup at the very end, when the rest is in place, as once that happens, the driver might get user space control requests via netlink or from debugfs that might translate into calls into `wimax_dev->op_*()`.

```
void wimax_dev_rm(struct wimax_dev * wimax_dev)
    Unregister an existing WiMAX device
```

## Parameters

**struct wimax\_dev \* wimax\_dev** WiMAX device descriptor

## Description

Unregisters a WiMAX device previously registered for use with `wimax_add_rm()`.

IMPORTANT! Must call before calling `unregister_netdev()`.

After this function returns, you will not get any more user space control requests (via netlink or debugfs) and thus to `wimax_dev->ops`.

Reentrancy control is ensured by setting the state to `__WIMAX_ST QUIESCING`. `rkill` operations coming through `wimax_*rkill*()` will be stopped by the quiescing state; `ops` coming from the `rkill` subsystem will be stopped by the support being removed by `wimax_rkill_rm()`.

**struct wimax\_dev**  
Generic WiMAX device

## Definition

```
struct wimax_dev {
    struct net_device * net_dev;
    struct list_head id_table_node;
    struct mutex mutex;
    struct mutex mutex_reset;
    enum wimax_st state;
    int (* op_msg_from_user) (struct wimax_dev *wimax_dev, const char *, const void *, size_t,
    ↪const struct genl_info *info);
    int (* op_rkill_sw_toggle) (struct wimax_dev *wimax_dev, enum wimax_rf_state);
    int (* op_reset) (struct wimax_dev *wimax_dev);
    struct rkill * rkill;
    unsigned int rf_hw;
    unsigned int rf_sw;
    char name;
    struct dentry * debugfs_dentry;
};
```

## Members

**net\_dev** [fill] Pointer to the *struct net\_device* this WiMAX device implements.

**id\_table\_node** [private] link to the list of wimax devices kept by id-table.c. Protected by it's own spinlock.

**mutex** [private] Serializes all concurrent access and execution of operations.

**mutex\_reset** [private] Serializes reset operations. Needs to be a different mutex because as part of the reset operation, the driver has to call back into the stack to do things such as state change, that require `wimax_dev->mutex`.

**state** [private] Current state of the WiMAX device.

**op\_msg\_from\_user** [fill] Driver-specific operation to handle a raw message from user space to the driver. The driver can send messages to user space using `wimax_msg_to_user()`.

**op\_rfkill\_sw\_toggle** [fill] Driver-specific operation to act on userspace (or any other agent) requesting the WiMAX device to change the RF Kill software switch (`WIMAX_RF_ON` or `WIMAX_RF_OFF`). If such hardware support is not present, it is assumed the radio cannot be switched off and it is always on (and the stack will error out when trying to switch it off). In such case, this function pointer can be left as `NULL`.

**op\_reset** [fill] Driver specific operation to reset the device. This operation should always attempt first a warm reset that does not disconnect the device from the bus and return 0. If that fails, it should resort to some sort of cold or bus reset (even if it implies a bus disconnection and device disappearance). In that case, `-ENODEV` should be returned to indicate the device is gone. This operation has to be synchronous, and return only when the reset is complete. In case of having had to resort to bus/cold reset implying a device disconnection, the call is allowed to return immediately.

**rfkill** [private] integration into the RF-Kill infrastructure.

**rf\_hw** [private] State of the hardware radio switch (OFF/ON)

**rf\_sw** [private] State of the software radio switch (OFF/ON)

**name** [fill] A way to identify this device. We need to register a name with many subsystems (rfkill, workqueue creation, etc). We can't use the network device name as that might change and in some instances we don't know it yet (until we don't call `register_netdev()`). So we generate an unique one using the driver name and device bus id, place it here and use it across the board. Recommended naming: `DRIVERNAME-BUSNAME:BUSID` (`dev->bus->name`, `dev->bus_id`).

**debugfs\_dentry** [private] Used to hook up a debugfs entry. This shows up in the debugfs root as `wimax:DEVICENAME`.

## NOTE

**wimax\_dev->mutex is NOT locked when this op is being called;** however, `wimax_dev->mutex_reset` IS locked to ensure serialization of calls to `wimax_reset()`. See `wimax_reset()`'s documentation.

## Description

This structure defines a common interface to access all WiMAX devices from different vendors and provides a common API as well as a free-form device-specific messaging channel.

## Usage:

1. Embed a *struct wimax\_dev* at the beginning the network device structure so that `net_dev_priv()` points to it.
2. `memset()` it to zero
3. Initialize with `wimax_dev_init()`. This will leave the WiMAX device in the `__WIMAX_ST_NULL` state.
4. Fill all the fields marked with [fill]; once called `wimax_dev_add()`, those fields CANNOT be modified.

5. Call `wimax_dev_add()` after registering the network device. This will leave the WiMAX device in the `WIMAX_ST_DOWN` state. Protect the driver's `net_device->:c:func:open()` against succeeding if the wimax device state is lower than `WIMAX_ST_DOWN`.
6. Select when the device is going to be turned on/initialized; for example, it could be initialized on 'ifconfig up' (when the netdev op 'open()' is called on the driver).

When the device is initialized (at *ifconfig up* time, or right after calling `wimax_dev_add()` from `_probe()`, make sure the following steps are taken

1. Move the device to `WIMAX_ST_UNINITIALIZED`. This is needed so some API calls that shouldn't work until the device is ready can be blocked.
2. Initialize the device. Make sure to turn the SW radio switch off and move the device to state `WIMAX_ST_RADIO_OFF` when done. When just initialized, a device should be left in `RADIO OFF` state until user space devices to turn it on.
3. Query the device for the state of the hardware rfkill switch and call `wimax_rfkill_report_hw()` and `wimax_rfkill_report_sw()` as needed. See below.

`wimax_dev_rm()` undoes before unregistering the network device. Once `wimax_dev_add()` is called, the driver can get called on the `wimax_dev->op_*` function pointers

#### CONCURRENCY:

The stack provides a mutex for each device that will disallow API calls happening concurrently; thus, op calls into the driver through the `wimax_dev->op*()` function pointers will always be serialized and *never* concurrent.

For locking, take `wimax_dev->mutex` is taken; (most) operations in the API have to check for `wimax_dev_is_ready()` to return 0 before continuing (this is done internally).

#### REFERENCE COUNTING:

The WiMAX device is reference counted by the associated network device. The only operation that can be used to reference the device is `wimax_dev_get_by_genl_info()`, and the reference it acquires has to be released with `dev_put(wimax_dev->net_dev)`.

#### RFKILL:

At startup, both HW and SW radio switchess are assumed to be off.

At initialization time [after calling `wimax_dev_add()`], have the driver query the device for the status of the software and hardware RF kill switches and call `wimax_report_rfkill_hw()` and `wimax_rfkill_report_sw()` to indicate their state. If any is missing, just call it to indicate it is ON (radio always on).

Whenever the driver detects a change in the state of the RF kill switches, it should call `wimax_report_rfkill_hw()` or `wimax_report_rfkill_sw()` to report it to the stack.

#### enum `wimax_st`

The different states of a WiMAX device

#### Constants

**\_\_WIMAX\_ST\_NULL** The device structure has been allocated and zeroed, but still `wimax_dev_add()` hasn't been called. There is no state.

**WIMAX\_ST\_DOWN** The device has been registered with the WiMAX and networking stacks, but it is not initialized (normally that is done with 'ifconfig DEV up' [or equivalent], which can upload firmware and enable communications with the device). In this state, the device is powered down and using as less power as possible. This state is the default after a call to `wimax_dev_add()`. It is ok to have drivers move directly to `WIMAX_ST_UNINITIALIZED` or `WIMAX_ST_RADIO_OFF` in `_probe()` after the call to `wimax_dev_add()`. It is recommended that the driver leaves this state when calling 'ifconfig DEV up' and enters it back on 'ifconfig DEV down'.

**\_\_WIMAX\_ST\_QUIESCING** The device is being torn down, so no API operations are allowed to proceed except the ones needed to complete the device clean up process.



**WIMAX\_ST\_UNINITIALIZED** [optional] Communication with the device is setup, but the device still requires some configuration before being operational. Some WiMAX API calls might work.

**WIMAX\_ST\_RADIO\_OFF** The device is fully up; radio is off (wether by hardware or software switches). It is recommended to always leave the device in this state after initialization.

**WIMAX\_ST\_READY** The device is fully up and radio is on.

**WIMAX\_ST\_SCANNING** [optional] The device has been instructed to scan. In this state, the device cannot be actively connected to a network.

**WIMAX\_ST\_CONNECTING** The device is connecting to a network. This state exists because in some devices, the connect process can include a number of negotiations between user space, kernel space and the device. User space needs to know what the device is doing. If the connect sequence in a device is atomic and fast, the device can transition directly to **CONNECTED**

**WIMAX\_ST\_CONNECTED** The device is connected to a network.

**\_\_WIMAX\_ST\_INVALID** This is an invalid state used to mark the maximum numeric value of states.

### Description

Transitions from one state to another one are atomic and can only be caused in kernel space with [wimax\\_state\\_change\(\)](#). To read the state, use [wimax\\_state\\_get\(\)](#).

States starting with `__` are internal and shall not be used or referred to by drivers or userspace. They look ugly, but that's the point – if any use is made non-internal to the stack, it is easier to catch on review.

All API operations [with well defined exceptions] will take the device mutex before starting and then check the state. If the state is `__WIMAX_ST_NULL`, `WIMAX_ST_DOWN`, `WIMAX_ST_UNINITIALIZED` or `__WIMAX_ST_QUIESCING`, it will drop the lock and quit with `-EINVAL`, `-ENOMEDIUM`, `-ENOTCONN` or `-ESHUTDOWN`.

The order of the definitions is important, so we can do numerical comparisons (eg: `< WIMAX_ST_RADIO_OFF` means the device is not ready to operate).

## Network device support

### Driver Support

void **dev\_add\_pack**(struct packet\_type \* pt)  
add packet handler

#### Parameters

**struct packet\_type \* pt** packet type declaration

#### Description

Add a protocol handler to the networking stack. The passed `packet_type` is linked into kernel lists and may not be freed until it has been removed from the kernel lists.

This call does not sleep therefore it can not guarantee all CPU's that are in middle of receiving packets will see the new packet type (until the next received packet).

void **\_\_dev\_remove\_pack**(struct packet\_type \* pt)  
remove packet handler

#### Parameters

**struct packet\_type \* pt** packet type declaration

#### Description

Remove a protocol handler that was previously added to the kernel protocol handlers by [dev\\_add\\_pack\(\)](#). The passed `packet_type` is removed from the kernel lists and can be freed or reused once this function returns.

**The packet type might still be in use by receivers** and must not be freed until after all the CPU's have gone through a quiescent state.

void **dev\_remove\_pack**(struct packet\_type \* *pt*)  
remove packet handler

#### Parameters

struct packet\_type \* **pt** packet type declaration

#### Description

Remove a protocol handler that was previously added to the kernel protocol handlers by [dev\\_add\\_pack\(\)](#). The passed packet\_type is removed from the kernel lists and can be freed or reused once this function returns.

This call sleeps to guarantee that no CPU is looking at the packet type after return.

void **dev\_add\_offload**(struct packet\_offload \* *po*)  
register offload handlers

#### Parameters

struct packet\_offload \* **po** protocol offload declaration

#### Description

Add protocol offload handlers to the networking stack. The passed proto\_offload is linked into kernel lists and may not be freed until it has been removed from the kernel lists.

This call does not sleep therefore it can not guarantee all CPU's that are in middle of receiving packets will see the new offload handlers (until the next received packet).

void **dev\_remove\_offload**(struct packet\_offload \* *po*)  
remove packet offload handler

#### Parameters

struct packet\_offload \* **po** packet offload declaration

#### Description

Remove a packet offload handler that was previously added to the kernel offload handlers by [dev\\_add\\_offload\(\)](#). The passed offload\_type is removed from the kernel lists and can be freed or reused once this function returns.

This call sleeps to guarantee that no CPU is looking at the packet type after return.

int **netdev\_boot\_setup\_check**(struct [net\\_device](#) \* *dev*)  
check boot time settings

#### Parameters

struct net\_device \* **dev** the netdevice

#### Description

Check boot time settings for the device. The found settings are set for the device to be used later in the device probing. Returns 0 if no settings found, 1 if they are.

int **dev\_get\_iflink**(const struct [net\\_device](#) \* *dev*)  
get 'iflink' value of a interface

#### Parameters

const struct net\_device \* **dev** targeted interface

#### Description

Indicates the ifindex the interface is linked to. Physical interfaces have the same 'ifindex' and 'iflink' values.

int **dev\_fill\_metadata\_dst**(struct *net\_device* \* *dev*, struct *sk\_buff* \* *skb*)  
Retrieve tunnel egress information.

#### Parameters

**struct net\_device \* dev** targeted interface

**struct sk\_buff \* skb** The packet.

#### Description

For better visibility of tunnel traffic OVS needs to retrieve egress tunnel information for a packet. Following API allows user to get this info.

struct *net\_device* \* **\_\_dev\_get\_by\_name**(struct net \* *net*, const char \* *name*)  
find a device by its name

#### Parameters

**struct net \* net** the applicable net namespace

**const char \* name** name to find

#### Description

Find an interface by name. Must be called under RTNL semaphore or **dev\_base\_lock**. If the name is found a pointer to the device is returned. If the name is not found then NULL is returned. The reference counters are not incremented so the caller must be careful with locks.

struct *net\_device* \* **dev\_get\_by\_name\_rcu**(struct net \* *net*, const char \* *name*)  
find a device by its name

#### Parameters

**struct net \* net** the applicable net namespace

**const char \* name** name to find

#### Description

Find an interface by name. If the name is found a pointer to the device is returned. If the name is not found then NULL is returned. The reference counters are not incremented so the caller must be careful with locks. The caller must hold RCU lock.

struct *net\_device* \* **dev\_get\_by\_name**(struct net \* *net*, const char \* *name*)  
find a device by its name

#### Parameters

**struct net \* net** the applicable net namespace

**const char \* name** name to find

#### Description

Find an interface by name. This can be called from any context and does its own locking. The returned handle has the usage count incremented and the caller must use *dev\_put()* to release it when it is no longer needed. NULL is returned if no matching device is found.

struct *net\_device* \* **\_\_dev\_get\_by\_index**(struct net \* *net*, int *ifindex*)  
find a device by its ifindex

#### Parameters

**struct net \* net** the applicable net namespace

**int ifindex** index of device

#### Description

Search for an interface by index. Returns NULL if the device is not found or a pointer to the device. The device has not had its reference counter increased so the caller must be careful about locking. The caller must hold either the RTNL semaphore or **dev\_base\_lock**.

struct *net\_device* \* **dev\_get\_by\_index\_rcu**(struct net \* *net*, int *ifindex*)  
find a device by its ifindex

#### Parameters

**struct net \* net** the applicable net namespace

**int ifindex** index of device

#### Description

Search for an interface by index. Returns NULL if the device is not found or a pointer to the device. The device has not had its reference counter increased so the caller must be careful about locking. The caller must hold RCU lock.

struct *net\_device* \* **dev\_get\_by\_index**(struct net \* *net*, int *ifindex*)  
find a device by its ifindex

#### Parameters

**struct net \* net** the applicable net namespace

**int ifindex** index of device

#### Description

Search for an interface by index. Returns NULL if the device is not found or a pointer to the device. The device returned has had a reference added and the pointer is safe until the user calls `dev_put` to indicate they have finished with it.

struct *net\_device* \* **dev\_get\_by\_napi\_id**(unsigned int *napi\_id*)  
find a device by napi\_id

#### Parameters

**unsigned int napi\_id** ID of the NAPI struct

#### Description

Search for an interface by NAPI ID. Returns NULL if the device is not found or a pointer to the device. The device has not had its reference counter increased so the caller must be careful about locking. The caller must hold RCU lock.

struct *net\_device* \* **dev\_getbyhwaddr\_rcu**(struct net \* *net*, unsigned short *type*, const char \* *ha*)  
find a device by its hardware address

#### Parameters

**struct net \* net** the applicable net namespace

**unsigned short type** media type of device

**const char \* ha** hardware address

#### Description

Search for an interface by MAC address. Returns NULL if the device is not found or a pointer to the device. The caller must hold RCU or RTNL. The returned device has not had its ref count increased and the caller must therefore be careful about locking

struct *net\_device* \* **\_\_dev\_get\_by\_flags**(struct net \* *net*, unsigned short *if\_flags*, unsigned short *mask*)  
find any device with given flags

#### Parameters

**struct net \* net** the applicable net namespace

**unsigned short if\_flags** IFF\_\* values

**unsigned short mask** bitmask of bits in if\_flags to check

#### Description

Search for any interface with the given flags. Returns NULL if a device is not found or a pointer to the device. Must be called inside `rtnl_lock()`, and result refcount is unchanged.

bool **dev\_valid\_name**(const char \* *name*)  
check if name is okay for network device

#### Parameters

const char \* **name** name string

#### Description

Network device names need to be valid file names to to allow sysfs to work. We also disallow any kind of whitespace.

int **dev\_alloc\_name**(struct *net\_device* \* *dev*, const char \* *name*)  
allocate a name for a device

#### Parameters

struct *net\_device* \* **dev** device

const char \* **name** name format string

#### Description

Passed a format string - eg "lt%d" it will try and find a suitable id. It scans list of devices to build up a free map, then chooses the first empty slot. The caller must hold the dev\_base or rtnl lock while allocating the name and adding the device in order to avoid duplicates. Limited to bits\_per\_byte \* page size devices (ie 32K on most platforms). Returns the number of the unit assigned or a negative errno code.

void **netdev\_features\_change**(struct *net\_device* \* *dev*)  
device changes features

#### Parameters

struct *net\_device* \* **dev** device to cause notification

#### Description

Called to indicate a device has changed features.

void **netdev\_state\_change**(struct *net\_device* \* *dev*)  
device changes state

#### Parameters

struct *net\_device* \* **dev** device to cause notification

#### Description

Called to indicate a device has changed state. This function calls the notifier chains for netdev\_chain and sends a NEWLINK message to the routing socket.

void **netdev\_notify\_peers**(struct *net\_device* \* *dev*)  
notify network peers about existence of **dev**

#### Parameters

struct *net\_device* \* **dev** network device

#### Description

Generate traffic such that interested network peers are aware of **dev**, such as by generating a gratuitous ARP. This may be used when a device wants to inform the rest of the network about some sort of reconfiguration such as a failover event or virtual machine migration.

int **dev\_open**(struct *net\_device* \* *dev*)  
prepare an interface for use.

#### Parameters

**struct net\_device \* dev** device to open

### Description

Takes a device from down to up state. The device's private open function is invoked and then the multicast lists are loaded. Finally the device is moved into the up state and a NETDEV\_UP message is sent to the netdev notifier chain.

Calling this function on an active interface is a nop. On a failure a negative errno code is returned.

void **dev\_close**(struct [net\\_device](#) \* dev)  
shutdown an interface.

### Parameters

**struct net\_device \* dev** device to shutdown

### Description

This function moves an active device into down state. A NETDEV\_GOING\_DOWN is sent to the netdev notifier chain. The device is then deactivated and finally a NETDEV\_DOWN is sent to the notifier chain.

void **dev\_disable\_lro**(struct [net\\_device](#) \* dev)  
disable Large Receive Offload on a device

### Parameters

**struct net\_device \* dev** device

### Description

Disable Large Receive Offload (LRO) on a net device. Must be called under RTNL. This is needed if received packets may be forwarded to another interface.

int **register\_netdevice\_notifier**(struct notifier\_block \* nb)  
register a network notifier block

### Parameters

**struct notifier\_block \* nb** notifier

### Description

Register a notifier to be called when network device events occur. The notifier passed is linked into the kernel structures and must not be reused until it has been unregistered. A negative errno code is returned on a failure.

When registered all registration and up events are replayed to the new notifier to allow device to have a race free view of the network device list.

int **unregister\_netdevice\_notifier**(struct notifier\_block \* nb)  
unregister a network notifier block

### Parameters

**struct notifier\_block \* nb** notifier

### Description

Unregister a notifier previously registered by [register\\_netdevice\\_notifier\(\)](#). The notifier is unlinked into the kernel structures and may then be reused. A negative errno code is returned on a failure.

After unregistering unregister and down device events are synthesized for all devices on the device list to the removed notifier to remove the need for special case cleanup code.

int **call\_netdevice\_notifiers**(unsigned long val, struct [net\\_device](#) \* dev)  
call all network notifier blocks

### Parameters

**unsigned long val** value passed unmodified to notifier function

**struct net\_device \* dev** net\_device pointer passed unmodified to notifier function

#### Description

Call all network notifier blocks. Parameters and return value are as for `raw_notifier_call_chain()`.

int **dev\_forward\_skb**(struct *net\_device* \* dev, struct *sk\_buff* \* skb)  
loopback an skb to another netif

#### Parameters

**struct net\_device \* dev** destination network device

**struct sk\_buff \* skb** buffer to forward

#### Description

**return values:** NET\_RX\_SUCCESS (no congestion) NET\_RX\_DROP (packet was dropped, but freed)

`dev_forward_skb` can be used for injecting an skb from the `start_xmit` function of one device into the receive queue of another device.

The receiving device may be in another namespace, so we have to clear all information in the skb that could impact namespace isolation.

int **netif\_set\_real\_num\_rx\_queues**(struct *net\_device* \* dev, unsigned int rxq)  
set actual number of RX queues used

#### Parameters

**struct net\_device \* dev** Network device

**unsigned int rxq** Actual number of RX queues

#### Description

This must be called either with the `rtnl_lock` held or before registration of the net device. Returns 0 on success, or a negative error code. If called before registration, it always succeeds.

int **netif\_get\_num\_default\_rss\_queues**(void)  
default number of RSS queues

#### Parameters

**void** no arguments

#### Description

This routine should set an upper limit on the number of RSS queues used by default by multiqueue devices.

void **netif\_device\_detach**(struct *net\_device* \* dev)  
mark device as removed

#### Parameters

**struct net\_device \* dev** network device

#### Description

Mark device as removed from system and therefore no longer available.

void **netif\_device\_attach**(struct *net\_device* \* dev)  
mark device as attached

#### Parameters

**struct net\_device \* dev** network device

## Description

Mark device as attached from system and restart if needed.

struct *sk\_buff* \* **skb\_mac\_gso\_segment**(struct *sk\_buff* \* *skb*, netdev\_features\_t *features*)  
mac layer segmentation handler.

## Parameters

struct *sk\_buff* \* **skb** buffer to segment

netdev\_features\_t **features** features for the output path (see dev->features)

struct *sk\_buff* \* **\_\_skb\_gso\_segment**(struct *sk\_buff* \* *skb*, netdev\_features\_t *features*, bool *tx\_path*)  
Perform segmentation on skb.

## Parameters

struct *sk\_buff* \* **skb** buffer to segment

netdev\_features\_t **features** features for the output path (see dev->features)

bool **tx\_path** whether it is called in TX path

## Description

This function segments the given skb and returns a list of segments.

It may return NULL if the skb requires no segmentation. This is only possible when GSO is used for verifying header integrity.

Segmentation preserves SKB\_SGO\_CB\_OFFSET bytes of previous skb cb.

int **dev\_loopback\_xmit**(struct net \* *net*, struct *sock* \* *sk*, struct *sk\_buff* \* *skb*)  
loop back **skb**

## Parameters

struct net \* **net** network namespace this loopback is happening in

struct sock \* **sk** sk needed to be a netfilter okfn

struct sk\_buff \* **skb** buffer to transmit

bool **rps\_may\_expire\_flow**(struct *net\_device* \* *dev*, u16 *rxq\_index*, u32 *flow\_id*, u16 *filter\_id*)  
check whether an RFS hardware filter may be removed

## Parameters

struct net\_device \* **dev** Device on which the filter was set

u16 **rxq\_index** RX queue index

u32 **flow\_id** Flow ID passed to ndo\_rx\_flow\_steer()

u16 **filter\_id** Filter ID returned by ndo\_rx\_flow\_steer()

## Description

Drivers that implement ndo\_rx\_flow\_steer() should periodically call this function for each installed filter and remove the filters for which it returns true.

int **netif\_rx**(struct *sk\_buff* \* *skb*)  
post buffer to the network code

## Parameters

struct sk\_buff \* **skb** buffer to post

## Description



This function receives a packet from a device driver and queues it for the upper (protocol) levels to process. It always succeeds. The buffer may be dropped during processing for congestion control or by the protocol layers.

return values: NET\_RX\_SUCCESS (no congestion) NET\_RX\_DROP (packet was dropped)

bool **netdev\_is\_rx\_handler\_busy**(struct *net\_device* \* dev)  
check if receive handler is registered

#### Parameters

**struct net\_device \* dev** device to check

#### Description

Check if a receive handler is already registered for a given device. Return true if there one.

The caller must hold the rtnl\_mutex.

int **netdev\_rx\_handler\_register**(struct *net\_device* \* dev, rx\_handler\_func\_t \* rx\_handler, void \* rx\_handler\_data)  
register receive handler

#### Parameters

**struct net\_device \* dev** device to register a handler for

**rx\_handler\_func\_t \* rx\_handler** receive handler to register

**void \* rx\_handler\_data** data pointer that is used by rx handler

#### Description

Register a receive handler for a device. This handler will then be called from `__netif_receive_skb`. A negative errno code is returned on a failure.

The caller must hold the rtnl\_mutex.

For a general description of rx\_handler, see enum rx\_handler\_result.

void **netdev\_rx\_handler\_unregister**(struct *net\_device* \* dev)  
unregister receive handler

#### Parameters

**struct net\_device \* dev** device to unregister a handler from

#### Description

Unregister a receive handler from a device.

The caller must hold the rtnl\_mutex.

int **netif\_receive\_skb**(struct *sk\_buff* \* skb)  
process receive buffer from network

#### Parameters

**struct sk\_buff \* skb** buffer to process

#### Description

`netif_receive_skb()` is the main receive data processing function. It always succeeds. The buffer may be dropped during processing for congestion control or by the protocol layers.

This function may only be called from softirq context and interrupts should be enabled.

Return values (usually ignored): NET\_RX\_SUCCESS: no congestion NET\_RX\_DROP: packet was dropped

void **\_\_napi\_schedule**(struct napi\_struct \* n)  
schedule for receive

#### Parameters

**struct napi\_struct \* n** entry to schedule

### Description

The entry's receive function will be scheduled to run. Consider using `__napi_schedule_irqoff()` if hard irqs are masked.

bool **napi\_schedule\_prep**(struct napi\_struct \* n)  
check if napi can be scheduled

### Parameters

**struct napi\_struct \* n** napi context

### Description

Test if NAPI routine is already running, and if not mark it as running. This is used as a condition variable insure only one NAPI poll instance runs. We also make sure there is no pending NAPI disable.

void **\_\_napi\_schedule\_irqoff**(struct napi\_struct \* n)  
schedule for receive

### Parameters

**struct napi\_struct \* n** entry to schedule

### Description

Variant of `__napi_schedule()` assuming hard irqs are masked

bool **netdev\_has\_upper\_dev**(struct *net\_device* \* dev, struct *net\_device* \* upper\_dev)  
Check if device is linked to an upper device

### Parameters

**struct net\_device \* dev** device

**struct net\_device \* upper\_dev** upper device to check

### Description

Find out if a device is linked to specified upper device and return true in case it is. Note that this checks only immediate upper device, not through a complete stack of devices. The caller must hold the RTNL lock.

bool **netdev\_has\_upper\_dev\_all\_rcu**(struct *net\_device* \* dev, struct *net\_device* \* upper\_dev)  
Check if device is linked to an upper device

### Parameters

**struct net\_device \* dev** device

**struct net\_device \* upper\_dev** upper device to check

### Description

Find out if a device is linked to specified upper device and return true in case it is. Note that this checks the entire upper device chain. The caller must hold rcu lock.

bool **netdev\_has\_any\_upper\_dev**(struct *net\_device* \* dev)  
Check if device is linked to some device

### Parameters

**struct net\_device \* dev** device

### Description

Find out if a device is linked to an upper device and return true in case it is. The caller must hold the RTNL lock.

struct *net\_device* \* **netdev\_master\_upper\_dev\_get**(struct *net\_device* \* dev)  
Get master upper device

**Parameters**

**struct net\_device \* dev** device

**Description**

Find a master upper device and return pointer to it or NULL in case it's not there. The caller must hold the RTNL lock.

**struct net\_device \* netdev\_upper\_get\_next\_dev\_rcu**(**struct net\_device \* dev**, **struct list\_head \*\* iter**)

Get the next dev from upper list

**Parameters**

**struct net\_device \* dev** device

**struct list\_head \*\* iter** list\_head \*\* of the current position

**Description**

Gets the next device from the dev's upper list, starting from iter position. The caller must hold RCU read lock.

**void \* netdev\_lower\_get\_next\_private**(**struct net\_device \* dev**, **struct list\_head \*\* iter**)

Get the next ->private from the lower neighbour list

**Parameters**

**struct net\_device \* dev** device

**struct list\_head \*\* iter** list\_head \*\* of the current position

**Description**

Gets the next netdev\_adjacent->private from the dev's lower neighbour list, starting from iter position. The caller must either hold the RTNL lock or its own locking that guarantees that the neighbour lower list will remain unchanged.

**void \* netdev\_lower\_get\_next\_private\_rcu**(**struct net\_device \* dev**, **struct list\_head \*\* iter**)

Get the next ->private from the lower neighbour list, RCU variant

**Parameters**

**struct net\_device \* dev** device

**struct list\_head \*\* iter** list\_head \*\* of the current position

**Description**

Gets the next netdev\_adjacent->private from the dev's lower neighbour list, starting from iter position. The caller must hold RCU read lock.

**void \* netdev\_lower\_get\_next**(**struct net\_device \* dev**, **struct list\_head \*\* iter**)

Get the next device from the lower neighbour list

**Parameters**

**struct net\_device \* dev** device

**struct list\_head \*\* iter** list\_head \*\* of the current position

**Description**

Gets the next netdev\_adjacent from the dev's lower neighbour list, starting from iter position. The caller must hold RTNL lock or its own locking that guarantees that the neighbour lower list will remain unchanged.

**void \* netdev\_lower\_get\_first\_private\_rcu**(**struct net\_device \* dev**)

Get the first ->private from the lower neighbour list, RCU variant

**Parameters**

**struct net\_device \* dev** device

### Description

Gets the first `netdev_adjacent->private` from the dev's lower neighbour list. The caller must hold RCU read lock.

`struct net_device * netdev_master_upper_dev_get_rcu(struct net_device * dev)`  
Get master upper device

### Parameters

**struct net\_device \* dev** device

### Description

Find a master upper device and return pointer to it or NULL in case it's not there. The caller must hold the RCU read lock.

`int netdev_upper_dev_link(struct net_device * dev, struct net_device * upper_dev)`  
Add a link to the upper device

### Parameters

**struct net\_device \* dev** device

**struct net\_device \* upper\_dev** new upper device

### Description

Adds a link to device which is upper to this one. The caller must hold the RTNL lock. On a failure a negative errno code is returned. On success the reference counts are adjusted and the function returns zero.

`int netdev_master_upper_dev_link(struct net_device * dev, struct net_device * upper_dev, void * upper_priv, void * upper_info)`  
Add a master link to the upper device

### Parameters

**struct net\_device \* dev** device

**struct net\_device \* upper\_dev** new upper device

**void \* upper\_priv** upper device private

**void \* upper\_info** upper info to be passed down via notifier

### Description

Adds a link to device which is upper to this one. In this case, only one master upper device can be linked, although other non-master devices might be linked as well. The caller must hold the RTNL lock. On a failure a negative errno code is returned. On success the reference counts are adjusted and the function returns zero.

`void netdev_upper_dev_unlink(struct net_device * dev, struct net_device * upper_dev)`  
Removes a link to upper device

### Parameters

**struct net\_device \* dev** device

**struct net\_device \* upper\_dev** new upper device

### Description

Removes a link to device which is upper to this one. The caller must hold the RTNL lock.

`void netdev_bonding_info_change(struct net_device * dev, struct netdev_bonding_info * bonding_info)`  
Dispatch event about slave change

### Parameters

```
struct net_device * dev device
```

```
struct netdev_bonding_info * bonding_info info to dispatch
```

### Description

Send NETDEV\_BONDING\_INFO to netdev notifiers with info. The caller must hold the RTNL lock.

```
void netdev_lower_state_changed(struct net_device * lower_dev, void * lower_state_info)
```

Dispatch event about lower device state change

### Parameters

```
struct net_device * lower_dev device
```

```
void * lower_state_info state to dispatch
```

### Description

Send NETDEV\_CHANGELOWERSTATE to netdev notifiers with info. The caller must hold the RTNL lock.

```
int dev_set_promiscuity(struct net_device * dev, int inc)
```

update promiscuity count on a device

### Parameters

```
struct net_device * dev device
```

```
int inc modifier
```

### Description

Add or remove promiscuity from a device. While the count in the device remains above zero the interface remains promiscuous. Once it hits zero the device reverts back to normal filtering operation. A negative inc value is used to drop promiscuity on the device. Return 0 if successful or a negative errno code on error.

```
int dev_set_allmulti(struct net_device * dev, int inc)
```

update allmulti count on a device

### Parameters

```
struct net_device * dev device
```

```
int inc modifier
```

### Description

Add or remove reception of all multicast frames to a device. While the count in the device remains above zero the interface remains listening to all interfaces. Once it hits zero the device reverts back to normal filtering operation. A negative inc value is used to drop the counter when releasing a resource needing all multicasts. Return 0 if successful or a negative errno code on error.

```
unsigned int dev_get_flags(const struct net_device * dev)
```

get flags reported to userspace

### Parameters

```
const struct net_device * dev device
```

### Description

Get the combination of flag bits exported through APIs to userspace.

```
int dev_change_flags(struct net_device * dev, unsigned int flags)
```

change device settings

### Parameters

```
struct net_device * dev device
```

```
unsigned int flags device state flags
```

### Description

Change settings on device based state flags. The flags are in the userspace exported format.

int **dev\_set\_mtu**(struct *net\_device* \* *dev*, int *new\_mtu*)  
Change maximum transfer unit

### Parameters

**struct net\_device \* dev** device  
**int new\_mtu** new transfer unit

### Description

Change the maximum transfer size of the network device.

void **dev\_set\_group**(struct *net\_device* \* *dev*, int *new\_group*)  
Change group this device belongs to

### Parameters

**struct net\_device \* dev** device  
**int new\_group** group this device should belong to  
int **dev\_set\_mac\_address**(struct *net\_device* \* *dev*, struct sockaddr \* *sa*)  
Change Media Access Control Address

### Parameters

**struct net\_device \* dev** device  
**struct sockaddr \* sa** new address

### Description

Change the hardware (MAC) address of the device

int **dev\_change\_carrier**(struct *net\_device* \* *dev*, bool *new\_carrier*)  
Change device carrier

### Parameters

**struct net\_device \* dev** device  
**bool new\_carrier** new value

### Description

Change device carrier

int **dev\_get\_phys\_port\_id**(struct *net\_device* \* *dev*, struct netdev\_phys\_item\_id \* *ppid*)  
Get device physical port ID

### Parameters

**struct net\_device \* dev** device  
**struct netdev\_phys\_item\_id \* ppid** port ID

### Description

Get device physical port ID

int **dev\_get\_phys\_port\_name**(struct *net\_device* \* *dev*, char \* *name*, size\_t *len*)  
Get device physical port name

### Parameters

**struct net\_device \* dev** device  
**char \* name** port name  
**size\_t len** limit of bytes to copy to name

**Description**

Get device physical port name

int **dev\_change\_proto\_down**(struct *net\_device* \* dev, bool *proto\_down*)  
update protocol port state information

**Parameters**

**struct net\_device \* dev** device

**bool proto\_down** new value

**Description**

This info can be used by switch drivers to set the phys state of the port.

void **netdev\_update\_features**(struct *net\_device* \* dev)  
recalculate device features

**Parameters**

**struct net\_device \* dev** the device to check

**Description**

Recalculate dev->features set and send notifications if it has changed. Should be called after driver or hardware dependent conditions might have changed that influence the features.

void **netdev\_change\_features**(struct *net\_device* \* dev)  
recalculate device features

**Parameters**

**struct net\_device \* dev** the device to check

**Description**

Recalculate dev->features set and send notifications even if they have not changed. Should be called instead of *netdev\_update\_features()* if also dev->vlan\_features might have changed to allow the changes to be propagated to stacked VLAN devices.

void **netif\_stacked\_transfer\_operstate**(const struct *net\_device* \* rootdev, struct *net\_device* \* dev)  
transfer operstate

**Parameters**

**const struct net\_device \* rootdev** the root or lower level device to transfer state from

**struct net\_device \* dev** the device to transfer operstate to

**Description**

Transfer operational state from root to device. This is normally called when a stacking relationship exists between the root device and the device(a leaf device).

int **register\_netdevice**(struct *net\_device* \* dev)  
register a network device

**Parameters**

**struct net\_device \* dev** device to register

**Description**

Take a completed network device structure and add it to the kernel interfaces. A NET\_DEV\_REGISTER message is sent to the netdev notifier chain. 0 is returned on success. A negative errno code is returned on a failure to set up the device, or if the name is a duplicate.

Callers must hold the rtnl semaphore. You may want *register\_netdev()* instead of this.

BUGS: The locking appears insufficient to guarantee two parallel registers will not get the same name.

int **init\_dummy\_netdev**(struct *net\_device* \* dev)  
init a dummy network device for NAPI

### Parameters

**struct net\_device \* dev** device to init

### Description

This takes a network device structure and initialize the minimum amount of fields so it can be used to schedule NAPI polls without registering a full blown interface. This is to be used by drivers that need to tie several hardware interfaces to a single NAPI poll scheduler due to HW limitations.

int **register\_netdev**(struct *net\_device* \* dev)  
register a network device

### Parameters

**struct net\_device \* dev** device to register

### Description

Take a completed network device structure and add it to the kernel interfaces. A NET\_DEV\_REGISTER message is sent to the netdev notifier chain. 0 is returned on success. A negative errno code is returned on a failure to set up the device, or if the name is a duplicate.

This is a wrapper around register\_netdevice that takes the rtnl semaphore and expands the device name if you passed a format string to alloc\_netdev.

struct rtnl\_link\_stats64 \* **dev\_get\_stats**(struct *net\_device* \* dev, struct rtnl\_link\_stats64 \* storage)  
get network device statistics

### Parameters

**struct net\_device \* dev** device to get statistics from

**struct rtnl\_link\_stats64 \* storage** place to store stats

### Description

Get network statistics from device. Return **storage**. The device driver may provide its own method by setting dev->netdev\_ops->get\_stats64 or dev->netdev\_ops->get\_stats; otherwise the internal statistics structure is used.

struct *net\_device* \* **alloc\_netdev\_mqs**(int sizeof\_priv, const char \* name, unsigned char name\_assign\_type, void (\*setup) (struct *net\_device* \*, unsigned int txqs, unsigned int rxqs)  
allocate network device

### Parameters

int **sizeof\_priv** size of private data to allocate space for

const char \* **name** device name format string

unsigned char **name\_assign\_type** origin of device name

void (\*)(struct *net\_device* \*) **setup** callback to initialize device

unsigned int **txqs** the number of TX subqueues to allocate

unsigned int **rxqs** the number of RX subqueues to allocate

### Description

Allocates a struct net\_device with private data area for driver use and performs basic initialization. Also allocates subqueue structs for each queue on the device.



void **free\_netdev**(struct *net\_device* \* *dev*)  
free network device

#### Parameters

**struct net\_device \* dev** device

#### Description

This function does the last stage of destroying an allocated device interface. The reference to the device object is released. If this is the last reference then it will be freed. Must be called in process context.

void **synchronize\_net**(void)  
Synchronize with packet receive processing

#### Parameters

**void** no arguments

#### Description

Wait for packets currently being received to be done. Does not block later packets from starting.

void **unregister\_netdevice\_queue**(struct *net\_device* \* *dev*, struct list\_head \* *head*)  
remove device from the kernel

#### Parameters

**struct net\_device \* dev** device

**struct list\_head \* head** list

#### Description

This function shuts down a device interface and removes it from the kernel tables. If head not NULL, device is queued to be unregistered later.

Callers must hold the rtnl semaphore. You may want *unregister\_netdev()* instead of this.

void **unregister\_netdevice\_many**(struct list\_head \* *head*)  
unregister many devices

#### Parameters

**struct list\_head \* head** list of devices

#### Note

**As most callers use a stack allocated list\_head**, we force a *list\_del()* to make sure stack wont be corrupted later.

void **unregister\_netdev**(struct *net\_device* \* *dev*)  
remove device from the kernel

#### Parameters

**struct net\_device \* dev** device

#### Description

This function shuts down a device interface and removes it from the kernel tables.

This is just a wrapper for *unregister\_netdevice* that takes the rtnl semaphore. In general you want to use this and not *unregister\_netdevice*.

int **dev\_change\_net\_namespace**(struct *net\_device* \* *dev*, struct net \* *net*, const char \* *pat*)  
move device to different nethost namespace

#### Parameters

**struct net\_device \* dev** device

**struct net \* net** network namespace

**const char \* pat** If not NULL name pattern to try if the current device name is already taken in the destination network namespace.

### Description

This function shuts down a device interface and moves it to a new network namespace. On success 0 is returned, on a failure a netagive errno code is returned.

Callers must hold the rtnl semaphore.

**netdev\_features\_t netdev\_increment\_features**(**netdev\_features\_t all**, **netdev\_features\_t one**, **netdev\_features\_t mask**)  
increment feature set by one

### Parameters

**netdev\_features\_t all** current feature set

**netdev\_features\_t one** new feature set

**netdev\_features\_t mask** mask feature set

### Description

Computes a new feature set after adding a device with feature set **one** to the master device with current feature set **all**. Will not enable anything that is off in **mask**. Returns the new feature set.

**int eth\_header**(**struct sk\_buff \* skb**, **struct net\_device \* dev**, **unsigned short type**, **const void \* daddr**, **const void \* saddr**, **unsigned int len**)  
create the Ethernet header

### Parameters

**struct sk\_buff \* skb** buffer to alter

**struct net\_device \* dev** source device

**unsigned short type** Ethernet type field

**const void \* daddr** destination address (NULL leave destination address)

**const void \* saddr** source address (NULL use device source address)

**unsigned int len** packet length (<= skb->len)

### Description

Set the protocol type. For a packet of type ETH\_P\_802\_3/2 we put the length in here instead.

**u32 eth\_get\_headlen**(**void \* data**, **unsigned int len**)  
determine the length of header for an ethernet frame

### Parameters

**void \* data** pointer to start of frame

**unsigned int len** total length of frame

### Description

Make a best effort attempt to pull the length for all of the headers for a given frame in a linear buffer.

**\_\_be16 eth\_type\_trans**(**struct sk\_buff \* skb**, **struct net\_device \* dev**)  
determine the packet's protocol ID.

### Parameters

**struct sk\_buff \* skb** received socket data

**struct net\_device \* dev** receiving network device

### Description

The rule here is that we assume 802.3 if the type field is short enough to be a length. This is normal practice and works for any 'now in use' protocol.

int **eth\_header\_parse**(const struct *sk\_buff* \* *skb*, unsigned char \* *haddr*)  
extract hardware address from packet

### Parameters

const struct *sk\_buff* \* *skb* packet to extract header from

unsigned char \* *haddr* destination buffer

int **eth\_header\_cache**(const struct neighbour \* *neigh*, struct hh\_cache \* *hh*, \_\_be16 type)  
fill cache entry from neighbour

### Parameters

const struct neighbour \* *neigh* source neighbour

struct hh\_cache \* *hh* destination cache entry

\_\_be16 type Ethernet type field

### Description

Create an Ethernet header template from the neighbour.

void **eth\_header\_cache\_update**(struct hh\_cache \* *hh*, const struct *net\_device* \* *dev*, const unsigned char \* *haddr*)  
update cache entry

### Parameters

struct hh\_cache \* *hh* destination cache entry

const struct *net\_device* \* *dev* network device

const unsigned char \* *haddr* new hardware address

### Description

Called by Address Resolution module to notify changes in address.

int **eth\_prepare\_mac\_addr\_change**(struct *net\_device* \* *dev*, void \* *p*)  
prepare for mac change

### Parameters

struct *net\_device* \* *dev* network device

void \* *p* socket address

void **eth\_commit\_mac\_addr\_change**(struct *net\_device* \* *dev*, void \* *p*)  
commit mac change

### Parameters

struct *net\_device* \* *dev* network device

void \* *p* socket address

int **eth\_mac\_addr**(struct *net\_device* \* *dev*, void \* *p*)  
set new Ethernet hardware address

### Parameters

struct *net\_device* \* *dev* network device

void \* *p* socket address

### Description

Change hardware address of device.

This doesn't change hardware matching, so needs to be overridden for most real devices.

```
int eth_change_mtu(struct net_device * dev, int new_mtu)  
    set new MTU size
```

### Parameters

**struct net\_device \* dev** network device

**int new\_mtu** new Maximum Transfer Unit

### Description

Allow changing MTU size. Needs to be overridden for devices supporting jumbo frames.

```
void ether_setup(struct net_device * dev)  
    setup Ethernet network device
```

### Parameters

**struct net\_device \* dev** network device

### Description

Fill in the fields of the device structure with Ethernet-generic values.

```
struct net_device * alloc_etherdev_mqs(int sizeof_priv, unsigned int txqs, unsigned int rxqs)  
    Allocates and sets up an Ethernet device
```

### Parameters

**int sizeof\_priv** Size of additional driver-private structure to be allocated for this Ethernet device

**unsigned int txqs** The number of TX queues this device has.

**unsigned int rxqs** The number of RX queues this device has.

### Description

Fill in the fields of the device structure with Ethernet-generic values. Basically does everything except registering the device.

Constructs a new net device, complete with a private data area of size (*sizeof\_priv*). A 32-byte (not bit) alignment is enforced for this private data area.

```
void netif_carrier_on(struct net_device * dev)  
    set carrier
```

### Parameters

**struct net\_device \* dev** network device

### Description

Device has detected that carrier.

```
void netif_carrier_off(struct net_device * dev)  
    clear carrier
```

### Parameters

**struct net\_device \* dev** network device

### Description

Device has detected loss of carrier.

```
bool is_link_local_ether_addr(const u8 * addr)  
    Determine if given Ethernet address is link-local
```

### Parameters

**const u8 \* addr** Pointer to a six-byte array containing the Ethernet address

### Description

Return true if address is link local reserved addr (01:80:c2:00:00:0X) per IEEE 802.1Q 8.6.3 Frame filtering.

Please note: addr must be aligned to u16.

bool **is\_zero\_ether\_addr**(const u8 \* *addr*)  
Determine if give Ethernet address is all zeros.

### Parameters

**const u8 \* addr** Pointer to a six-byte array containing the Ethernet address

### Description

Return true if the address is all zeroes.

Please note: addr must be aligned to u16.

bool **is\_multicast\_ether\_addr**(const u8 \* *addr*)  
Determine if the Ethernet address is a multicast.

### Parameters

**const u8 \* addr** Pointer to a six-byte array containing the Ethernet address

### Description

Return true if the address is a multicast address. By definition the broadcast address is also a multicast address.

bool **is\_local\_ether\_addr**(const u8 \* *addr*)  
Determine if the Ethernet address is locally-assigned one (IEEE 802).

### Parameters

**const u8 \* addr** Pointer to a six-byte array containing the Ethernet address

### Description

Return true if the address is a local address.

bool **is\_broadcast\_ether\_addr**(const u8 \* *addr*)  
Determine if the Ethernet address is broadcast

### Parameters

**const u8 \* addr** Pointer to a six-byte array containing the Ethernet address

### Description

Return true if the address is the broadcast address.

Please note: addr must be aligned to u16.

bool **is\_unicast\_ether\_addr**(const u8 \* *addr*)  
Determine if the Ethernet address is unicast

### Parameters

**const u8 \* addr** Pointer to a six-byte array containing the Ethernet address

### Description

Return true if the address is a unicast address.

bool **is\_valid\_ether\_addr**(const u8 \* *addr*)  
Determine if the given Ethernet address is valid

### Parameters

**const u8 \* addr** Pointer to a six-byte array containing the Ethernet address

### Description

Check that the Ethernet address (MAC) is not 00:00:00:00:00:00, is not a multicast address, and is not FF:FF:FF:FF:FF:FF.

Return true if the address is valid.

Please note: `addr` must be aligned to `u16`.

bool **eth\_proto\_is\_802\_3**(*\_\_be16 proto*)

Determine if a given Ethertype/length is a protocol

### Parameters

*\_\_be16 proto* Ethertype/length value to be tested

### Description

Check that the value from the Ethertype/length field is a valid Ethertype.

Return true if the valid is an 802.3 supported Ethertype.

void **eth\_random\_addr**(*u8 \* addr*)

Generate software assigned random Ethernet address

### Parameters

*u8 \* addr* Pointer to a six-byte array containing the Ethernet address

### Description

Generate a random Ethernet address (MAC) that is not multicast and has the local assigned bit set.

void **eth\_broadcast\_addr**(*u8 \* addr*)

Assign broadcast address

### Parameters

*u8 \* addr* Pointer to a six-byte array containing the Ethernet address

### Description

Assign the broadcast address to the given address array.

void **eth\_zero\_addr**(*u8 \* addr*)

Assign zero address

### Parameters

*u8 \* addr* Pointer to a six-byte array containing the Ethernet address

### Description

Assign the zero address to the given address array.

void **eth\_hw\_addr\_random**(*struct net\_device \* dev*)

Generate software assigned random Ethernet and set device flag

### Parameters

*struct net\_device \* dev* pointer to `net_device` structure

### Description

Generate a random Ethernet address (MAC) to be used by a net device and set `addr_assign_type` so the state can be read by `sysfs` and be used by userspace.

void **ether\_addr\_copy**(*u8 \* dst, const u8 \* src*)

Copy an Ethernet address

### Parameters

*u8 \* dst* Pointer to a six-byte array Ethernet address destination

*const u8 \* src* Pointer to a six-byte array Ethernet address source

### Description

Please note: `dst` & `src` must both be aligned to u16.

```
void eth_hw_addr_inherit(struct net_device * dst, struct net_device * src)  
    Copy dev_addr from another net_device
```

### Parameters

**struct net\_device \* dst** pointer to net\_device to copy dev\_addr to

**struct net\_device \* src** pointer to net\_device to copy dev\_addr from

### Description

Copy the Ethernet address from one net\_device to another along with the address attributes (`addr_assign_type`).

```
bool ether_addr_equal(const u8 * addr1, const u8 * addr2)  
    Compare two Ethernet addresses
```

### Parameters

**const u8 \* addr1** Pointer to a six-byte array containing the Ethernet address

**const u8 \* addr2** Pointer other six-byte array containing the Ethernet address

### Description

Compare two Ethernet addresses, returns true if equal

Please note: `addr1` & `addr2` must both be aligned to u16.

```
bool ether_addr_equal_64bits(const u8 addr1, const u8 addr2)  
    Compare two Ethernet addresses
```

### Parameters

**const u8 addr1** Pointer to an array of 8 bytes

**const u8 addr2** Pointer to an other array of 8 bytes

### Description

Compare two Ethernet addresses, returns true if equal, false otherwise.

The function doesn't need any conditional branches and possibly uses word memory accesses on CPU allowing cheap unaligned memory reads. `arrays = { byte1, byte2, byte3, byte4, byte5, byte6, pad1, pad2 }`

Please note that alignment of `addr1` & `addr2` are only guaranteed to be 16 bits.

```
bool ether_addr_equal_unaligned(const u8 * addr1, const u8 * addr2)  
    Compare two not u16 aligned Ethernet addresses
```

### Parameters

**const u8 \* addr1** Pointer to a six-byte array containing the Ethernet address

**const u8 \* addr2** Pointer other six-byte array containing the Ethernet address

### Description

Compare two Ethernet addresses, returns true if equal

Please note: Use only when any Ethernet address may not be u16 aligned.

```
bool ether_addr_equal_masked(const u8 * addr1, const u8 * addr2, const u8 * mask)  
    Compare two Ethernet addresses with a mask
```

### Parameters

**const u8 \* addr1** Pointer to a six-byte array containing the 1st Ethernet address

**const u8 \* addr2** Pointer to a six-byte array containing the 2nd Ethernet address

**const u8 \* mask** Pointer to a six-byte array containing the Ethernet address bitmask

### Description

Compare two Ethernet addresses with a mask, returns true if for every bit set in the bitmask the equivalent bits in the ethernet addresses are equal. Using a mask with all bits set is a slower `ether_addr_equal`.

**u64 ether\_addr\_to\_u64**(const u8 \* *addr*)  
Convert an Ethernet address into a u64 value.

### Parameters

**const u8 \* addr** Pointer to a six-byte array containing the Ethernet address

### Description

Return a u64 value of the address

**void u64\_to\_ether\_addr**(u64 *u*, u8 \* *addr*)  
Convert a u64 to an Ethernet address.

### Parameters

**u64 u** u64 to convert to an Ethernet MAC address

**u8 \* addr** Pointer to a six-byte array to contain the Ethernet address

**void eth\_addr\_dec**(u8 \* *addr*)  
Decrement the given MAC address

### Parameters

**u8 \* addr** Pointer to a six-byte array containing Ethernet address to decrement

**bool is\_etherdev\_addr**(const struct *net\_device* \* *dev*, const u8 *addr*)  
Tell if given Ethernet address belongs to the device.

### Parameters

**const struct net\_device \* dev** Pointer to a device structure

**const u8 addr** Pointer to a six-byte array containing the Ethernet address

### Description

Compare passed address with all addresses of the device. Return true if the address if one of the device addresses.

Note that this function calls `ether_addr_equal_64bits()` so take care of the right padding.

**unsigned long compare\_ether\_header**(const void \* *a*, const void \* *b*)  
Compare two Ethernet headers

### Parameters

**const void \* a** Pointer to Ethernet header

**const void \* b** Pointer to Ethernet header

### Description

Compare two Ethernet headers, returns 0 if equal. This assumes that the network header (i.e., IP header) is 4-byte aligned OR the platform can handle unaligned access. This is the case for all packets coming into `netif_receive_skb` or similar entry points.

**int eth\_skb\_pad**(struct *sk\_buff* \* *skb*)  
Pad buffer to minimum number of octets for Ethernet frame

### Parameters

**struct sk\_buff \* skb** Buffer to pad



**Description**

An Ethernet frame should have a minimum size of 60 bytes. This function takes short frames and pads them with zeros up to the 60 byte limit.

```
void napi_schedule(struct napi_struct * n)
    schedule NAPI poll
```

**Parameters**

```
struct napi_struct * n NAPI context
```

**Description**

Schedule NAPI poll routine to be called if it is not already running.

```
void napi_schedule_irqoff(struct napi_struct * n)
    schedule NAPI poll
```

**Parameters**

```
struct napi_struct * n NAPI context
```

**Description**

Variant of *napi\_schedule()*, assuming hard irqs are masked.

```
bool napi_complete(struct napi_struct * n)
    NAPI processing complete
```

**Parameters**

```
struct napi_struct * n NAPI context
```

**Description**

Mark NAPI processing as complete. Consider using *napi\_complete\_done()* instead. Return false if device should avoid rearming interrupts.

```
bool napi_hash_del(struct napi_struct * napi)
    remove a NAPI from global table
```

**Parameters**

```
struct napi_struct * napi NAPI context
```

**Description**

Warning: caller must observe RCU grace period before freeing memory containing **napi**, if this function returns true.

**Note**

core networking stack automatically calls it from *netif\_napi\_del()*. Drivers might want to call this helper to combine all the needed RCU grace periods into a single one.

```
void napi_disable(struct napi_struct * n)
    prevent NAPI from scheduling
```

**Parameters**

```
struct napi_struct * n NAPI context
```

**Description**

Stop NAPI from being scheduled on this context. Waits till any outstanding processing completes.

```
void napi_enable(struct napi_struct * n)
    enable NAPI scheduling
```

**Parameters**

```
struct napi_struct * n NAPI context
```

## Description

Resume NAPI from being scheduled on this context. Must be paired with `napi_disable`.

```
void napi_synchronize(const struct napi_struct * n)
    wait until NAPI is not running
```

## Parameters

`const struct napi_struct * n` NAPI context

## Description

Wait until NAPI is done being scheduled on this context. Waits till any outstanding processing completes but does not disable future activations.

```
enum netdev_priv_flags
    struct net_device priv_flags
```

## Constants

**IFF\_802\_1Q\_VLAN** 802.1Q VLAN device

**IFF\_EBRIDGE** Ethernet bridging device

**IFF\_BONDING** bonding master or slave

**IFF\_ISATAP** ISATAP interface (RFC4214)

**IFF\_WAN\_HDLC** WAN HDLC device

**IFF\_XMIT\_DST\_RELEASE** `dev_hard_start_xmit()` is allowed to release `skb->dst`

**IFF\_DONT\_BRIDGE** disallow bridging this ether dev

**IFF\_DISABLE\_NETPOLL** disable netpoll at run-time

**IFF\_MACVLAN\_PORT** device used as macvlan port

**IFF\_BRIDGE\_PORT** device used as bridge port

**IFF\_OVS\_DATAPATH** device used as Open vSwitch datapath port

**IFF\_TX\_SKB\_SHARING** The interface supports sharing skbs on transmit

**IFF\_UNICAST\_FLT** Supports unicast filtering

**IFF\_TEAM\_PORT** device used as team port

**IFF\_SUPP\_NOFCS** device supports sending custom FCS

**IFF\_LIVE\_ADDR\_CHANGE** device supports hardware address change when it's running

**IFF\_MACVLAN** Macvlan device

**IFF\_XMIT\_DST\_RELEASE\_PERM** **IFF\_XMIT\_DST\_RELEASE** not taking into account underlying stacked devices

**IFF\_IPVLAN\_MASTER** IPvlan master device

**IFF\_IPVLAN\_SLAVE** IPvlan slave device

**IFF\_L3MDEV\_MASTER** device is an L3 master device

**IFF\_NO\_QUEUE** device can run without qdisc attached

**IFF\_OPENVSWITCH** device is a Open vSwitch master

**IFF\_L3MDEV\_SLAVE** device is enslaved to an L3 master device

**IFF\_TEAM** device is a team device

**IFF\_RXFH\_CONFIGURED** device has had Rx Flow indirection table configured

**IFF\_PHONY\_HEADROOM** the headroom value is controlled by an external entity (i.e. the master device for bridged veth)

**IFF\_MACSEC** device is a MACsec device

### Description

These are the *struct net\_device*, they are only set internally by drivers and used in the kernel. These flags are invisible to userspace; this means that the order of these flags can change during any kernel release.

You should have a pretty good reason to be extending these flags.

struct **net\_device**

The DEVICE structure.

### Definition

```
struct net_device {
    char name;
    struct hlist_node name_hlist;
    char * ifalias;
    unsigned long mem_end;
    unsigned long mem_start;
    unsigned long base_addr;
    int irq;
    atomic_t carrier_changes;
    unsigned long state;
    struct list_head dev_list;
    struct list_head napi_list;
    struct list_head unreg_list;
    struct list_head close_list;
    struct list_head ptype_all;
    struct list_head ptype_specific;
    struct {unnamed_struct};
#ifdef IS_ENABLED(CONFIG_GARP)
    struct garp_port __rcu * garp_port;
#endif
#ifdef IS_ENABLED(CONFIG_MRP)
    struct mrp_port __rcu * mrp_port;
#endif
    struct device dev;
    const struct attribute_group * sysfs_groups;
    const struct attribute_group * sysfs_rx_queue_group;
    const struct rtnl_link_ops * rtnl_link_ops;
#define GSO_MAX_SIZE      65536
    unsigned int gso_max_size;
#define GSO_MAX_SEGS      65535
    u16 gso_max_segs;
#ifdef CONFIG_DCB
    const struct dcbnl_rtnl_ops * dcbnl_ops;
#endif
    u8 num_tc;
    struct netdev_tc_txq tc_to_txq;
    u8 prio_tc_map;
#ifdef IS_ENABLED(CONFIG_FCOE)
    unsigned int fcoe_ddp_xid;
#endif
#ifdef IS_ENABLED(CONFIG_CGROUP_NET_PRIO)
    struct netprio_map __rcu * priomap;
#endif
    struct phy_device * phydev;
    struct lock_class_key * qdisc_tx_busylock;
    struct lock_class_key * qdisc_running_key;
    bool proto_down;
};
```

### Members

**name** This is the first field of the “visible” part of this structure (i.e. as seen by users in the “Space.c” file). It is the name of the interface.

**name\_hlist** Device name hash chain, please keep it close to name[]

**ifalias** SNMP alias

**mem\_end** Shared memory end

**mem\_start** Shared memory start

**base\_addr** Device I/O address

**irq** Device IRQ number

**carrier\_changes** Stats to monitor carrier on<->off transitions

**state** Generic network queuing layer state, see netdev\_state\_t

**dev\_list** The global list of network devices

**napi\_list** List entry used for polling NAPI devices

**unreg\_list** List entry when we are unregistering the device; see the function unregister\_netdev

**close\_list** List entry used when we are closing the device

**ptype\_all** Device-specific packet handlers for all protocols

**ptype\_specific** Device-specific, protocol-specific packet handlers

**{unnamed\_struct}** anonymous

**garp\_port** GARP

**mrp\_port** MRP

**dev** Class/net/name entry

**sysfs\_groups** Space for optional device, statistics and wireless sysfs groups

**sysfs\_rx\_queue\_group** Space for optional per-rx queue attributes

**rtnl\_link\_ops** Rtnl\_link\_ops

**gso\_max\_size** Maximum size of generic segmentation offload

**gso\_max\_segs** Maximum number of segments that can be passed to the NIC for GSO

**dcbnl\_ops** Data Center Bridging netlink ops

**num\_tc** Number of traffic classes in the net device

**tc\_to\_txq** XXX: need comments on this one

**prio\_tc\_map** XXX: need comments on this one

**fcoe\_ddp\_xid** Max exchange id for FCoE LRO by ddp

**priomap** XXX: need comments on this one

**phydev** Physical device may attach itself for hardware timestamping

**qdisc\_tx\_busylock** lockdep class annotating Qdisc->busylock spinlock

**qdisc\_running\_key** lockdep class annotating Qdisc->running seqcount

**proto\_down** protocol port state information can be sent to the switch driver and used to set the phys state of the switch port.

## Description

Actually, this whole structure is a big mistake. It mixes I/O data with strictly “high-level” data, and it has to know about almost every data structure used in the INET module.

interface address info:

FIXME: cleanup struct `net_device` such that network protocol info moves out.

void \* **netdev\_priv**(const struct *net\_device* \* *dev*)  
access network device private data

#### Parameters

const struct *net\_device* \* *dev* network device

#### Description

Get network device private data

void **netif\_napi\_add**(struct *net\_device* \* *dev*, struct *napi\_struct* \* *napi*, int (\**poll*) (struct *napi\_struct* \*, int, int *weight*)  
initialize a NAPI context

#### Parameters

struct *net\_device* \* *dev* network device

struct *napi\_struct* \* *napi* NAPI context

int (\*)(struct *napi\_struct* \*,int) *poll* polling function

int *weight* default weight

#### Description

*netif\_napi\_add()* must be used to initialize a NAPI context prior to calling *any* of the other NAPI-related functions.

void **netif\_tx\_napi\_add**(struct *net\_device* \* *dev*, struct *napi\_struct* \* *napi*, int (\**poll*) (struct *napi\_struct* \*, int, int *weight*)  
initialize a NAPI context

#### Parameters

struct *net\_device* \* *dev* network device

struct *napi\_struct* \* *napi* NAPI context

int (\*)(struct *napi\_struct* \*,int) *poll* polling function

int *weight* default weight

#### Description

This variant of *netif\_napi\_add()* should be used from drivers using NAPI to exclusively poll a TX queue. This will avoid we add it into *napi\_hash[]*, thus polluting this hash table.

void **netif\_napi\_del**(struct *napi\_struct* \* *napi*)  
remove a NAPI context

#### Parameters

struct *napi\_struct* \* *napi* NAPI context

#### Description

*netif\_napi\_del()* removes a NAPI context from the network device NAPI list

void **netif\_start\_queue**(struct *net\_device* \* *dev*)  
allow transmit

#### Parameters

struct *net\_device* \* *dev* network device

#### Description

Allow upper layers to call the device *hard\_start\_xmit* routine.

void **netif\_wake\_queue**(struct *net\_device* \* *dev*)  
restart transmit

#### Parameters

**struct net\_device \* dev** network device

#### Description

Allow upper layers to call the device `hard_start_xmit` routine. Used for flow control when transmit resources are available.

void **netif\_stop\_queue**(struct *net\_device* \* *dev*)  
stop transmitted packets

#### Parameters

**struct net\_device \* dev** network device

#### Description

Stop upper layers calling the device `hard_start_xmit` routine. Used for flow control when transmit resources are unavailable.

bool **netif\_queue\_stopped**(const struct *net\_device* \* *dev*)  
test if transmit queue is flowblocked

#### Parameters

**const struct net\_device \* dev** network device

#### Description

Test if transmit queue on device is currently unable to send.

void **netdev\_txq\_bql\_enqueue\_prefetchw**(struct *netdev\_queue* \* *dev\_queue*)  
prefetch bql data for write

#### Parameters

**struct netdev\_queue \* dev\_queue** pointer to transmit queue

#### Description

BQL enabled drivers might use this helper in their `ndo_start_xmit()`, to give appropriate hint to the CPU.

void **netdev\_txq\_bql\_complete\_prefetchw**(struct *netdev\_queue* \* *dev\_queue*)  
prefetch bql data for write

#### Parameters

**struct netdev\_queue \* dev\_queue** pointer to transmit queue

#### Description

BQL enabled drivers might use this helper in their TX completion path, to give appropriate hint to the CPU.

void **netdev\_sent\_queue**(struct *net\_device* \* *dev*, unsigned int *bytes*)  
report the number of bytes queued to hardware

#### Parameters

**struct net\_device \* dev** network device

**unsigned int bytes** number of bytes queued to the hardware device queue

#### Description

Report the number of bytes queued for sending/completion to the network device hardware queue. **bytes** should be a good approximation and should exactly match `netdev_completed_queue()` **bytes**

void **netdev\_completed\_queue**(struct *net\_device* \* *dev*, unsigned int *pkts*, unsigned int *bytes*)  
report bytes and packets completed by device

#### Parameters

**struct net\_device \* dev** network device  
**unsigned int pkts** actual number of packets sent over the medium  
**unsigned int bytes** actual number of bytes sent over the medium

#### Description

Report the number of bytes and packets transmitted by the network device hardware queue over the physical medium, **bytes** must exactly match the **bytes** amount passed to *netdev\_sent\_queue()*

void **netdev\_reset\_queue**(struct *net\_device* \* *dev\_queue*)  
reset the packets and bytes count of a network device

#### Parameters

**struct net\_device \* dev\_queue** network device

#### Description

Reset the bytes and packet count of a network device and clear the software flow control OFF bit for this network device

u16 **netdev\_cap\_txqueue**(struct *net\_device* \* *dev*, u16 *queue\_index*)  
check if selected tx queue exceeds device queues

#### Parameters

**struct net\_device \* dev** network device  
**u16 queue\_index** given tx queue index

#### Description

Returns 0 if given tx queue index  $\geq$  number of device tx queues, otherwise returns the originally passed tx queue index.

bool **netif\_running**(const struct *net\_device* \* *dev*)  
test if up

#### Parameters

**const struct net\_device \* dev** network device

#### Description

Test if the device has been brought up.

void **netif\_start\_subqueue**(struct *net\_device* \* *dev*, u16 *queue\_index*)  
allow sending packets on subqueue

#### Parameters

**struct net\_device \* dev** network device  
**u16 queue\_index** sub queue index

#### Description

Start individual transmit queue of a device with multiple transmit queues.

void **netif\_stop\_subqueue**(struct *net\_device* \* *dev*, u16 *queue\_index*)  
stop sending packets on subqueue

#### Parameters

**struct net\_device \* dev** network device

**u16 queue\_index** sub queue index

### Description

Stop individual transmit queue of a device with multiple transmit queues.

bool **\_\_netif\_subqueue\_stopped**(const struct *net\_device* \* dev, u16 *queue\_index*)  
test status of subqueue

### Parameters

const struct *net\_device* \* **dev** network device

**u16 queue\_index** sub queue index

### Description

Check individual transmit queue of a device with multiple transmit queues.

void **netif\_wake\_subqueue**(struct *net\_device* \* dev, u16 *queue\_index*)  
allow sending packets on subqueue

### Parameters

struct *net\_device* \* **dev** network device

**u16 queue\_index** sub queue index

### Description

Resume individual transmit queue of a device with multiple transmit queues.

bool **netif\_is\_multiqueue**(const struct *net\_device* \* dev)  
test if device has multiple transmit queues

### Parameters

const struct *net\_device* \* **dev** network device

### Description

Check if device has multiple transmit queues

void **dev\_put**(struct *net\_device* \* dev)  
release reference to device

### Parameters

struct *net\_device* \* **dev** network device

### Description

Release reference to device to allow it to be freed.

void **dev\_hold**(struct *net\_device* \* dev)  
get reference to device

### Parameters

struct *net\_device* \* **dev** network device

### Description

Hold reference to device to keep it from being freed.

bool **netif\_carrier\_ok**(const struct *net\_device* \* dev)  
test if carrier present

### Parameters

const struct *net\_device* \* **dev** network device

### Description

Check if carrier is present on device



void **netif\_dormant\_on**(struct *net\_device* \* dev)  
mark device as dormant.

**Parameters**

**struct net\_device \* dev** network device

**Description**

Mark device as dormant (as per RFC2863).

The dormant state indicates that the relevant interface is not actually in a condition to pass packets (i.e., it is not 'up') but is in a "pending" state, waiting for some external event. For "on-demand" interfaces, this new state identifies the situation where the interface is waiting for events to place it in the up state.

void **netif\_dormant\_off**(struct *net\_device* \* dev)  
set device as not dormant.

**Parameters**

**struct net\_device \* dev** network device

**Description**

Device is not in dormant state.

bool **netif\_dormant**(const struct *net\_device* \* dev)  
test if device is dormant

**Parameters**

**const struct net\_device \* dev** network device

**Description**

Check if device is dormant.

bool **netif\_oper\_up**(const struct *net\_device* \* dev)  
test if device is operational

**Parameters**

**const struct net\_device \* dev** network device

**Description**

Check if carrier is operational

bool **netif\_device\_present**(struct *net\_device* \* dev)  
is device available or removed

**Parameters**

**struct net\_device \* dev** network device

**Description**

Check if device has not been removed from system.

void **netif\_tx\_lock**(struct *net\_device* \* dev)  
grab network device transmit lock

**Parameters**

**struct net\_device \* dev** network device

**Description**

Get network device transmit lock

int **\_\_dev\_uc\_sync**(struct *net\_device* \* dev, int (\*sync) (struct *net\_device* \*, const unsigned char \*,  
int (\*unsync) (struct *net\_device* \*, const unsigned char \*))  
Synchronize device's unicast list

### Parameters

**struct net\_device \* dev** device to sync

**int (\*)(struct net\_device \*,const unsigned char \*) sync** function to call if address should be added

**int (\*)(struct net\_device \*,const unsigned char \*) unsync** function to call if address should be removed

### Description

Add newly added addresses to the interface, and release addresses that have been deleted.

**void \_\_dev\_uc\_unsync**(struct *net\_device* \* dev, int (\*unsync) (struct *net\_device* \*, const unsigned char \*))

Remove synchronized addresses from device

### Parameters

**struct net\_device \* dev** device to sync

**int (\*)(struct net\_device \*,const unsigned char \*) unsync** function to call if address should be removed

### Description

Remove all addresses that were added to the device by `dev_uc_sync()`.

**int \_\_dev\_mc\_sync**(struct *net\_device* \* dev, int (\*sync) (struct *net\_device* \*, const unsigned char \*, int (\*unsync) (struct *net\_device* \*, const unsigned char \*))

Synchronize device's multicast list

### Parameters

**struct net\_device \* dev** device to sync

**int (\*)(struct net\_device \*,const unsigned char \*) sync** function to call if address should be added

**int (\*)(struct net\_device \*,const unsigned char \*) unsync** function to call if address should be removed

### Description

Add newly added addresses to the interface, and release addresses that have been deleted.

**void \_\_dev\_mc\_unsync**(struct *net\_device* \* dev, int (\*unsync) (struct *net\_device* \*, const unsigned char \*))

Remove synchronized addresses from device

### Parameters

**struct net\_device \* dev** device to sync

**int (\*)(struct net\_device \*,const unsigned char \*) unsync** function to call if address should be removed

### Description

Remove all addresses that were added to the device by `dev_mc_sync()`.

## PHY Support

**void phy\_print\_status**(struct phy\_device \* phydev)

Convenience function to print out the current phy status

### Parameters

**struct phy\_device \* phydev** the phy\_device struct

int **phy\_restart\_aneg**(struct phy\_device \* *phydev*)  
restart auto-negotiation

#### Parameters

**struct phy\_device \* phydev** target phy\_device struct

#### Description

Restart the autonegotiation on **phydev**. Returns  $\geq 0$  on success or negative errno on error.

int **phy\_aneg\_done**(struct phy\_device \* *phydev*)  
return auto-negotiation status

#### Parameters

**struct phy\_device \* phydev** target phy\_device struct

#### Description

Return the auto-negotiation status from this **phydev**. Returns  $> 0$  on success or  $< 0$  on error. 0 means that auto-negotiation is still pending.

int **phy\_ethtool\_sset**(struct phy\_device \* *phydev*, struct ethtool\_cmd \* *cmd*)  
generic ethtool sset function, handles all the details

#### Parameters

**struct phy\_device \* phydev** target phy\_device struct

**struct ethtool\_cmd \* cmd** ethtool\_cmd

#### Description

A few notes about parameter checking:

- We don't set port or transceiver, so we don't care what they were set to.
- [\*phy\\_start\\_aneg\(\)\*](#) will make sure forced settings are sane, and choose the next best ones from the ones selected, so we don't care if ethtool tries to give us bad values.

int **phy\_mii\_ioctl**(struct phy\_device \* *phydev*, struct ifreq \* *ifr*, int *cmd*)  
generic PHY MII ioctl interface

#### Parameters

**struct phy\_device \* phydev** the phy\_device struct

**struct ifreq \* ifr** struct ifreq for socket ioctl's

**int cmd** ioctl cmd to execute

#### Description

Note that this function is currently incompatible with the PHYCONTROL layer. It changes registers without regard to current state. Use at own risk.

int **phy\_start\_aneg**(struct phy\_device \* *phydev*)  
start auto-negotiation for this PHY device

#### Parameters

**struct phy\_device \* phydev** the phy\_device struct

#### Description

**Sanitizes the settings (if we're not autonegotiating** them), and then calls the driver's `config_aneg` function. If the PHYCONTROL Layer is operating, we change the state to reflect the beginning of Auto-negotiation or forcing.

void **phy\_start\_machine**(struct phy\_device \* *phydev*)  
start PHY state machine tracking

#### Parameters

**struct phy\_device \* phydev** the phy\_device struct

### Description

**The PHY infrastructure can run a state machine** which tracks whether the PHY is starting up, negotiating, etc. This function starts the delayed workqueue which tracks the state of the PHY. If you want to maintain your own state machine, do not call this function.

int **phy\_start\_interrupts**(struct phy\_device \* *phydev*)  
request and enable interrupts for a PHY device

### Parameters

**struct phy\_device \* phydev** target phy\_device struct

### Description

**Request the interrupt for the given PHY.** If this fails, then we set irq to PHY\_POLL. Otherwise, we enable the interrupts in the PHY. This should only be called with a valid IRQ number. Returns 0 on success or < 0 on error.

int **phy\_stop\_interrupts**(struct phy\_device \* *phydev*)  
disable interrupts from a PHY device

### Parameters

**struct phy\_device \* phydev** target phy\_device struct

void **phy\_stop**(struct phy\_device \* *phydev*)  
Bring down the PHY link, and stop checking the status

### Parameters

**struct phy\_device \* phydev** target phy\_device struct

void **phy\_start**(struct phy\_device \* *phydev*)  
start or restart a PHY device

### Parameters

**struct phy\_device \* phydev** target phy\_device struct

### Description

**Indicates the attached device's readiness to** handle PHY-related work. Used during startup to start the PHY, and after a call to [phy\\_stop\(\)](#) to resume operation. Also used to indicate the MDIO bus has cleared an error condition.

void **phy\_mac\_interrupt**(struct phy\_device \* *phydev*, int *new\_link*)  
MAC says the link has changed

### Parameters

**struct phy\_device \* phydev** phy\_device struct with changed link

int **new\_link** Link is Up/Down.

### Description

**The MAC layer is able indicate there has been a change** in the PHY link status. Set the new link status, and trigger the state machine, work a work queue.

int **phy\_init\_eee**(struct phy\_device \* *phydev*, bool *clk\_stop\_enable*)  
init and check the EEE feature

### Parameters

**struct phy\_device \* phydev** target phy\_device struct

bool **clk\_stop\_enable** PHY may stop the clock during LPI

**Description**

it checks if the Energy-Efficient Ethernet (EEE) is supported by looking at the MMD registers 3.20 and 7.60/61 and it programs the MMD register 3.0 setting the "Clock stop enable" bit if required.

int **phy\_get\_eee\_err**(struct phy\_device \* *phydev*)  
report the EEE wake error count

**Parameters**

**struct phy\_device \* phydev** target phy\_device struct

**Description**

it is to report the number of time where the PHY failed to complete its normal wake sequence.

int **phy\_ethtool\_get\_eee**(struct phy\_device \* *phydev*, struct ethtool\_eee \* *data*)  
get EEE supported and status

**Parameters**

**struct phy\_device \* phydev** target phy\_device struct

**struct ethtool\_eee \* data** ethtool\_eee data

**Description**

it reports the Supported/Advertisement/LP Advertisement capabilities.

int **phy\_ethtool\_set\_eee**(struct phy\_device \* *phydev*, struct ethtool\_eee \* *data*)  
set EEE supported and status

**Parameters**

**struct phy\_device \* phydev** target phy\_device struct

**struct ethtool\_eee \* data** ethtool\_eee data

**Description**

it is to program the Advertisement EEE register.

int **phy\_clear\_interrupt**(struct phy\_device \* *phydev*)  
Ack the phy device's interrupt

**Parameters**

**struct phy\_device \* phydev** the phy\_device struct

**Description**

If the **phydev** driver has an `ack_interrupt` function, call it to ack and clear the phy device's interrupt.

Returns 0 on success or < 0 on error.

int **phy\_config\_interrupt**(struct phy\_device \* *phydev*, u32 *interrupts*)  
configure the PHY device for the requested interrupts

**Parameters**

**struct phy\_device \* phydev** the phy\_device struct

**u32 interrupts** interrupt flags to configure for this **phydev**

**Description**

Returns 0 on success or < 0 on error.

const struct phy\_setting \* **phy\_find\_valid**(int *speed*, int *duplex*, u32 *supported*)  
find a PHY setting that matches the requested parameters

**Parameters**

**int speed** desired speed

**int duplex** desired duplex

**u32 supported** mask of supported link modes

### Description

Locate a supported phy setting that is, in priority order: - an exact match for the specified speed and duplex mode - a match for the specified speed, or slower speed - the slowest supported speed Returns the matched phy\_setting entry, or NULL if no supported phy settings were found.

unsigned int **phy\_supported\_speeds**(struct phy\_device \* *phy*, unsigned int \* *speeds*, unsigned int *size*)  
return all speeds currently supported by a phy device

### Parameters

**struct phy\_device \* phy** The phy device to return supported speeds of.

**unsigned int \* speeds** buffer to store supported speeds in.

**unsigned int size** size of speeds buffer.

### Description

Returns the number of supported speeds, and fills the speeds buffer with the supported speeds. If speeds buffer is too small to contain all currently supported speeds, will return as many speeds as can fit.

bool **phy\_check\_valid**(int *speed*, int *duplex*, u32 *features*)  
check if there is a valid PHY setting which matches speed, duplex, and feature mask

### Parameters

**int speed** speed to match

**int duplex** duplex to match

**u32 features** A mask of the valid settings

### Description

Returns true if there is a valid setting, false otherwise.

void **phy\_sanitize\_settings**(struct phy\_device \* *phydev*)  
make sure the PHY is set to supported speed and duplex

### Parameters

**struct phy\_device \* phydev** the target phy\_device struct

### Description

**Make sure the PHY is set to supported speeds and duplexes.** Drop down by one in this order: 1000/FULL, 1000/HALF, 100/FULL, 100/HALF, 10/FULL, 10/HALF.

int **phy\_start\_aneg\_priv**(struct phy\_device \* *phydev*, bool *sync*)  
start auto-negotiation for this PHY device

### Parameters

**struct phy\_device \* phydev** the phy\_device struct

**bool sync** indicate whether we should wait for the workqueue cancelation

### Description

**Sanitizes the settings (if we're not autonegotiating** them), and then calls the driver's config\_aneg function. If the PHYCONTROL Layer is operating, we change the state to reflect the beginning of Auto-negotiation or forcing.

void **phy\_trigger\_machine**(struct phy\_device \* *phydev*, bool *sync*)  
trigger the state machine to run

### Parameters

**struct phy\_device \* phydev** the phy\_device struct

**bool sync** indicate whether we should wait for the workqueue cancelation

### Description

**There has been a change in state which requires that the** state machine runs.

void **phy\_stop\_machine**(struct phy\_device \* *phydev*)  
stop the PHY state machine tracking

### Parameters

**struct phy\_device \* phydev** target phy\_device struct

### Description

**Stops the state machine delayed workqueue, sets the** state to UP (unless it wasn't up yet). This function must be called BEFORE phy\_detach.

void **phy\_error**(struct phy\_device \* *phydev*)  
enter HALTED state for this PHY device

### Parameters

**struct phy\_device \* phydev** target phy\_device struct

### Description

Moves the PHY to the HALTED state in response to a read or write error, and tells the controller the link is down. Must not be called from interrupt context, or while the phydev->lock is held.

irqreturn\_t **phy\_interrupt**(int *irq*, void \* *phy\_dat*)  
PHY interrupt handler

### Parameters

**int irq** interrupt line

**void \* phy\_dat** phy\_device pointer

### Description

When a PHY interrupt occurs, the handler disables interrupts, and uses phy\_change to handle the interrupt.

int **phy\_enable\_interrupts**(struct phy\_device \* *phydev*)  
Enable the interrupts from the PHY side

### Parameters

**struct phy\_device \* phydev** target phy\_device struct

int **phy\_disable\_interrupts**(struct phy\_device \* *phydev*)  
Disable the PHY interrupts from the PHY side

### Parameters

**struct phy\_device \* phydev** target phy\_device struct

void **phy\_change**(struct phy\_device \* *phydev*)  
Called by the phy\_interrupt to handle PHY changes

### Parameters

**struct phy\_device \* phydev** phy\_device struct that interrupted

void **phy\_change\_work**(struct work\_struct \* *work*)  
Scheduled by the phy\_mac\_interrupt to handle PHY changes

### Parameters

**struct work\_struct \* work** work\_struct that describes the work to be done

void **phy\_state\_machine**(struct work\_struct \* *work*)  
Handle the state machine

#### Parameters

**struct work\_struct \* work** work\_struct that describes the work to be done

int **phy\_register\_fixup**(const char \* *bus\_id*, u32 *phy\_uid*, u32 *phy\_uid\_mask*, int (\*run) (struct phy\_device \*))  
creates a new phy\_fixup and adds it to the list

#### Parameters

**const char \* bus\_id** A string which matches phydev->mdio.dev.bus\_id (or PHY\_ANY\_ID)

**u32 phy\_uid** Used to match against phydev->phy\_id (the UID of the PHY) It can also be PHY\_ANY\_UID

**u32 phy\_uid\_mask** Applied to phydev->phy\_id and fixup->phy\_uid before comparison

**int (\*)(struct phy\_device \*) run** The actual code to be run when a matching PHY is found

int **phy\_unregister\_fixup**(const char \* *bus\_id*, u32 *phy\_uid*, u32 *phy\_uid\_mask*)  
remove a phy\_fixup from the list

#### Parameters

**const char \* bus\_id** A string matches fixup->bus\_id (or PHY\_ANY\_ID) in phy\_fixup\_list

**u32 phy\_uid** A phy id matches fixup->phy\_id (or PHY\_ANY\_UID) in phy\_fixup\_list

**u32 phy\_uid\_mask** Applied to phy\_uid and fixup->phy\_uid before comparison

struct phy\_device \* **get\_phy\_device**(struct mii\_bus \* *bus*, int *addr*, bool *is\_c45*)  
reads the specified PHY device and returns its **phy\_device** struct

#### Parameters

**struct mii\_bus \* bus** the target MII bus

**int addr** PHY address on the MII bus

**bool is\_c45** If true the PHY uses the 802.3 clause 45 protocol

#### Description

**Reads the ID registers of the PHY at addr on the bus**, then allocates and returns the phy\_device to represent it.

int **phy\_device\_register**(struct phy\_device \* *phydev*)  
Register the phy device on the MDIO bus

#### Parameters

**struct phy\_device \* phydev** phy\_device structure to be added to the MDIO bus

void **phy\_device\_remove**(struct phy\_device \* *phydev*)  
Remove a previously registered phy device from the MDIO bus

#### Parameters

**struct phy\_device \* phydev** phy\_device structure to remove

#### Description

This doesn't free the phy\_device itself, it merely reverses the effects of *phy\_device\_register()*. Use *phy\_device\_free()* to free the device after calling this function.

struct phy\_device \* **phy\_find\_first**(struct mii\_bus \* *bus*)  
finds the first PHY device on the bus

#### Parameters

**struct mii\_bus \* bus** the target MII bus



int **phy\_connect\_direct**(struct *net\_device* \* *dev*, struct phy\_device \* *phydev*, void (\*handler) (struct *net\_device* \*, phy\_interface\_t *interface*)  
connect an ethernet device to a specific phy\_device

#### Parameters

**struct net\_device \* dev** the network device to connect

**struct phy\_device \* phydev** the pointer to the phy device

**void (\*)(struct net\_device \*) handler** callback function for state change notifications

**phy\_interface\_t interface** PHY device's interface

struct phy\_device \* **phy\_connect**(struct *net\_device* \* *dev*, const char \* *bus\_id*, void (\*handler) (struct *net\_device* \*, phy\_interface\_t *interface*)  
connect an ethernet device to a PHY device

#### Parameters

**struct net\_device \* dev** the network device to connect

**const char \* bus\_id** the id string of the PHY device to connect

**void (\*)(struct net\_device \*) handler** callback function for state change notifications

**phy\_interface\_t interface** PHY device's interface

#### Description

**Convenience function for connecting ethernet** devices to PHY devices. The default behavior is for the PHY infrastructure to handle everything, and only notify the connected driver when the link status changes. If you don't want, or can't use the provided functionality, you may choose to call only the subset of functions which provide the desired functionality.

void **phy\_disconnect**(struct phy\_device \* *phydev*)  
disable interrupts, stop state machine, and detach a PHY device

#### Parameters

**struct phy\_device \* phydev** target phy\_device struct

int **phy\_attach\_direct**(struct *net\_device* \* *dev*, struct phy\_device \* *phydev*, u32 *flags*,  
phy\_interface\_t *interface*)  
attach a network device to a given PHY device pointer

#### Parameters

**struct net\_device \* dev** network device to attach

**struct phy\_device \* phydev** Pointer to phy\_device to attach

**u32 flags** PHY device's dev\_flags

**phy\_interface\_t interface** PHY device's interface

#### Description

**Called by drivers to attach to a particular PHY** device. The phy\_device is found, and properly hooked up to the phy\_driver. If no driver is attached, then a generic driver is used. The phy\_device is given a ptr to the attaching device, and given a callback for link status change. The phy\_device is returned to the attaching driver. This function takes a reference on the phy device.

struct phy\_device \* **phy\_attach**(struct *net\_device* \* *dev*, const char \* *bus\_id*,  
phy\_interface\_t *interface*)  
attach a network device to a particular PHY device

#### Parameters

**struct net\_device \* dev** network device to attach

**const char \* bus\_id** Bus ID of PHY device to attach

**phy\_interface\_t interface** PHY device's interface

### Description

Same as [phy\\_attach\\_direct\(\)](#) except that a **PHY bus\_id** string is passed instead of a pointer to a struct `phy_device`.

void **phy\_detach**(struct `phy_device` \* *phydev*)  
detach a PHY device from its network device

### Parameters

struct `phy_device` \* **phydev** target `phy_device` struct

### Description

This detaches the phy device from its network device and the phy driver, and drops the reference count taken in [phy\\_attach\\_direct\(\)](#).

int **genphy\_setup\_forced**(struct `phy_device` \* *phydev*)  
configures/forces speed/duplex from **phydev**

### Parameters

struct `phy_device` \* **phydev** target `phy_device` struct

### Description

**Configures MII\_BMCR to force speed/duplex** to the values in `phydev`. Assumes that the values are valid. Please see [phy\\_sanitize\\_settings\(\)](#).

int **genphy\_restart\_aneg**(struct `phy_device` \* *phydev*)  
Enable and Restart Autonegotiation

### Parameters

struct `phy_device` \* **phydev** target `phy_device` struct

int **genphy\_config\_aneg**(struct `phy_device` \* *phydev*)  
restart auto-negotiation or write BMCR

### Parameters

struct `phy_device` \* **phydev** target `phy_device` struct

### Description

**If auto-negotiation is enabled, we configure the** advertising, and then restart auto-negotiation. If it is not enabled, then we write the BMCR.

int **genphy\_aneg\_done**(struct `phy_device` \* *phydev*)  
return auto-negotiation status

### Parameters

struct `phy_device` \* **phydev** target `phy_device` struct

### Description

**Reads the status register and returns 0 either if** auto-negotiation is incomplete, or if there was an error. Returns `BMSR_ANEGCOMPLETE` if auto-negotiation is done.

int **genphy\_update\_link**(struct `phy_device` \* *phydev*)  
update link status in **phydev**

### Parameters

struct `phy_device` \* **phydev** target `phy_device` struct

### Description

**Update the value in phydev->link to reflect the** current link value. In order to do this, we need to read the status register twice, keeping the second value.

int **genphy\_read\_status**(struct phy\_device \* *phydev*)  
check the link status and update current link state

#### Parameters

**struct phy\_device \* phydev** target phy\_device struct

#### Description

**Check the link, then figure out the current state** by comparing what we advertise with what the link partner advertises. Start by checking the gigabit possibilities, then move on to 10/100.

int **genphy\_soft\_reset**(struct phy\_device \* *phydev*)  
software reset the PHY via BMCR\_RESET bit

#### Parameters

**struct phy\_device \* phydev** target phy\_device struct

#### Description

Perform a software PHY reset using the standard BMCR\_RESET bit and poll for the reset bit to be cleared.

#### Return

0 on success, < 0 on failure

int **phy\_driver\_register**(struct phy\_driver \* *new\_driver*, struct module \* *owner*)  
register a phy\_driver with the PHY layer

#### Parameters

**struct phy\_driver \* new\_driver** new phy\_driver to register

**struct module \* owner** module owning this PHY

int **get\_phy\_c45\_ids**(struct mii\_bus \* *bus*, int *addr*, u32 \* *phy\_id*, struct phy\_c45\_device\_ids \* *c45\_ids*)  
reads the specified *addr* for its 802.3-c45 IDs.

#### Parameters

**struct mii\_bus \* bus** the target MII bus

**int addr** PHY address on the MII bus

**u32 \* phy\_id** where to store the ID retrieved.

**struct phy\_c45\_device\_ids \* c45\_ids** where to store the c45 ID information.

#### Description

If the PHY devices-in-package appears to be valid, it and the corresponding identifiers are stored in **c45\_ids**, zero is stored in **phy\_id**. Otherwise 0xffffffff is stored in **phy\_id**. Returns zero on success.

int **get\_phy\_id**(struct mii\_bus \* *bus*, int *addr*, u32 \* *phy\_id*, bool *is\_c45*, struct phy\_c45\_device\_ids \* *c45\_ids*)  
reads the specified *addr* for its ID.

#### Parameters

**struct mii\_bus \* bus** the target MII bus

**int addr** PHY address on the MII bus

**u32 \* phy\_id** where to store the ID retrieved.

**bool is\_c45** If true the PHY uses the 802.3 clause 45 protocol

**struct phy\_c45\_device\_ids \* c45\_ids** where to store the c45 ID information.

#### Description

In the case of a 802.3-c22 PHY, reads the ID registers of the PHY at **addr** on the **bus**, stores it in **phy\_id** and returns zero on success.

In the case of a 802.3-c45 PHY, *get\_phy\_c45\_ids()* is invoked, and its return value is in turn returned.

void **phy\_prepare\_link**(struct phy\_device \* *phydev*, void (\*handler) (struct *net\_device* \*))  
prepares the PHY layer to monitor link status

#### Parameters

struct phy\_device \* **phydev** target phy\_device struct

void (\*)(struct net\_device \*) **handler** callback function for link status change notifications

#### Description

**Tells the PHY infrastructure to handle the** gory details on monitoring link status (whether through polling or an interrupt), and to call back to the connected device driver when the link status changes. If you want to monitor your own link state, don't call this function.

int **phy\_poll\_reset**(struct phy\_device \* *phydev*)  
Safely wait until a PHY reset has properly completed

#### Parameters

struct phy\_device \* **phydev** The PHY device to poll

#### Description

**According to IEEE 802.3, Section 2, Subsection 22.2.4.1.1, as** published in 2008, a PHY reset may take up to 0.5 seconds. The MII BMCR register must be polled until the BMCR\_RESET bit clears.

Furthermore, any attempts to write to PHY registers may have no effect or even generate MDIO bus errors until this is complete.

Some PHYs (such as the Marvell 88E1111) don't entirely conform to the standard and do not fully reset after the BMCR\_RESET bit is set, and may even *REQUIRE* a soft-reset to properly restart autonegotiation. In an effort to support such broken PHYs, this function is separate from the standard *phy\_init\_hw()* which will zero all the other bits in the BMCR and reapply all driver-specific and board-specific fixups.

int **genphy\_config\_advert**(struct phy\_device \* *phydev*)  
sanitize and advertise auto-negotiation parameters

#### Parameters

struct phy\_device \* **phydev** target phy\_device struct

#### Description

**Writes MII\_ADVERTISE with the appropriate values,** after sanitizing the values to make sure we only advertise what is supported. Returns < 0 on error, 0 if the PHY's advertisement hasn't changed, and > 0 if it has changed.

int **genphy\_config\_eee\_advert**(struct phy\_device \* *phydev*)  
disable unwanted eee mode advertisement

#### Parameters

struct phy\_device \* **phydev** target phy\_device struct

#### Description

**Writes MDIO\_AN\_EEE\_ADV after disabling unsupported energy** efficient ethernet modes. Returns 0 if the PHY's advertisement hasn't changed, and 1 if it has changed.

int **phy\_probe**(struct device \* *dev*)  
probe and init a PHY device

#### Parameters

struct device \* **dev** device to probe and init

## Description

**Take care of setting up the phy\_device structure**, set the state to READY (the driver's init function should set it to STARTING if needed).

```
struct mii_bus * mdiobus_alloc_size(size_t size)  
    allocate a mii_bus structure
```

## Parameters

**size\_t size** extra amount of memory to allocate for private storage. If non-zero, then bus->priv is points to that memory.

## Description

called by a bus driver to allocate an mii\_bus structure to fill in.

```
struct mii_bus * devm_mdiobus_alloc_size(struct device * dev, int sizeof_priv)  
    Resource-managed mdiobus\_alloc\_size\(\)
```

## Parameters

**struct device \* dev** Device to allocate mii\_bus for

**int sizeof\_priv** Space to allocate for private structure.

## Description

Managed mdiobus\_alloc\_size. mii\_bus allocated with this function is automatically freed on driver detach.

If an mii\_bus allocated with this function needs to be freed separately, [devm\\_mdiobus\\_free\(\)](#) must be used.

## Return

Pointer to allocated mii\_bus on success, NULL on failure.

```
void devm_mdiobus_free(struct device * dev, struct mii_bus * bus)  
    Resource-managed mdiobus\_free\(\)
```

## Parameters

**struct device \* dev** Device this mii\_bus belongs to

**struct mii\_bus \* bus** the mii\_bus associated with the device

## Description

Free mii\_bus allocated with [devm\\_mdiobus\\_alloc\\_size\(\)](#).

```
struct mii_bus * of_mdio_find_bus(struct device_node * mdio_bus_np)  
    Given an mii_bus node, find the mii_bus.
```

## Parameters

**struct device\_node \* mdio\_bus\_np** Pointer to the mii\_bus.

## Description

Returns a reference to the mii\_bus, or NULL if none found. The embedded struct device will have its reference count incremented, and this must be put once the bus is finished with.

Because the association of a device\_node and mii\_bus is made via of\_mdio\_register(), the mii\_bus cannot be found before it is registered with of\_mdio\_register().

```
int __mdiobus_register(struct mii_bus * bus, struct module * owner)  
    bring up all the PHYs on a given bus and attach them to bus
```

## Parameters

**struct mii\_bus \* bus** target mii\_bus

**struct module \* owner** module containing bus accessor functions

## Description

**Called by a bus driver to bring up all the PHYs** on a given bus, and attach them to the bus. Drivers should use `mdiobus_register()` rather than `__mdiobus_register()` unless they need to pass a specific owner module. MDIO devices which are not PHYs will not be brought up by this function. They are expected to be explicitly listed in DT and instantiated by `of_mdiobus_register()`.

Returns 0 on success or < 0 on error.

```
void mdiobus_free(struct mii_bus * bus)
    free a struct mii_bus
```

## Parameters

**struct mii\_bus \* bus** mii\_bus to free

## Description

This function releases the reference to the underlying device object in the `mii_bus`. If this is the last reference, the `mii_bus` will be freed.

```
struct phy_device * mdiobus_scan(struct mii_bus * bus, int addr)
    scan a bus for MDIO devices.
```

## Parameters

**struct mii\_bus \* bus** mii\_bus to scan

**int addr** address on bus to scan

## Description

This function scans the MDIO bus, looking for devices which can be identified using a vendor/product ID in registers 2 and 3. Not all MDIO devices have such registers, but PHY devices typically do. Hence this function assumes anything found is a PHY, or can be treated as a PHY. Other MDIO devices, such as switches, will probably not be found during the scan.

```
int mdiobus_read_nested(struct mii_bus * bus, int addr, u32 regnum)
    Nested version of the mdiobus_read function
```

## Parameters

**struct mii\_bus \* bus** the mii\_bus struct

**int addr** the phy address

**u32 regnum** register number to read

## Description

In case of nested MDIO bus access avoid lockdep false positives by using `mutex_lock_nested()`.

## NOTE

MUST NOT be called from interrupt context, because the bus read/write functions may wait for an interrupt to conclude the operation.

```
int mdiobus_read(struct mii_bus * bus, int addr, u32 regnum)
    Convenience function for reading a given MII mgmt register
```

## Parameters

**struct mii\_bus \* bus** the mii\_bus struct

**int addr** the phy address

**u32 regnum** register number to read

## NOTE

MUST NOT be called from interrupt context, because the bus read/write functions may wait for an interrupt to conclude the operation.

int **mdiobus\_write\_nested**(struct mii\_bus \* *bus*, int *addr*, u32 *regnum*, u16 *val*)  
Nested version of the mdiobus\_write function

**Parameters**

**struct mii\_bus \* bus** the mii\_bus struct

**int addr** the phy address

**u32 regnum** register number to write

**u16 val** value to write to **regnum**

**Description**

In case of nested MDIO bus access avoid lockdep false positives by using `mutex_lock_nested()`.

**NOTE**

MUST NOT be called from interrupt context, because the bus read/write functions may wait for an interrupt to conclude the operation.

int **mdiobus\_write**(struct mii\_bus \* *bus*, int *addr*, u32 *regnum*, u16 *val*)  
Convenience function for writing a given MII mgmt register

**Parameters**

**struct mii\_bus \* bus** the mii\_bus struct

**int addr** the phy address

**u32 regnum** register number to write

**u16 val** value to write to **regnum**

**NOTE**

MUST NOT be called from interrupt context, because the bus read/write functions may wait for an interrupt to conclude the operation.

void **mdiobus\_release**(struct device \* *d*)  
mii\_bus device release callback

**Parameters**

**struct device \* d** the target struct device that contains the mii\_bus

**Description**

called when the last reference to an mii\_bus is dropped, to free the underlying memory.

int **mdiobus\_create\_device**(struct mii\_bus \* *bus*, struct mdio\_board\_info \* *bi*)  
create a full MDIO device given a mdio\_board\_info structure

**Parameters**

**struct mii\_bus \* bus** MDIO bus to create the devices on

**struct mdio\_board\_info \* bi** mdio\_board\_info structure describing the devices

**Description**

Returns 0 on success or < 0 on error.

int **mdio\_bus\_match**(struct device \* *dev*, struct device\_driver \* *drv*)  
determine if given MDIO driver supports the given MDIO device

**Parameters**

**struct device \* dev** target MDIO device

**struct device\_driver \* drv** given MDIO driver

**Description**

**Given a MDIO device, and a MDIO driver, return 1 if** the driver supports the device. Otherwise, return 0. This may require calling the devices own match function, since different classes of MDIO devices have different match criteria.



## Z8530 PROGRAMMING GUIDE

**Author** Alan Cox

### Introduction

The Z85x30 family synchronous/asynchronous controller chips are used on a large number of cheap network interface cards. The kernel provides a core interface layer that is designed to make it easy to provide WAN services using this chip.

The current driver only support synchronous operation. Merging the asynchronous driver support into this code to allow any Z85x30 device to be used as both a tty interface and as a synchronous controller is a project for Linux post the 2.4 release

### Driver Modes

The Z85230 driver layer can drive Z8530, Z85C30 and Z85230 devices in three different modes. Each mode can be applied to an individual channel on the chip (each chip has two channels).

The PIO synchronous mode supports the most common Z8530 wiring. Here the chip is interface to the I/O and interrupt facilities of the host machine but not to the DMA subsystem. When running PIO the Z8530 has extremely tight timing requirements. Doing high speeds, even with a Z85230 will be tricky. Typically you should expect to achieve at best 9600 baud with a Z8C530 and 64Kbits with a Z85230.

The DMA mode supports the chip when it is configured to use dual DMA channels on an ISA bus. The better cards tend to support this mode of operation for a single channel. With DMA running the Z85230 tops out when it starts to hit ISA DMA constraints at about 512Kbits. It is worth noting here that many PC machines hang or crash when the chip is driven fast enough to hold the ISA bus solid.

Transmit DMA mode uses a single DMA channel. The DMA channel is used for transmission as the transmit FIFO is smaller than the receive FIFO. it gives better performance than pure PIO mode but is nowhere near as ideal as pure DMA mode.

### Using the Z85230 driver

The Z85230 driver provides the back end interface to your board. To configure a Z8530 interface you need to detect the board and to identify its ports and interrupt resources. It is also your problem to verify the resources are available.

Having identified the chip you need to fill in a struct `z8530_dev`, which describes each chip. This object must exist until you finally shutdown the board. Firstly zero the active field. This ensures nothing goes off without you intending it. The `irq` field should be set to the interrupt number of the chip. (Each chip has a single interrupt source rather than each channel). You are responsible for allocating the interrupt line. The interrupt handler should be set to `z8530_interrupt()`. The device id should be set to the

z8530\_dev structure pointer. Whether the interrupt can be shared or not is board dependent, and up to you to initialise.

The structure holds two channel structures. Initialise chanA.ctrlrio and chanA.dataio with the address of the control and data ports. You can or this with Z8530\_PORT\_SLEEP to indicate your interface needs the 5uS delay for chip settling done in software. The PORT\_SLEEP option is architecture specific. Other flags may become available on future platforms, eg for MMIO. Initialise the chanA.irqs to &z8530\_nop to start the chip up as disabled and discarding interrupt events. This ensures that stray interrupts will be mopped up and not hang the bus. Set chanA.dev to point to the device structure itself. The private and name field you may use as you wish. The private field is unused by the Z85230 layer. The name is used for error reporting and it may thus make sense to make it match the network name.

Repeat the same operation with the B channel if your chip has both channels wired to something useful. This isn't always the case. If it is not wired then the I/O values do not matter, but you must initialise chanB.dev.

If your board has DMA facilities then initialise the txdma and rxdma fields for the relevant channels. You must also allocate the ISA DMA channels and do any necessary board level initialisation to configure them. The low level driver will do the Z8530 and DMA controller programming but not board specific magic.

Having initialised the device you can then call `z8530_init()`. This will probe the chip and reset it into a known state. An identification sequence is then run to identify the chip type. If the checks fail to pass the function returns a non zero error code. Typically this indicates that the port given is not valid. After this call the type field of the z8530\_dev structure is initialised to either Z8530, Z85C30 or Z85230 according to the chip found.

Once you have called `z8530_init` you can also make use of the utility function `z8530_describe()`. This provides a consistent reporting format for the Z8530 devices, and allows all the drivers to provide consistent reporting.

## Attaching Network Interfaces

If you wish to use the network interface facilities of the driver, then you need to attach a network device to each channel that is present and in use. In addition to use the generic HDLC you need to follow some additional plumbing rules. They may seem complex but a look at the example `hostess_sv11` driver should reassure you.

The network device used for each channel should be pointed to by the `netdevice` field of each channel. The `hdlc->priv` field of the network device points to your private data - you will need to be able to find your private data from this.

The way most drivers approach this particular problem is to create a structure holding the Z8530 device definition and put that into the private field of the network device. The network device fields of the channels then point back to the network devices.

If you wish to use the generic HDLC then you need to register the HDLC device.

Before you register your network device you will also need to provide suitable handlers for most of the network device callbacks. See the network device documentation for more details on this.

## Configuring And Activating The Port

The Z85230 driver provides helper functions and tables to load the port registers on the Z8530 chips. When programming the register settings for a channel be aware that the documentation recommends initialisation orders. Strange things happen when these are not followed.

`z8530_channel_load()` takes an array of pairs of initialisation values in an array of `u8` type. The first value is the Z8530 register number. Add 16 to indicate the alternate register bank on the later chips. The array is terminated by a 255.

The driver provides a pair of public tables. The `z8530_hdlc_kilostream` table is for the UK 'Kilostream' service and also happens to cover most other end host configurations. The `z8530_hdlc_kilostream_85230` table is the same configuration using the enhancements of the 85230 chip. The configuration loaded is standard NRZ encoded synchronous data with HDLC bitstuffing. All of the timing is taken from the other end of the link.

When writing your own tables be aware that the driver internally tracks register values. It may need to reload values. You should therefore be sure to set registers 1-7, 9-11, 14 and 15 in all configurations. Where the register settings depend on DMA selection the driver will update the bits itself when you open or close. Loading a new table with the interface open is not recommended.

There are three standard configurations supported by the core code. In PIO mode the interface is programmed up to use interrupt driven PIO. This places high demands on the host processor to avoid latency. The driver is written to take account of latency issues but it cannot avoid latencies caused by other drivers, notably IDE in PIO mode. Because the drivers allocate buffers you must also prevent MTU changes while the port is open.

Once the port is open it will call the `rx_function` of each channel whenever a completed packet arrived. This is invoked from interrupt context and passes you the channel and a network buffer (`struct sk_buff`) holding the data. The data includes the CRC bytes so most users will want to trim the last two bytes before processing the data. This function is very timing critical. When you wish to simply discard data the support code provides the function `z8530_null_rx()` to discard the data.

To active PIO mode sending and receiving the `z8530_sync_open` is called. This expects to be passed the network device and the channel. Typically this is called from your network device open callback. On a failure a non zero error status is returned. The `z8530_sync_close()` function shuts down a PIO channel. This must be done before the channel is opened again and before the driver shuts down and unloads.

The ideal mode of operation is dual channel DMA mode. Here the kernel driver will configure the board for DMA in both directions. The driver also handles ISA DMA issues such as controller programming and the memory range limit for you. This mode is activated by calling the `z8530_sync_dma_open()` function. On failure a non zero error value is returned. Once this mode is activated it can be shut down by calling the `z8530_sync_dma_close()`. You must call the close function matching the open mode you used.

The final supported mode uses a single DMA channel to drive the transmit side. As the Z85C30 has a larger FIFO on the receive channel this tends to increase the maximum speed a little. This is activated by calling the `z8530_sync_txdma_open`. This returns a non zero error code on failure. The `z8530_sync_txdma_close()` function closes down the Z8530 interface from this mode.

## Network Layer Functions

The Z8530 layer provides functions to queue packets for transmission. The driver internally buffers the frame currently being transmitted and one further frame (in order to keep back to back transmission running). Any further buffering is up to the caller.

The function `z8530_queue_xmit()` takes a network buffer in `sk_buff` format and queues it for transmission. The caller must provide the entire packet with the exception of the bitstuffing and CRC. This is normally done by the caller via the generic HDLC interface layer. It returns 0 if the buffer has been queued and non zero values for queue full. If the function accepts the buffer it becomes property of the Z8530 layer and the caller should not free it.

The function `z8530_get_stats()` returns a pointer to an internally maintained per interface statistics block. This provides most of the interface code needed to implement the network layer `get_stats` callback.

## Porting The Z8530 Driver

The Z8530 driver is written to be portable. In DMA mode it makes assumptions about the use of ISA DMA. These are probably warranted in most cases as the Z85230 in particular was designed to glue to PC type machines. The PIO mode makes no real assumptions.

Should you need to retarget the Z8530 driver to another architecture the only code that should need changing are the port I/O functions. At the moment these assume PC I/O port accesses. This may not be appropriate for all platforms. Replacing `z8530_read_port()` and `z8530_write_port` is intended to be all that is required to port this driver layer.

## Known Bugs And Assumptions

**Interrupt Locking** The locking in the driver is done via the global cli/sti lock. This makes for relatively poor SMP performance. Switching this to use a per device spin lock would probably materially improve performance.

**Occasional Failures** We have reports of occasional failures when run for very long periods of time and the driver starts to receive junk frames. At the moment the cause of this is not clear.

## Public Functions Provided

`irqreturn_t z8530_interrupt(int irq, void * dev_id)`  
Handle an interrupt from a Z8530

### Parameters

**int irq** Interrupt number

**void \* dev\_id** The Z8530 device that is interrupting.

### Description

A Z85[2]30 device has stuck its hand in the air for attention. We scan both the channels on the chip for events and then call the channel specific call backs for each channel that has events. We have to use callback functions because the two channels can be in different modes.

Locking is done for the handlers. Note that locking is done at the chip level (the 5uS delay issue is per chip not per channel). `c->lock` for both channels points to `dev->lock`

`int z8530_sync_open(struct net_device * dev, struct z8530_channel * c)`  
Open a Z8530 channel for PIO

### Parameters

**struct net\_device \* dev** The network interface we are using

**struct z8530\_channel \* c** The Z8530 channel to open in synchronous PIO mode

### Description

Switch a Z8530 into synchronous mode without DMA assist. We raise the RTS/DTR and commence network operation.

`int z8530_sync_close(struct net_device * dev, struct z8530_channel * c)`  
Close a PIO Z8530 channel

### Parameters

**struct net\_device \* dev** Network device to close

**struct z8530\_channel \* c** Z8530 channel to disassociate and move to idle

### Description

Close down a Z8530 interface and switch its interrupt handlers to discard future events.

`int z8530_sync_dma_open(struct net_device * dev, struct z8530_channel * c)`  
Open a Z8530 for DMA I/O

### Parameters

**struct net\_device \* dev** The network device to attach

**struct z8530\_channel \* c** The Z8530 channel to configure in sync DMA mode.

#### Description

Set up a Z85x30 device for synchronous DMA in both directions. Two ISA DMA channels must be available for this to work. We assume ISA DMA driven I/O and PC limits on access.

int **z8530\_sync\_dma\_close**(struct *net\_device* \* dev, struct z8530\_channel \* c)  
Close down DMA I/O

#### Parameters

**struct net\_device \* dev** Network device to detach

**struct z8530\_channel \* c** Z8530 channel to move into discard mode

#### Description

Shut down a DMA mode synchronous interface. Halt the DMA, and free the buffers.

int **z8530\_sync\_txdma\_open**(struct *net\_device* \* dev, struct z8530\_channel \* c)  
Open a Z8530 for TX driven DMA

#### Parameters

**struct net\_device \* dev** The network device to attach

**struct z8530\_channel \* c** The Z8530 channel to configure in sync DMA mode.

#### Description

Set up a Z85x30 device for synchronous DMA transmission. One ISA DMA channel must be available for this to work. The receive side is run in PIO mode, but then it has the bigger FIFO.

int **z8530\_sync\_txdma\_close**(struct *net\_device* \* dev, struct z8530\_channel \* c)  
Close down a TX driven DMA channel

#### Parameters

**struct net\_device \* dev** Network device to detach

**struct z8530\_channel \* c** Z8530 channel to move into discard mode

#### Description

Shut down a DMA/PIO split mode synchronous interface. Halt the DMA, and free the buffers.

void **z8530\_describe**(struct z8530\_dev \* dev, char \* *mapping*, unsigned long *io*)  
Uniformly describe a Z8530 port

#### Parameters

**struct z8530\_dev \* dev** Z8530 device to describe

**char \* mapping** string holding mapping type (eg "I/O" or "Mem")

**unsigned long io** the port value in question

#### Description

Describe a Z8530 in a standard format. We must pass the I/O as the port offset isn't predictable. The main reason for this function is to try and get a common format of report.

int **z8530\_init**(struct z8530\_dev \* dev)  
Initialise a Z8530 device

#### Parameters

**struct z8530\_dev \* dev** Z8530 device to initialise.

#### Description

Configure up a Z8530/Z85C30 or Z85230 chip. We check the device is present, identify the type and then program it to hopefully keep quite and behave. This matters a lot, a Z8530 in the wrong state will sometimes get into stupid modes generating 10Khz interrupt streams and the like.

We set the interrupt handler up to discard any events, in case we get them during reset or setp.

Return 0 for success, or a negative value indicating the problem in errno form.

int **z8530\_shutdown**(struct z8530\_dev \* *dev*)  
Shutdown a Z8530 device

#### Parameters

**struct z8530\_dev \* dev** The Z8530 chip to shutdown

#### Description

We set the interrupt handlers to silence any interrupts. We then reset the chip and wait 100uS to be sure the reset completed. Just in case the caller then tries to do stuff.

This is called without the lock held

int **z8530\_channel\_load**(struct z8530\_channel \* *c*, u8 \* *rtable*)  
Load channel data

#### Parameters

**struct z8530\_channel \* c** Z8530 channel to configure

**u8 \* rtable** table of register, value pairs FIXME: ioctl to allow user uploaded tables

#### Description

Load a Z8530 channel up from the system data. We use +16 to indicate the “prime” registers. The value 255 terminates the table.

void **z8530\_null\_rx**(struct z8530\_channel \* *c*, struct *sk\_buff* \* *skb*)  
Discard a packet

#### Parameters

**struct z8530\_channel \* c** The channel the packet arrived on

**struct sk\_buff \* skb** The buffer

#### Description

We point the receive handler at this function when idle. Instead of processing the frames we get to throw them away.

netdev\_tx\_t **z8530\_queue\_xmit**(struct z8530\_channel \* *c*, struct *sk\_buff* \* *skb*)  
Queue a packet

#### Parameters

**struct z8530\_channel \* c** The channel to use

**struct sk\_buff \* skb** The packet to kick down the channel

#### Description

Queue a packet for transmission. Because we have rather hard to hit interrupt latencies for the Z85230 per packet even in DMA mode we do the flip to DMA buffer if needed here not in the IRQ.

Called from the network code. The lock is not held at this point.

## Internal Functions

int **z8530\_read\_port**(unsigned long *p*)  
Architecture specific interface function

### Parameters

unsigned long *p* port to read

### Description

Provided port access methods. The Control SV11 requires no delays between accesses and uses PC I/O. Some drivers may need a 5uS delay

In the longer term this should become an architecture specific section so that this can become a generic driver interface for all platforms. For now we only handle PC I/O ports with or without the dread 5uS sanity delay.

The caller must hold sufficient locks to avoid violating the horrible 5uS delay rule.

void **z8530\_write\_port**(unsigned long *p*, u8 *d*)  
Architecture specific interface function

### Parameters

unsigned long *p* port to write

u8 *d* value to write

### Description

Write a value to a port with delays if need be. Note that the caller must hold locks to avoid read/writes from other contexts violating the 5uS rule

In the longer term this should become an architecture specific section so that this can become a generic driver interface for all platforms. For now we only handle PC I/O ports with or without the dread 5uS sanity delay.

u8 **read\_zsreg**(struct z8530\_channel \* *c*, u8 *reg*)  
Read a register from a Z85230

### Parameters

struct z8530\_channel \* *c* Z8530 channel to read from (2 per chip)

u8 *reg* Register to read FIXME: Use a spinlock.

Most of the Z8530 registers are indexed off the control registers. A read is done by writing to the control register and reading the register back. The caller must hold the lock

u8 **read\_zsdata**(struct z8530\_channel \* *c*)  
Read the data port of a Z8530 channel

### Parameters

struct z8530\_channel \* *c* The Z8530 channel to read the data port from

### Description

The data port provides fast access to some things. We still have all the 5uS delays to worry about.

void **write\_zsreg**(struct z8530\_channel \* *c*, u8 *reg*, u8 *val*)  
Write to a Z8530 channel register

### Parameters

struct z8530\_channel \* *c* The Z8530 channel

u8 *reg* Register number

**u8 val** Value to write

### Description

Write a value to an indexed register. The caller must hold the lock to honour the irritating delay rules. We know about register 0 being fast to access.

Assumes c->lock is held.

void **write\_zsctrl**(struct z8530\_channel \* c, u8 val)  
Write to a Z8530 control register

### Parameters

**struct z8530\_channel \* c** The Z8530 channel

**u8 val** Value to write

### Description

Write directly to the control register on the Z8530

void **write\_zsdata**(struct z8530\_channel \* c, u8 val)  
Write to a Z8530 control register

### Parameters

**struct z8530\_channel \* c** The Z8530 channel

**u8 val** Value to write

### Description

Write directly to the data register on the Z8530

void **z8530\_flush\_fifo**(struct z8530\_channel \* c)  
Flush on chip RX FIFO

### Parameters

**struct z8530\_channel \* c** Channel to flush

### Description

Flush the receive FIFO. There is no specific option for this, we blindly read bytes and discard them. Reading when there is no data is harmless. The 8530 has a 4 byte FIFO, the 85230 has 8 bytes.

All locking is handled for the caller. On return data may still be present if it arrived during the flush.

void **z8530\_rtsdtr**(struct z8530\_channel \* c, int set)  
Control the outgoing DTS/RTS line

### Parameters

**struct z8530\_channel \* c** The Z8530 channel to control;

**int set** 1 to set, 0 to clear

### Description

Sets or clears DTR/RTS on the requested line. All locking is handled by the caller. For now we assume all boards use the actual RTS/DTR on the chip. Apparently one or two don't. We'll scream about them later.

void **z8530\_rx**(struct z8530\_channel \* c)  
Handle a PIO receive event

### Parameters

**struct z8530\_channel \* c** Z8530 channel to process

### Description



Receive handler for receiving in PIO mode. This is much like the async one but not quite the same or as complex

#### Note

**Its intended that this handler can easily be separated from** the main code to run realtime. That'll be needed for some machines (eg to ever clock 64kbits on a sparc ;)).

The RT\_LOCK macros don't do anything now. Keep the code covered by them as short as possible in all circumstances - clocks cost baud. The interrupt handler is assumed to be atomic w.r.t. to other code - this is true in the RT case too.

We only cover the sync cases for this. If you want 2Mbit async do it yourself but consider medical assistance first. This non DMA synchronous mode is portable code. The DMA mode assumes PCI like ISA DMA

Called with the device lock held

```
void z8530_tx(struct z8530_channel * c)  
    Handle a PIO transmit event
```

#### Parameters

**struct z8530\_channel \* c** Z8530 channel to process

#### Description

Z8530 transmit interrupt handler for the PIO mode. The basic idea is to attempt to keep the FIFO fed. We fill as many bytes in as possible, its quite possible that we won't keep up with the data rate otherwise.

```
void z8530_status(struct z8530_channel * chan)  
    Handle a PIO status exception
```

#### Parameters

**struct z8530\_channel \* chan** Z8530 channel to process

#### Description

A status event occurred in PIO synchronous mode. There are several reasons the chip will bother us here. A transmit underrun means we failed to feed the chip fast enough and just broke a packet. A DCD change is a line up or down.

```
void z8530_dma_rx(struct z8530_channel * chan)  
    Handle a DMA RX event
```

#### Parameters

**struct z8530\_channel \* chan** Channel to handle

#### Description

Non bus mastering DMA interfaces for the Z8x30 devices. This is really pretty PC specific. The DMA mode means that most receive events are handled by the DMA hardware. We get a kick here only if a frame ended.

```
void z8530_dma_tx(struct z8530_channel * chan)  
    Handle a DMA TX event
```

#### Parameters

**struct z8530\_channel \* chan** The Z8530 channel to handle

#### Description

We have received an interrupt while doing DMA transmissions. It shouldn't happen. Scream loudly if it does.

```
void z8530_dma_status(struct z8530_channel * chan)  
    Handle a DMA status exception
```

### Parameters

**struct z8530\_channel \* chan** Z8530 channel to process

A status event occurred on the Z8530. We receive these for two reasons when in DMA mode. Firstly if we finished a packet transfer we get one and kick the next packet out. Secondly we may see a DCD change.

void **z8530\_rx\_clear**(struct z8530\_channel \* c)  
Handle RX events from a stopped chip

### Parameters

**struct z8530\_channel \* c** Z8530 channel to shut up

### Description

Receive interrupt vectors for a Z8530 that is in 'parked' mode. For machines with PCI Z85x30 cards, or level triggered interrupts (eg the MacII) we must clear the interrupt cause or die.

void **z8530\_tx\_clear**(struct z8530\_channel \* c)  
Handle TX events from a stopped chip

### Parameters

**struct z8530\_channel \* c** Z8530 channel to shut up

### Description

Transmit interrupt vectors for a Z8530 that is in 'parked' mode. For machines with PCI Z85x30 cards, or level triggered interrupts (eg the MacII) we must clear the interrupt cause or die.

void **z8530\_status\_clear**(struct z8530\_channel \* chan)  
Handle status events from a stopped chip

### Parameters

**struct z8530\_channel \* chan** Z8530 channel to shut up

### Description

Status interrupt vectors for a Z8530 that is in 'parked' mode. For machines with PCI Z85x30 cards, or level triggered interrupts (eg the MacII) we must clear the interrupt cause or die.

void **z8530\_tx\_begin**(struct z8530\_channel \* c)  
Begin packet transmission

### Parameters

**struct z8530\_channel \* c** The Z8530 channel to kick

### Description

This is the speed sensitive side of transmission. If we are called and no buffer is being transmitted we commence the next buffer. If nothing is queued we idle the sync.

### Note

**We are handling this code path in the interrupt path, keep it** fast or bad things will happen.

Called with the lock held.

void **z8530\_tx\_done**(struct z8530\_channel \* c)  
TX complete callback

### Parameters

**struct z8530\_channel \* c** The channel that completed a transmit.

### Description

This is called when we complete a packet send. We wake the queue, start the next packet going and then free the buffer of the existing packet. This code is fairly timing sensitive.

Called with the register lock held.

```
void z8530_rx_done(struct z8530_channel * c)  
    Receive completion callback
```

#### Parameters

**struct z8530\_channel \* c** The channel that completed a receive

#### Description

A new packet is complete. Our goal here is to get back into receive mode as fast as possible. On the Z85230 we could change to using ESCC mode, but on the older chips we have no choice. We flip to the new buffer immediately in DMA mode so that the DMA of the next frame can occur while we are copying the previous buffer to an `sk_buff`

Called with the lock held

```
int spans_boundary(struct sk_buff * skb)  
    Check a packet can be ISA DMA'd
```

#### Parameters

**struct sk\_buff \* skb** The buffer to check

#### Description

Returns true if the buffer cross a DMA boundary on a PC. The poor thing can only DMA within a 64K block not across the edges of it.



## Symbols

- `__alloc_skb` (C function), 32
- `__dev_alloc_page` (C function), 18
- `__dev_alloc_pages` (C function), 17
- `__dev_get_by_flags` (C function), 80
- `__dev_get_by_index` (C function), 79
- `__dev_get_by_name` (C function), 79
- `__dev_mc_sync` (C function), 110
- `__dev_mc_unsync` (C function), 110
- `__dev_remove_pack` (C function), 77
- `__dev_uc_sync` (C function), 109
- `__dev_uc_unsync` (C function), 110
- `kfree_skb` (C function), 33
- `mdiobus_register` (C function), 121
- `napi_alloc_skb` (C function), 33
- `napi_schedule` (C function), 85
- `napi_schedule_irqoff` (C function), 86
- `netdev_alloc_skb` (C function), 32
- `netif_subqueue_stopped` (C function), 108
- `pskb_copy_fclone` (C function), 35
- `pskb_pull_tail` (C function), 37
- `sk_mem_raise_allocated` (C function), 45
- `sk_mem_reclaim` (C function), 46
- `sk_mem_reduce_allocated` (C function), 46
- `sk_mem_schedule` (C function), 46
- `skb_fill_page_desc` (C function), 15
- `skb_frag_ref` (C function), 18
- `skb_frag_set_page` (C function), 19
- `skb_frag_unref` (C function), 18
- `skb_gso_segment` (C function), 84
- `skb_header_release` (C function), 12
- `skb_pad` (C function), 35
- `skb_put_padto` (C function), 20
- `skb_queue_after` (C function), 14
- `skb_queue_head_init` (C function), 13
- `skb_try_rcv_datagram` (C function), 46
- `copy_from_pages` (C function), 55

## A

- `alloc_etherdev_mqs` (C function), 96
- `alloc_netdev_mqs` (C function), 92
- `alloc_skb_with_frags` (C function), 44

## B

- `bpf_prog_create` (C function), 49
- `bpf_prog_create_from_user` (C function), 49

## C

- `call_netdevice_notifiers` (C function), 82
- `compare_ether_header` (C function), 100
- `consume_skb` (C function), 33
- `csum_partial_copy_to_xdr` (C function), 64

## D

- `datagram_poll` (C function), 48
- `dev_add_offload` (C function), 78
- `dev_add_pack` (C function), 77
- `dev_alloc_name` (C function), 81
- `dev_change_carrier` (C function), 90
- `dev_change_flags` (C function), 89
- `dev_change_net_namespace` (C function), 93
- `dev_change_proto_down` (C function), 91
- `dev_close` (C function), 82
- `dev_disable_lro` (C function), 82
- `dev_fill_metadata_dst` (C function), 78
- `dev_forward_skb` (C function), 83
- `dev_get_by_index` (C function), 80
- `dev_get_by_index_rcu` (C function), 80
- `dev_get_by_name` (C function), 79
- `dev_get_by_name_rcu` (C function), 79
- `dev_get_by_napi_id` (C function), 80
- `dev_get_flags` (C function), 89
- `dev_get_iflink` (C function), 78
- `dev_get_phys_port_id` (C function), 90
- `dev_get_phys_port_name` (C function), 90
- `dev_get_stats` (C function), 92
- `dev_getbyhwaddr_rcu` (C function), 80
- `dev_hold` (C function), 108
- `dev_loopback_xmit` (C function), 84
- `dev_open` (C function), 81
- `dev_put` (C function), 108
- `dev_remove_offload` (C function), 78
- `dev_remove_pack` (C function), 78
- `dev_set_allmulti` (C function), 89
- `dev_set_group` (C function), 90
- `dev_set_mac_address` (C function), 90
- `dev_set_mtu` (C function), 90
- `dev_set_promiscuity` (C function), 89
- `dev_valid_name` (C function), 81
- `devm_mdiobus_alloc_size` (C function), 121
- `devm_mdiobus_free` (C function), 121

## E

`eth_addr_dec` (C function), [100](#)  
`eth_broadcast_addr` (C function), [98](#)  
`eth_change_mtu` (C function), [96](#)  
`eth_commit_mac_addr_change` (C function), [95](#)  
`eth_get_headlen` (C function), [94](#)  
`eth_header` (C function), [94](#)  
`eth_header_cache` (C function), [95](#)  
`eth_header_cache_update` (C function), [95](#)  
`eth_header_parse` (C function), [95](#)  
`eth_hw_addr_inherit` (C function), [99](#)  
`eth_hw_addr_random` (C function), [98](#)  
`eth_mac_addr` (C function), [95](#)  
`eth_prepare_mac_addr_change` (C function), [95](#)  
`eth_proto_is_802_3` (C function), [98](#)  
`eth_random_addr` (C function), [98](#)  
`eth_skb_pad` (C function), [100](#)  
`eth_type_trans` (C function), [94](#)  
`eth_zero_addr` (C function), [98](#)  
`ether_addr_copy` (C function), [98](#)  
`ether_addr_equal` (C function), [99](#)  
`ether_addr_equal_64bits` (C function), [99](#)  
`ether_addr_equal_masked` (C function), [99](#)  
`ether_addr_equal_unaligned` (C function), [99](#)  
`ether_addr_to_u64` (C function), [100](#)  
`ether_setup` (C function), [96](#)

## F

`free_netdev` (C function), [92](#)

## G

`gen_estimator_active` (C function), [54](#)  
`gen_kill_estimator` (C function), [54](#)  
`gen_new_estimator` (C function), [53](#)  
`gen_replace_estimator` (C function), [54](#)  
`genphy_aneg_done` (C function), [118](#)  
`genphy_config_advert` (C function), [120](#)  
`genphy_config_aneg` (C function), [118](#)  
`genphy_config_eee_advert` (C function), [120](#)  
`genphy_read_status` (C function), [118](#)  
`genphy_restart_aneg` (C function), [118](#)  
`genphy_setup_forced` (C function), [118](#)  
`genphy_soft_reset` (C function), [119](#)  
`genphy_update_link` (C function), [118](#)  
`get_phy_c45_ids` (C function), [119](#)  
`get_phy_device` (C function), [116](#)  
`get_phy_id` (C function), [119](#)  
`gnet_estimator` (C type), [51](#)  
`gnet_stats_basic` (C type), [50](#)  
`gnet_stats_copy_app` (C function), [53](#)  
`gnet_stats_copy_basic` (C function), [52](#)  
`gnet_stats_copy_queue` (C function), [52](#)  
`gnet_stats_copy_rate_est` (C function), [52](#)  
`gnet_stats_finish_copy` (C function), [53](#)  
`gnet_stats_queue` (C type), [51](#)  
`gnet_stats_rate_est` (C type), [50](#)  
`gnet_stats_rate_est64` (C type), [50](#)

`gnet_stats_start_copy` (C function), [52](#)  
`gnet_stats_start_copy_compat` (C function), [51](#)

## I

`init_dummy_netdev` (C function), [92](#)  
`is_broadcast_ether_addr` (C function), [97](#)  
`is_etherdev_addr` (C function), [100](#)  
`is_link_local_ether_addr` (C function), [96](#)  
`is_local_ether_addr` (C function), [97](#)  
`is_multicast_ether_addr` (C function), [97](#)  
`is_unicast_ether_addr` (C function), [97](#)  
`is_valid_ether_addr` (C function), [97](#)  
`is_zero_ether_addr` (C function), [97](#)

## K

`kernel_recvmsg` (C function), [31](#)  
`kfree_skb` (C function), [33](#)

## L

`lock_sock_fast` (C function), [46](#)

## M

`mdio_bus_match` (C function), [123](#)  
`mdiobus_alloc_size` (C function), [121](#)  
`mdiobus_create_device` (C function), [123](#)  
`mdiobus_free` (C function), [122](#)  
`mdiobus_read` (C function), [122](#)  
`mdiobus_read_nested` (C function), [122](#)  
`mdiobus_release` (C function), [123](#)  
`mdiobus_scan` (C function), [122](#)  
`mdiobus_write` (C function), [123](#)  
`mdiobus_write_nested` (C function), [122](#)

## N

`napi_complete` (C function), [101](#)  
`napi_disable` (C function), [101](#)  
`napi_enable` (C function), [101](#)  
`napi_hash_del` (C function), [101](#)  
`napi_schedule` (C function), [101](#)  
`napi_schedule_irqoff` (C function), [101](#)  
`napi_schedule_prep` (C function), [86](#)  
`napi_synchronize` (C function), [102](#)  
`net_device` (C type), [103](#)  
`netdev_alloc_frag` (C function), [32](#)  
`netdev_alloc_skb` (C function), [17](#)  
`netdev_bonding_info_change` (C function), [88](#)  
`netdev_boot_setup_check` (C function), [78](#)  
`netdev_cap_txqueue` (C function), [107](#)  
`netdev_change_features` (C function), [91](#)  
`netdev_completed_queue` (C function), [106](#)  
`netdev_features_change` (C function), [81](#)  
`netdev_has_any_upper_dev` (C function), [86](#)  
`netdev_has_upper_dev` (C function), [86](#)  
`netdev_has_upper_dev_all_rcu` (C function), [86](#)  
`netdev_increment_features` (C function), [94](#)  
`netdev_is_rx_handler_busy` (C function), [85](#)  
`netdev_lower_get_first_private_rcu` (C function), [87](#)

[netdev\\_lower\\_get\\_next](#) (C function), 87  
[netdev\\_lower\\_get\\_next\\_private](#) (C function), 87  
[netdev\\_lower\\_get\\_next\\_private\\_rcu](#) (C function), 87  
[netdev\\_lower\\_state\\_changed](#) (C function), 89  
[netdev\\_master\\_upper\\_dev\\_get](#) (C function), 86  
[netdev\\_master\\_upper\\_dev\\_get\\_rcu](#) (C function), 88  
[netdev\\_master\\_upper\\_dev\\_link](#) (C function), 88  
[netdev\\_notify\\_peers](#) (C function), 81  
[netdev\\_priv](#) (C function), 105  
[netdev\\_priv\\_flags](#) (C type), 102  
[netdev\\_reset\\_queue](#) (C function), 107  
[netdev\\_rx\\_handler\\_register](#) (C function), 85  
[netdev\\_rx\\_handler\\_unregister](#) (C function), 85  
[netdev\\_sent\\_queue](#) (C function), 106  
[netdev\\_state\\_change](#) (C function), 81  
[netdev\\_txq\\_bql\\_complete\\_prefetchw](#) (C function), 106  
[netdev\\_txq\\_bql\\_enqueue\\_prefetchw](#) (C function), 106  
[netdev\\_update\\_features](#) (C function), 91  
[netdev\\_upper\\_dev\\_link](#) (C function), 88  
[netdev\\_upper\\_dev\\_unlink](#) (C function), 88  
[netdev\\_upper\\_get\\_next\\_dev\\_rcu](#) (C function), 87  
[netif\\_carrier\\_off](#) (C function), 96  
[netif\\_carrier\\_ok](#) (C function), 108  
[netif\\_carrier\\_on](#) (C function), 96  
[netif\\_device\\_attach](#) (C function), 83  
[netif\\_device\\_detach](#) (C function), 83  
[netif\\_device\\_present](#) (C function), 109  
[netif\\_dormant](#) (C function), 109  
[netif\\_dormant\\_off](#) (C function), 109  
[netif\\_dormant\\_on](#) (C function), 108  
[netif\\_get\\_num\\_default\\_rss\\_queues](#) (C function), 83  
[netif\\_is\\_multiqueue](#) (C function), 108  
[netif\\_napi\\_add](#) (C function), 105  
[netif\\_napi\\_del](#) (C function), 105  
[netif\\_oper\\_up](#) (C function), 109  
[netif\\_queue\\_stopped](#) (C function), 106  
[netif\\_receive\\_skb](#) (C function), 85  
[netif\\_running](#) (C function), 107  
[netif\\_rx](#) (C function), 84  
[netif\\_set\\_real\\_num\\_rx\\_queues](#) (C function), 83  
[netif\\_stacked\\_transfer\\_operstate](#) (C function), 91  
[netif\\_start\\_queue](#) (C function), 105  
[netif\\_start\\_subqueue](#) (C function), 107  
[netif\\_stop\\_queue](#) (C function), 106  
[netif\\_stop\\_subqueue](#) (C function), 107  
[netif\\_tx\\_lock](#) (C function), 109  
[netif\\_tx\\_napi\\_add](#) (C function), 105  
[netif\\_wake\\_queue](#) (C function), 105  
[netif\\_wake\\_subqueue](#) (C function), 108  
[phy\\_attach\\_direct](#) (C function), 117  
[phy\\_change](#) (C function), 115  
[phy\\_change\\_work](#) (C function), 115  
[phy\\_check\\_valid](#) (C function), 114  
[phy\\_clear\\_interrupt](#) (C function), 113  
[phy\\_config\\_interrupt](#) (C function), 113  
[phy\\_connect](#) (C function), 117  
[phy\\_connect\\_direct](#) (C function), 116  
[phy\\_detach](#) (C function), 118  
[phy\\_device\\_register](#) (C function), 116  
[phy\\_device\\_remove](#) (C function), 116  
[phy\\_disable\\_interrupts](#) (C function), 115  
[phy\\_disconnect](#) (C function), 117  
[phy\\_driver\\_register](#) (C function), 119  
[phy\\_enable\\_interrupts](#) (C function), 115  
[phy\\_error](#) (C function), 115  
[phy\\_ethtool\\_get\\_eee](#) (C function), 113  
[phy\\_ethtool\\_set\\_eee](#) (C function), 113  
[phy\\_ethtool\\_sset](#) (C function), 111  
[phy\\_find\\_first](#) (C function), 116  
[phy\\_find\\_valid](#) (C function), 113  
[phy\\_get\\_eee\\_err](#) (C function), 113  
[phy\\_init\\_eee](#) (C function), 112  
[phy\\_interrupt](#) (C function), 115  
[phy\\_mac\\_interrupt](#) (C function), 112  
[phy\\_mii\\_ioctl](#) (C function), 111  
[phy\\_poll\\_reset](#) (C function), 120  
[phy\\_prepare\\_link](#) (C function), 120  
[phy\\_print\\_status](#) (C function), 110  
[phy\\_probe](#) (C function), 120  
[phy\\_register\\_fixup](#) (C function), 116  
[phy\\_restart\\_aneg](#) (C function), 110  
[phy\\_sanitize\\_settings](#) (C function), 114  
[phy\\_start](#) (C function), 112  
[phy\\_start\\_aneg](#) (C function), 111  
[phy\\_start\\_aneg\\_priv](#) (C function), 114  
[phy\\_start\\_interrupts](#) (C function), 112  
[phy\\_start\\_machine](#) (C function), 111  
[phy\\_state\\_machine](#) (C function), 115  
[phy\\_stop](#) (C function), 112  
[phy\\_stop\\_interrupts](#) (C function), 112  
[phy\\_stop\\_machine](#) (C function), 115  
[phy\\_supported\\_speeds](#) (C function), 114  
[phy\\_trigger\\_machine](#) (C function), 114  
[phy\\_unregister\\_fixup](#) (C function), 116  
[pskb\\_expand\\_head](#) (C function), 35  
[pskb\\_put](#) (C function), 36  
[pskb\\_trim\\_rcsum](#) (C function), 22  
[pskb\\_trim\\_unique](#) (C function), 16

## R

[read\\_zsdata](#) (C function), 131  
[read\\_zsreg](#) (C function), 131  
[register\\_netdev](#) (C function), 92  
[register\\_netdevice](#) (C function), 91  
[register\\_netdevice\\_notifier](#) (C function), 82  
[rpc\\_add\\_pipe\\_dir\\_object](#) (C function), 65  
[rpc\\_alloc\\_iostats](#) (C function), 64

## O

[of\\_mdio\\_find\\_bus](#) (C function), 121

## P

[phy\\_aneg\\_done](#) (C function), 111  
[phy\\_attach](#) (C function), 117

`rpc_bind_new_program` (C function), 67  
`rpc_call_async` (C function), 68  
`rpc_call_sync` (C function), 68  
`rpc_clnt_add_xprt` (C function), 70  
`rpc_clnt_iterate_for_each_xprt` (C function), 67  
`rpc_clnt_setup_test_and_add_xprt` (C function), 69  
`rpc_clnt_test_and_add_xprt` (C function), 69  
`rpc_clone_client` (C function), 66  
`rpc_clone_client_set_auth` (C function), 67  
`rpc_count_iostats` (C function), 64  
`rpc_count_iostats_metrics` (C function), 64  
`rpc_create` (C function), 66  
`rpc_find_or_alloc_pipe_dir_object` (C function), 66  
`rpc_force_rebind` (C function), 69  
`rpc_free` (C function), 63  
`rpc_free_iostats` (C function), 64  
`rpc_init_pipe_dir_head` (C function), 65  
`rpc_init_pipe_dir_object` (C function), 65  
`rpc_localaddr` (C function), 68  
`rpc_malloc` (C function), 63  
`rpc_max_bc_payload` (C function), 69  
`rpc_max_payload` (C function), 69  
`rpc_mkpipe_dentry` (C function), 65  
`rpc_net_ns` (C function), 69  
`rpc_peeraddr` (C function), 68  
`rpc_peeraddr2str` (C function), 68  
`rpc_protocol` (C function), 69  
`rpc_queue_upcall` (C function), 64  
`rpc_remove_pipe_dir_object` (C function), 66  
`rpc_run_task` (C function), 67  
`rpc_switch_client_transport` (C function), 67  
`rpc_unlink` (C function), 65  
`rpc_wake_up` (C function), 63  
`rpc_wake_up_status` (C function), 63  
`rpcb_getport_async` (C function), 66  
`rps_may_expire_flow` (C function), 84

## S

`sk_alloc` (C function), 45  
`sk_attach_filter` (C function), 50  
`sk_buff` (C type), 8  
`sk_capable` (C function), 44  
`sk_clone_lock` (C function), 45  
`sk_eat_skb` (C function), 30  
`sk_filter_trim_cap` (C function), 49  
`sk_for_each_entry_offset_rcu` (C function), 28  
`sk_has_allocations` (C function), 29  
`sk_net_capable` (C function), 44  
`sk_ns_capable` (C function), 44  
`sk_page_frag` (C function), 30  
`sk_rmem_alloc_get` (C function), 29  
`sk_set_memalloc` (C function), 44  
`sk_state_load` (C function), 30  
`sk_state_store` (C function), 30  
`sk_stream_wait_connect` (C function), 49  
`sk_stream_wait_memory` (C function), 49  
`sk_wait_data` (C function), 45  
`sk_wmem_alloc_get` (C function), 29

`skb_abort_seq_read` (C function), 40  
`skb_append` (C function), 39  
`skb_append_datato_frags` (C function), 41  
`skb_availroom` (C function), 16  
`skb_checksum_complete` (C function), 23  
`skb_checksum_none_assert` (C function), 24  
`skb_checksum_setup` (C function), 42  
`skb_checksum_trimmed` (C function), 42  
`skb_clone` (C function), 34  
`skb_clone_sk` (C function), 42  
`skb_clone_writable` (C function), 20  
`skb_cloned` (C function), 11  
`skb_complete_tx_timestamp` (C function), 23  
`skb_complete_wifi_ack` (C function), 23  
`skb_copy` (C function), 34  
`skb_copy_and_csum_datagram_msg` (C function), 48  
`skb_copy_bits` (C function), 37  
`skb_copy_datagram_from_iter` (C function), 48  
`skb_copy_datagram_iter` (C function), 47  
`skb_copy_expand` (C function), 35  
`skb_copy_ubufs` (C function), 34  
`skb_cow` (C function), 20  
`skb_cow_data` (C function), 42  
`skb_cow_head` (C function), 20  
`skb_dequeue` (C function), 15, 38  
`skb_dequeue_tail` (C function), 15, 38  
`skb_dst` (C function), 9  
`skb_dst_is_noref` (C function), 10  
`skb_dst_set` (C function), 9  
`skb_dst_set_noref` (C function), 9  
`skb_fclone_busy` (C function), 10  
`skb_fill_page_desc` (C function), 15  
`skb_find_text` (C function), 40  
`skb_frag_address` (C function), 19  
`skb_frag_address_safe` (C function), 19  
`skb_frag_dma_map` (C function), 19  
`skb_frag_foreach_page` (C function), 8  
`skb_frag_page` (C function), 18  
`skb_frag_ref` (C function), 18  
`skb_frag_set_page` (C function), 19  
`skb_frag_unref` (C function), 18  
`skb_get` (C function), 11  
`skb_get_timestamp` (C function), 22  
`skb_gso_network_seglen` (C function), 24  
`skb_gso_transport_seglen` (C function), 43  
`skb_gso_validate_mtu` (C function), 43  
`skb_has_shared_frag` (C function), 21  
`skb_head_is_locked` (C function), 24  
`skb_header_cloned` (C function), 11  
`skb_header_release` (C function), 11  
`skb_headroom` (C function), 15  
`skb_insert` (C function), 39  
`skb_kill_datagram` (C function), 47  
`skb_linearize` (C function), 21  
`skb_linearize_cow` (C function), 21  
`skb_mac_gso_segment` (C function), 84  
`skb_morph` (C function), 34



skb\_needs\_linearize (C function), 22  
skb\_orphan (C function), 17  
skb\_orphan\_frags (C function), 17  
skb\_pad (C function), 10  
skb\_padto (C function), 20  
skb\_page\_frag\_refill (C function), 45  
skb\_partial\_csum\_set (C function), 42  
skb\_peek (C function), 12  
skb\_peek\_next (C function), 13  
skb\_peek\_tail (C function), 13  
skb\_postpull\_rcsum (C function), 21  
skb\_postpush\_rcsum (C function), 22  
skb\_prepare\_seq\_read (C function), 40  
skb\_propagate\_pfmemalloc (C function), 18  
skb\_pull (C function), 36  
skb\_pull\_rcsum (C function), 41  
skb\_push (C function), 36  
skb\_push\_rcsum (C function), 22  
skb\_put (C function), 36  
skb\_put\_padto (C function), 21  
skb\_queue\_empty (C function), 10  
skb\_queue\_head (C function), 14, 38  
skb\_queue\_is\_first (C function), 10  
skb\_queue\_is\_last (C function), 10  
skb\_queue\_len (C function), 13  
skb\_queue\_next (C function), 11  
skb\_queue\_prev (C function), 11  
skb\_queue\_purge (C function), 17, 38  
skb\_queue\_splice (C function), 13  
skb\_queue\_splice\_init (C function), 14  
skb\_queue\_splice\_tail (C function), 14  
skb\_queue\_splice\_tail\_init (C function), 14  
skb\_queue\_tail (C function), 14, 39  
skb\_reserve (C function), 16  
skb\_scrub\_packet (C function), 43  
skb\_segment (C function), 41  
skb\_seq\_read (C function), 40  
skb\_share\_check (C function), 12  
skb\_shared (C function), 12  
skb\_shared\_hwtstamps (C type), 8  
skb\_split (C function), 39  
skb\_store\_bits (C function), 37  
skb\_tailroom (C function), 16  
skb\_tailroom\_reserve (C function), 16  
skb\_to\_sgvec (C function), 41  
skb\_trim (C function), 36  
skb\_try\_coalesce (C function), 43  
skb\_tstamp\_tx (C function), 23  
skb\_tx\_error (C function), 33  
skb\_tx\_timestamp (C function), 23  
skb\_unlink (C function), 39  
skb\_unshare (C function), 12  
skb\_zerocopy (C function), 37  
skwq\_has\_sleeper (C function), 29  
sock (C type), 24  
sock\_alloc (C function), 31  
sock\_common (C type), 24  
sock\_poll\_wait (C function), 30

sock\_register (C function), 32  
sock\_release (C function), 31  
sock\_tx\_timestamp (C function), 30  
sock\_type (C type), 7  
sock\_unregister (C function), 32  
socket (C type), 7  
sockfd\_lookup (C function), 31  
spans\_boundary (C function), 135  
svc\_find\_xprt (C function), 59  
svc\_print\_addr (C function), 58  
svc\_reserve (C function), 59  
svc\_xprt\_names (C function), 59  
synchronize\_net (C function), 93

## U

u64\_to\_ether\_addr (C function), 100  
unlock\_sock\_fast (C function), 28  
unregister\_netdev (C function), 93  
unregister\_netdevice\_many (C function), 93  
unregister\_netdevice\_notifier (C function), 82  
unregister\_netdevice\_queue (C function), 93

## W

wimax\_dev (C type), 74  
wimax\_dev\_add (C function), 74  
wimax\_dev\_init (C function), 73  
wimax\_dev\_rm (C function), 74  
wimax\_msg (C function), 71  
wimax\_msg\_alloc (C function), 70  
wimax\_msg\_data (C function), 71  
wimax\_msg\_data\_len (C function), 70  
wimax\_msg\_len (C function), 71  
wimax\_msg\_send (C function), 71  
wimax\_report\_rfkil\_hw (C function), 72  
wimax\_report\_rfkil\_sw (C function), 72  
wimax\_reset (C function), 72  
wimax\_rfkil (C function), 73  
wimax\_st (C type), 76  
wimax\_state\_change (C function), 73  
wimax\_state\_get (C function), 73  
write\_zsctrl (C function), 132  
write\_zsdata (C function), 132  
write\_zsreg (C function), 131

## X

xdr\_buf\_subsegment (C function), 58  
xdr\_buf\_trim (C function), 58  
xdr\_commit\_encode (C function), 55  
xdr\_encode\_opaque (C function), 55  
xdr\_encode\_opaque\_fixed (C function), 54  
xdr\_enter\_page (C function), 57  
xdr\_init\_decode (C function), 57  
xdr\_init\_decode\_pages (C function), 57  
xdr\_init\_encode (C function), 55  
xdr\_inline\_decode (C function), 57  
xdr\_partial\_copy\_from\_skb (C function), 64  
xdr\_read\_pages (C function), 57  
xdr\_reserve\_space (C function), 56

xdr\_restrict\_buflen (C function), [56](#)  
xdr\_set\_scratch\_buffer (C function), [57](#)  
xdr\_skb\_read\_bits (C function), [63](#)  
xdr\_stream\_decode\_string\_dup (C function), [58](#)  
xdr\_stream\_pos (C function), [55](#)  
xdr\_terminate\_string (C function), [55](#)  
xdr\_truncate\_encode (C function), [56](#)  
xdr\_write\_pages (C function), [56](#)  
xpirt\_adjust\_cwnd (C function), [61](#)  
xpirt\_complete\_rqst (C function), [62](#)  
xpirt\_disconnect\_done (C function), [62](#)  
xpirt\_force\_disconnect (C function), [62](#)  
xpirt\_get (C function), [63](#)  
xpirt\_load\_transport (C function), [60](#)  
xpirt\_lookup\_rqst (C function), [62](#)  
xpirt\_pin\_rqst (C function), [62](#)  
xpirt\_put (C function), [63](#)  
xpirt\_register\_transport (C function), [59](#)  
xpirt\_release\_rqst\_cong (C function), [60](#)  
xpirt\_release\_xpirt (C function), [60](#)  
xpirt\_release\_xpirt\_cong (C function), [60](#)  
xpirt\_reserve\_xpirt (C function), [60](#)  
xpirt\_set\_retrans\_timeout\_def (C function), [61](#)  
xpirt\_set\_retrans\_timeout\_rtt (C function), [62](#)  
xpirt\_unpin\_rqst (C function), [62](#)  
xpirt\_unregister\_transport (C function), [60](#)  
xpirt\_wait\_for\_buffer\_space (C function), [61](#)  
xpirt\_wake\_pending\_tasks (C function), [61](#)  
xpirt\_write\_space (C function), [61](#)

z8530\_tx\_done (C function), [134](#)  
z8530\_write\_port (C function), [131](#)  
zerocopy\_sg\_from\_iter (C function), [48](#)

## Z

z8530\_channel\_load (C function), [130](#)  
z8530\_describe (C function), [129](#)  
z8530\_dma\_rx (C function), [133](#)  
z8530\_dma\_status (C function), [133](#)  
z8530\_dma\_tx (C function), [133](#)  
z8530\_flush\_fifo (C function), [132](#)  
z8530\_init (C function), [129](#)  
z8530\_interrupt (C function), [128](#)  
z8530\_null\_rx (C function), [130](#)  
z8530\_queue\_xmit (C function), [130](#)  
z8530\_read\_port (C function), [131](#)  
z8530\_rtsdtr (C function), [132](#)  
z8530\_rx (C function), [132](#)  
z8530\_rx\_clear (C function), [134](#)  
z8530\_rx\_done (C function), [135](#)  
z8530\_shutdown (C function), [130](#)  
z8530\_status (C function), [133](#)  
z8530\_status\_clear (C function), [134](#)  
z8530\_sync\_close (C function), [128](#)  
z8530\_sync\_dma\_close (C function), [129](#)  
z8530\_sync\_dma\_open (C function), [128](#)  
z8530\_sync\_open (C function), [128](#)  
z8530\_sync\_txdma\_close (C function), [129](#)  
z8530\_sync\_txdma\_open (C function), [129](#)  
z8530\_tx (C function), [133](#)  
z8530\_tx\_begin (C function), [134](#)  
z8530\_tx\_clear (C function), [134](#)