

**NAME**

**ifnet, ifaddr, ifqueue, if\_data** - kernel interfaces for manipulating network interfaces

**SYNOPSIS**

```
#include <sys/param.h>
#include <sys/time.h>
#include <sys/socket.h>
#include <net/if.h>
#include <net/if_var.h>
#include <net/if_types.h>
```

**Interface Manipulation Functions**

*struct ifnet \**

**if\_alloc**(*u\_char type*);

*struct ifnet \**

**if\_alloc\_dev**(*u\_char type, device\_t dev*);

*struct ifnet \**

**if\_alloc\_domain**(*u\_char type, int numa\_domain*);

*void*

**if\_attach**(*struct ifnet \*ifp*);

*void*

**if\_detach**(*struct ifnet \*ifp*);

*void*

**if\_free**(*struct ifnet \*ifp*);

*void*

**if\_free\_type**(*struct ifnet \*ifp, u\_char type*);

*void*

**if\_down**(*struct ifnet \*ifp*);

*int*

**ifioctl**(*struct socket \*so, u\_long cmd, caddr\_t data, struct thread \*td*);

*int*

**ifpromisc**(*struct ifnet \*ifp, int pswitch*);

*int*

**if\_allmulti**(*struct ifnet \*ifp, int amswitch*);

*struct ifnet \**

**ifunit**(*const char \*name*);

*struct ifnet \**

**ifunit\_ref**(*const char \*name*);

*void*

**if\_up**(*struct ifnet \*ifp*);

### Interface Address Functions

*struct ifaddr \**

**ifaddr\_byindex**(*u\_short idx*);

*struct ifaddr \**

**ifa\_ifwithaddr**(*struct sockaddr \*addr*);

*struct ifaddr \**

**ifa\_ifwithstaddr**(*struct sockaddr \*addr, int fib*);

*struct ifaddr \**

**ifa\_ifwithnet**(*struct sockaddr \*addr, int ignore\_ptp, int fib*);

*struct ifaddr \**

**ifaof\_ifpforaddr**(*struct sockaddr \*addr, struct ifnet \*ifp*);

*void*

**ifa\_ref**(*struct ifaddr \*ifa*);

*void*

**ifa\_free**(*struct ifaddr \*ifa*);

### Interface Multicast Address Functions

*int*

**if\_addmulti**(*struct ifnet \*ifp, struct sockaddr \*sa, struct ifmultiaddr \*\*ifmap*);

*int*

**if\_delmulti**(*struct ifnet \*ifp, struct sockaddr \*sa*);

*struct ifmultiaddr \**

**if\_findmulti**(*struct ifnet \*ifp, struct sockaddr \*sa*);

### Output queue macros

**IF\_DEQUEUE**(*struct ifqueue \*ifq, struct mbuf \*m*);

### struct ifnet Member Functions

*void*

**(\*if\_input)**(*struct ifnet \*ifp, struct mbuf \*m*);

*int*

**(\*if\_output)**(*struct ifnet \*ifp, struct mbuf \*m, const struct sockaddr \*dst, struct route \*ro*);

*void*

**(\*if\_start)**(*struct ifnet \*ifp*);

*int*

**(\*if\_transmit)**(*struct ifnet \*ifp, struct mbuf \*m*);

*void*

**(\*if\_qflush)**(*struct ifnet \*ifp*);

*int*

**(\*if\_ioctl)**(*struct ifnet \*ifp, u\_long cmd, caddr\_t data*);

*void*

**(\*if\_init)**(*void \*if\_softc*);

*int*

**(\*if\_resolvemulti)**(*struct ifnet \*ifp, struct sockaddr \*\*retsa, struct sockaddr \*addr*);

### struct ifaddr member function

*void*

**(\*ifa\_rtrequest)**(*int cmd, struct rentry \*rt, struct rt\_addrinfo \*info*);

### Global Variables

*extern struct ifnethead ifnet;*

```
extern int if_index;
extern int ifqmaxlen;
```

## DATA STRUCTURES

The kernel mechanisms for handling network interfaces reside primarily in the *ifnet*, *if\_data*, *ifaddr*, and *ifmultiaddr* structures in `<net/if.h>` and `<net/if_var.h>` and the functions named above and defined in `/sys/net/if.c`. Those interfaces which are intended to be used by user programs are defined in `<net/if.h>`; these include the interface flags, the *if\_data* structure, and the structures defining the appearance of interface-related messages on the route(4) routing socket and in sysctl(3). The header file `<net/if_var.h>` defines the kernel-internal interfaces, including the *ifnet*, *ifaddr*, and *ifmultiaddr* structures and the functions which manipulate them. (A few user programs will need `<net/if_var.h>` because it is the prerequisite of some other header file like `<netinet/if_ether.h>`. Most references to those two files in particular can be replaced by `<net/ethernet.h>`.)

The system keeps a linked list of interfaces using the TAILQ macros defined in queue(3); this list is headed by a *struct ifnethead* called *ifnet*. The elements of this list are of type *struct ifnet*, and most kernel routines which manipulate interface as such accept or return pointers to these structures. Each interface structure contains an *if\_data* structure used for statistics and information. Each interface also has a TAILQ of interface addresses, described by *ifaddr* structures. An AF\_LINK address (see link\_addr(3)) describing the link layer implemented by the interface (if any) is accessed by the **ifaddr\_byindex()** function or *if\_addr* structure. (Some trivial interfaces do not provide any link layer addresses; this structure, while still present, serves only to identify the interface name and index.)

Finally, those interfaces supporting reception of multicast datagrams have a TAILQ of multicast group memberships, described by *ifmultiaddr* structures. These memberships are reference-counted.

Interfaces are also associated with an output queue, defined as a *struct ifqueue*; this structure is used to hold packets while the interface is in the process of sending another.

### The ifnet Structure

The fields of *struct ifnet* are as follows:

<i>if_softc</i>	(void *) A pointer to the driver's private state block. (Initialized by driver.)
<i>if_l2com</i>	(void *) A pointer to the common data for the interface's layer 2 protocol. (Initialized by <b>if_alloc()</b> .)
<i>if_vnet</i>	( <i>struct vnet</i> *) A pointer to the virtual network stack instance. (Initialized by <b>if_attach()</b> .)

- if\_home\_vnet* (*struct vnet \**) A pointer to the parent virtual network stack, where this *struct ifnet* originates from. (Initialized by **if\_attach()**.)
- if\_link* (**TAILQ\_ENTRY(ifnet)**) queue(3) macro glue.
- if\_xname* (*char \**) The name of the interface, (e.g., "fxp0" or "lo0"). (Initialized by driver (usually via **if\_initname()**.)
- if\_dname* (*const char \**) The name of the driver. (Initialized by driver (usually via **if\_initname()**.)
- if\_dunit* (*int*) A unique number assigned to each interface managed by a particular driver. Drivers may choose to set this to IF\_DUNIT\_NONE if a unit number is not associated with the device. (Initialized by driver (usually via **if\_initname()**.)
- if\_refcount* (*u\_int*) The reference count. (Initialized by **if\_alloc()**.)
- if\_addrhead* (*struct ifaddrhead*) The head of the queue(3) TAILQ containing the list of addresses assigned to this interface.
- if\_pcount* (*int*) A count of promiscuous listeners on this interface, used to reference-count the IFF\_PROMISC flag.
- if\_carp* (*struct carp\_if \**) A pointer to the CARP interface structure, carp(4). (Initialized by the driver-specific **if\_ioctl()** routine.)
- if\_bpf* (*struct bpf\_if \**) Opaque per-interface data for the packet filter, bpf(4). (Initialized by **bpf\_attach()**.)
- if\_index* (*u\_short*) A unique number assigned to each interface in sequence as it is attached. This number can be used in a *struct sockaddr\_dl* to refer to a particular interface by index (see link\_addr(3)). (Initialized by **if\_alloc()**.)
- if\_vlantrunk* (*struct ifvlantrunk \**) A pointer to 802.1Q trunk structure, vlan(4). (Initialized by the driver-specific **if\_ioctl()** routine.)
- if\_flags* (*int*) Flags describing operational parameters of this interface (see below). (Manipulated by generic code.)
- if\_drv\_flags* (*int*) Flags describing operational status of this interface (see below).

(Manipulated by driver.)

*if\_capabilities* (*int*) Flags describing the capabilities the interface supports (see below).

*if\_capenable* (*int*) Flags describing the enabled capabilities of the interface (see below).

*if\_linkmib* (*void \**) A pointer to an interface-specific MIB structure exported by `ifmib(4)`.  
(Initialized by driver.)

*if\_linkmiblen* (*size\_t*) The size of said structure. (Initialized by driver.)

*if\_data* (*struct if\_data*) More statistics and information; see *The if\_data structure*, below.  
(Initialized by driver, manipulated by both driver and generic code.)

*if\_multiaddrs* (*struct ifmultihead*) The head of the queue(3) TAILQ containing the list of multicast addresses assigned to this interface.

*if\_amcount* (*int*) A number of multicast requests on this interface, used to reference-count the IFF\_ALLMULTI flag.

*if\_addr* (*struct ifaddr \**) A pointer to the link-level interface address. (Initialized by **if\_alloc()**.)

*if\_snd* (*struct ifaltq*) The output queue. (Manipulated by driver.)

*if\_broadcastaddr*  
(*const u\_int8\_t \**) A link-level broadcast bytestring for protocols with variable address length.

*if\_bridge* (*void \**) A pointer to the bridge interface structure, `if_bridge(4)`. (Initialized by the driver-specific **if\_ioctl()** routine.)

*if\_label* (*struct label \**) A pointer to the MAC Framework label structure, `mac(4)`.  
(Initialized by **if\_alloc()**.)

*if\_afdata* (*void \**) An address family dependent data region.

*if\_afdata\_initialized*  
(*int*) Used to track the current state of address family initialization.

- if\_afdata\_lock* (*struct rwlock*) An *rwlock(9)* lock used to protect *if\_afdata* internals.
- if\_linktask* (*struct task*) A *taskqueue(9)* task scheduled for link state change events of the interface.
- if\_addr\_lock* (*struct rwlock*) An *rwlock(9)* lock used to protect interface-related address lists.
- if\_clones* (**LIST\_ENTRY**(*ifnet*)) *queue(3)* macro glue for the list of clonable network interfaces.
- if\_groups* (**TAILQ\_HEAD**(, *ifg\_list*)) The head of the *queue(3)* TAILQ containing the list of groups per interface.
- if\_pf\_kif* (*void \**) A pointer to the structure used for interface abstraction by *pf(4)*.
- if\_lagg* (*void \**) A pointer to the *lagg(4)* interface structure.
- if\_alloctype* (*u\_char*) The type of the interface as it was at the time of its allocation. It is used to cache the type passed to **if\_alloc()**, but unlike *if\_type*, it would not be changed by drivers.
- if\_numa\_domain* (*uint8\_t*) The NUMA domain of the hardware device associated with the interface. This is filled in with a wildcard value unless the kernel is NUMA aware, the system is a NUMA system, and the *ifnet* is allocated using **if\_alloc\_dev()** or **if\_alloc\_domain()**.

References to *ifnet* structures are gained by calling the **if\_ref()** function and released by calling the **if\_rele()** function. They are used to allow kernel code walking global interface lists to release the *ifnet* lock yet keep the *ifnet* structure stable.

There are in addition a number of function pointers which the driver must initialize to complete its interface with the generic interface layer:

### **if\_input()**

Pass a packet to an appropriate upper layer as determined from the link-layer header of the packet. This routine is to be called from an interrupt handler or used to emulate reception of a packet on this interface. A single function implementing **if\_input()** can be shared among multiple drivers utilizing the same link-layer framing, e.g., Ethernet.

**if\_output()**

Output a packet on interface *ifp*, or queue it on the output queue if the interface is already active.

**if\_transmit()**

Transmit a packet on an interface or queue it if the interface is in use. This function will return ENOBUFS if the devices software and hardware queues are both full. This function must be installed after **if\_attach()** to override the default implementation. This function is exposed in order to allow drivers to manage their own queues and to reduce the latency caused by a frequently gratuitous enqueue / dequeue pair to ifq. The suggested internal software queuing mechanism is *buf\_ring*.

**if\_qflush()**

Free mbufs in internally managed queues when the interface is marked down. This function must be installed after **if\_attach()** to override the default implementation. This function is exposed in order to allow drivers to manage their own queues and to reduce the latency caused by a frequently gratuitous enqueue / dequeue pair to ifq. The suggested internal software queuing mechanism is *buf\_ring*.

**if\_start()**

Start queued output on an interface. This function is exposed in order to provide for some interface classes to share a **if\_output()** among all drivers. **if\_start()** may only be called when the IFF\_DRV\_OACTIVE flag is not set. (Thus, IFF\_DRV\_OACTIVE does not literally mean that output is active, but rather that the device's internal output queue is full.) Please note that this function will soon be deprecated.

**if\_ioctl()**

Process interface-related ioctl(2) requests (defined in *<sys/sockio.h>*). Preliminary processing is done by the generic routine **ifioctl()** to check for appropriate privileges, locate the interface being manipulated, and perform certain generic operations like twiddling flags and flushing queues. See the description of **ifioctl()** below for more information.

**if\_init()**

Initialize and bring up the hardware, e.g., reset the chip and enable the receiver unit. Should mark the interface running, but not active (IFF\_DRV\_RUNNING, ~IFF\_DRV\_OACTIVE).

**if\_resolvemulti()**

Check the requested multicast group membership, *addr*, for validity, and if necessary compute a link-layer group which corresponds to that address which is returned in *\*retsa*. Returns zero on success, or an error code on failure.



## Interface Flags

Interface flags are used for a number of different purposes. Some flags simply indicate information about the type of interface and its capabilities; others are dynamically manipulated to reflect the current state of the interface. Flags of the former kind are marked <S> in this table; the latter are marked <D>. Flags which begin with "IFF\_DRV\_" are stored in *if\_drv\_flags*; all other flags are stored in *if\_flags*.

The macro IFF\_CANTCHANGE defines the bits which cannot be set by a user program using the SIOCSIFFLAGS command to ioctl(2); these are indicated by an asterisk ('\*') in the following listing.

IFF_UP	<D> The interface has been configured up by the user-level code.
IFF_BROADCAST	<S*> The interface supports broadcast.
IFF_DEBUG	<D> Used to enable/disable driver debugging code.
IFF_LOOPBACK	<S> The interface is a loopback device.
IFF_POINTOPOINT	<S*> The interface is point-to-point; "broadcast" address is actually the address of the other end.
IFF_DRV_RUNNING	<D*> The interface has been configured and dynamic resources were successfully allocated. Probably only useful internal to the interface.
IFF_NOARP	<D> Disable network address resolution on this interface.
IFF_PROMISC	<D*> This interface is in promiscuous mode.
IFF_PPROMISC	<D> This interface is in the permanently promiscuous mode (implies IFF_PROMISC).
IFF_ALLMULTI	<D*> This interface is in all-multicasts mode (used by multicast routers).
IFF_DRV_OACTIVE	<D*> The interface's hardware output queue (if any) is full; output packets are to be queued.
IFF_SIMPLEX	<S*> The interface cannot hear its own transmissions.
IFF_LINK0	
IFF_LINK1	
IFF_LINK2	<D> Control flags for the link layer. (Currently abused to select among multiple physical layers on some devices.)
IFF_MULTICAST	<S*> This interface supports multicast.
IFF_CANTCONFIG	<S*> The interface is not configurable in a meaningful way. Primarily useful for IFT_USB interfaces registered at the interface list.
IFF_MONITOR	<D> This interface blocks transmission of packets and discards incoming packets after BPF processing. Used to monitor network traffic but not interact with the network in question.
IFF_STATICARP	<D> Used to enable/disable ARP requests on this interface.
IFF_DYING	<D*> Set when the <i>ifnet</i> structure of this interface is being released and still has <i>if_refcount</i> references.

IFF\_RENAMING <D> Set when this interface is being renamed.

### Interface Capabilities Flags

Interface capabilities are specialized features an interface may or may not support. These capabilities are very hardware-specific and allow, when enabled, to offload specific network processing to the interface or to offer a particular feature for use by other kernel parts.

It should be stressed that a capability can be completely uncontrolled (i.e., stay always enabled with no way to disable it) or allow limited control over itself (e.g., depend on another capability's state.) Such peculiarities are determined solely by the hardware and driver of a particular interface. Only the driver possesses the knowledge on whether and how the interface capabilities can be controlled. Consequently, capabilities flags in *if\_capenable* should never be modified directly by kernel code other than the interface driver. The command SIOCSIFCAP to **ioctl()** is the dedicated means to attempt altering *if\_capenable* on an interface. Userland code shall use **ioctl(2)**.

The following capabilities are currently supported by the system:

IFCAP_RXCSUM	This interface can do checksum validation on receiving data. Some interfaces do not have sufficient buffer storage to store frames above a certain MTU-size completely. The driver for the interface might disable hardware checksum validation if the MTU is set above the hardcoded limit.
IFCAP_TXCSUM	This interface can do checksum calculation on transmitting data.
IFCAP_HWCSUM	A shorthand for (IFCAP_RXCSUM   IFCAP_TXCSUM).
IFCAP_NETCONS	This interface can be a network console.
IFCAP_VLAN_MTU	The vlan(4) driver can operate over this interface in software tagging mode without having to decrease MTU on vlan(4) interfaces below 1500 bytes. This implies the ability of this interface to cope with frames somewhat longer than permitted by the Ethernet specification.
IFCAP_VLAN_HWTAGGING	This interface can do VLAN tagging on output and demultiplex frames by their VLAN tag on input.
IFCAP_JUMBO_MTU	This Ethernet interface can transmit and receive frames up to 9000 bytes long.

IFCAP_POLLING	This interface supports polling(4). See below for details.
IFCAP_VLAN_HWCSUM	This interface can do checksum calculation on both transmitting and receiving data on vlan(4) interfaces (implies IFCAP_HWCSUM).
IFCAP_TSO4	This Ethernet interface supports TCP4 Segmentation offloading.
IFCAP_TSO6	This Ethernet interface supports TCP6 Segmentation offloading.
IFCAP_TSO	A shorthand for (IFCAP_TSO4   IFCAP_TSO6).
IFCAP_TOE4	This Ethernet interface supports TCP4 Offload Engine.
IFCAP_TOE6	This Ethernet interface supports TCP6 Offload Engine.
IFCAP_TOE	A shorthand for (IFCAP_TOE4   IFCAP_TOE6).
IFCAP_WOL_UCAST	This Ethernet interface supports waking up on any Unicast packet.
IFCAP_WOL_MCAST	This Ethernet interface supports waking up on any Multicast packet.
IFCAP_WOL_MAGIC	This Ethernet interface supports waking up on any Magic packet such as those sent by wake(8).
IFCAP_WOL	A shorthand for (IFCAP_WOL_UCAST   IFCAP_WOL_MCAST   IFCAP_WOL_MAGIC).
IFCAP_VLAN_HWFILTER	This interface supports frame filtering in hardware on vlan(4) interfaces.
IFCAP_VLAN_HWTSO	This interface supports TCP Segmentation offloading on vlan(4) interfaces (implies IFCAP_TSO).
IFCAP_LINKSTATE	This Ethernet interface supports dynamic link state changes.
IFCAP_NETMAP	This Ethernet interface supports netmap(4).

The ability of advanced network interfaces to offload certain computational tasks from the host CPU to the board is limited mostly to TCP/IP. Therefore a separate field associated with an interface (see *ifnet.if\_data.if\_hwassist* below) keeps a detailed description of its enabled capabilities specific to TCP/IP processing. The TCP/IP module consults the field to see which tasks can be done on an *outgoing* packet by the interface. The flags defined for that field are a superset of those for *mbuf.m\_pkthdr.csum\_flags*, namely:

CSUM_IP	The interface will compute IP checksums.
CSUM_TCP	The interface will compute TCP checksums.
CSUM_UDP	The interface will compute UDP checksums.

An interface notifies the TCP/IP module about the tasks the former has performed on an *incoming* packet by setting the corresponding flags in the field *mbuf.m\_pkthdr.csum\_flags* of the *mbuf chain* containing the packet. See *mbuf(9)* for details.

The capability of a network interface to operate in polling(4) mode involves several flags in different global variables and per-interface fields. The capability flag `IFCAP_POLLING` set in interface's *if\_capabilities* indicates support for polling(4) on the particular interface. If set in *if\_capabilities*, the same flag can be marked or cleared in the interface's *if\_capenable* within `ifioctl()`, thus initiating switch of the interface to polling(4) mode or interrupt mode, respectively. The actual mode change is managed by the driver-specific `if_ioctl()` routine. The polling(4) handler returns the number of packets processed.

### The *if\_data* Structure

The *if\_data* structure contains statistics and identifying information used by management programs, and which is exported to user programs by way of the `ifmib(4)` branch of the `sysctl(3)` MIB. The following elements of the *if\_data* structure are initialized by the interface and are not expected to change significantly over the course of normal operation:

<i>ifi_type</i>	( <i>u_char</i> ) The type of the interface, as defined in <code>&lt;net/if_types.h&gt;</code> and described below in the <i>Interface Types</i> section.
<i>ifi_physical</i>	( <i>u_char</i> ) Intended to represent a selection of physical layers on devices which support more than one; never implemented.
<i>ifi_addrlen</i>	( <i>u_char</i> ) Length of a link-layer address on this device, or zero if there are none. Used to initialize the address length field in <i>sockaddr_dl</i> structures referring to this interface.

- ifi\_hdrlen* (*u\_char*) Maximum length of any link-layer header which might be prepended by the driver to a packet before transmission. The generic code computes the maximum over all interfaces and uses that value to influence the placement of data in *mbufs* to attempt to ensure that there is always sufficient space to prepend a link-layer header without allocating an additional *mbuf*.
- ifi\_datalen* (*u\_char*) Length of the *if\_data* structure. Allows some stabilization of the routing socket ABI in the face of increases in the length of *struct ifdata*.
- ifi\_mtu* (*u\_long*) The maximum transmission unit of the medium, exclusive of any link-layer overhead.
- ifi\_metric* (*u\_long*) A dimensionless metric interpreted by a user-mode routing process.
- ifi\_baudrate* (*u\_long*) The line rate of the interface, in bits per second.
- ifi\_hwassist* (*u\_long*) A detailed interpretation of the capabilities to offload computational tasks for *outgoing* packets. The interface driver must keep this field in accord with the current value of *if\_capenable*.
- ifi\_epoch* (*time\_t*) The system uptime when interface was attached or the statistics below were reset. This is intended to be used to set the SNMP variable *ifCounterDiscontinuityTime*. It may also be used to determine if two successive queries for an interface of the same index have returned results for the same interface.

The structure additionally contains generic statistics applicable to a variety of different interface types (except as noted, all members are of type *u\_long*):

- ifi\_link\_state* (*u\_char*) The current link state of Ethernet interfaces. See the *Interface Link States* section for possible values.
- ifi\_ipackets* Number of packets received.
- ifi\_ierrors* Number of receive errors detected (e.g., FCS errors, DMA overruns, etc.). More detailed breakdowns can often be had by way of a link-specific MIB.
- ifi\_opackets* Number of packets transmitted.
- ifi\_oerrors* Number of output errors detected (e.g., late collisions, DMA overruns, etc.). More

detailed breakdowns can often be had by way of a link-specific MIB.

<i>ifi_collisions</i>	Total number of collisions detected on output for CSMA interfaces. (This member is sometimes [ab]used by other types of interfaces for other output error counts.)
<i>ifi_ibytes</i>	Total traffic received, in bytes.
<i>ifi_obytes</i>	Total traffic transmitted, in bytes.
<i>ifi_imcasts</i>	Number of packets received which were sent by link-layer multicast.
<i>ifi_omcasts</i>	Number of packets sent by link-layer multicast.
<i>ifi_iqdrops</i>	Number of packets dropped on input. Rarely implemented.
<i>ifi_oqdrops</i>	Number of packets dropped on output.
<i>ifi_noproto</i>	Number of packets received for unknown network-layer protocol.
<i>ifi_lastchange</i>	( <i>struct timeval</i> ) The time of the last administrative change to the interface (as required for SNMP).

### Interface Types

The header file `<net/if_types.h>` defines symbolic constants for a number of different types of interfaces. The most common are:

IFT_OTHER	none of the following
IFT_ETHER	Ethernet
IFT_ISO88023	ISO 8802-3 CSMA/CD
IFT_ISO88024	ISO 8802-4 Token Bus
IFT_ISO88025	ISO 8802-5 Token Ring
IFT_ISO88026	ISO 8802-6 DQDB MAN
IFT_FDDI	FDDI
IFT_PPP	Internet Point-to-Point Protocol (ppp(8))
IFT_LOOP	The loopback (lo(4)) interface
IFT_SLIP	Serial Line IP
IFT_PARA	Parallel-port IP ("PLIP")
IFT_ATM	Asynchronous Transfer Mode
IFT_USB	USB Interface

## Interface Link States

The following link states are currently defined:

LINK_STATE_UNKNOWN	The link is in an invalid or unknown state.
LINK_STATE_DOWN	The link is down.
LINK_STATE_UP	The link is up.

## The ifaddr Structure

Every interface is associated with a list (or, rather, a TAILQ) of addresses, rooted at the interface structure's *if\_addrhead* member. The first element in this list is always an AF\_LINK address representing the interface itself; multi-access network drivers should complete this structure by filling in their link-layer addresses after calling **if\_attach()**. Other members of the structure represent network-layer addresses which have been configured by means of the SIOCAIFADDR command to **ioctl(2)**, called on a socket of the appropriate protocol family. The elements of this list consist of *ifaddr* structures. Most protocols will declare their own protocol-specific interface address structures, but all begin with a *struct ifaddr* which provides the most-commonly-needed functionality across all protocols. Interface addresses are reference-counted.

The members of *struct ifaddr* are as follows:

<i>ifa_addr</i>	( <i>struct sockaddr</i> *) The local address of the interface.
<i>ifa_dstaddr</i>	( <i>struct sockaddr</i> *) The remote address of point-to-point interfaces, and the broadcast address of broadcast interfaces. ( <i>ifa_broadaddr</i> is a macro for <i>ifa_dstaddr</i> .)
<i>ifa_netmask</i>	( <i>struct sockaddr</i> *) The network mask for multi-access interfaces, and the confusion generator for point-to-point interfaces.
<i>ifa_ifp</i>	( <i>struct ifnet</i> *) A link back to the interface structure.
<i>ifa_link</i>	( <b>TAILQ_ENTRY</b> ( <i>ifaddr</i> )) queue(3) glue for list of addresses on each interface.
<i>ifa_rtrequest</i>	See below.
<i>ifa_flags</i>	( <i>u_short</i> ) Some of the flags which would be used for a route representing this address in the route table.
<i>ifa_refcnt</i>	( <i>short</i> ) The reference count.

References to *ifaddr* structures are gained by calling the **ifa\_ref()** function and released by calling the **ifa\_free()** function.

**ifa\_rtrequest()** is a pointer to a function which receives callouts from the routing code (**rtrequest()**) to perform link-layer-specific actions upon requests to add, or delete routes. The *cmd* argument indicates the request in question: RTM\_ADD, or RTM\_DELETE. The *rt* argument is the route in question; the *info* argument contains the specific destination being manipulated.

## FUNCTIONS

The functions provided by the generic interface code can be divided into two groups: those which manipulate interfaces, and those which manipulate interface addresses. In addition to these functions, there may also be link-layer support routines which are used by a number of drivers implementing a specific link layer over different hardware; see the documentation for that link layer for more details.

### The ifmultiaddr Structure

Every multicast-capable interface is associated with a list of multicast group memberships, which indicate at a low level which link-layer multicast addresses (if any) should be accepted, and at a high level, in which network-layer multicast groups a user process has expressed interest.

The elements of the structure are as follows:

- ifma\_link*      (**LIST\_ENTRY**(*ifmultiaddr*)) queue(3) macro glue.
- ifma\_addr*      (*struct sockaddr* \*) A pointer to the address which this record represents. The memberships for various address families are stored in arbitrary order.
- ifma\_lladdr*    (*struct sockaddr* \*) A pointer to the link-layer multicast address, if any, to which the network-layer multicast address in *ifma\_addr* is mapped, else a null pointer. If this element is non-nil, this membership also holds an invisible reference to another membership for that link-layer address.
- ifma\_refcount* (*u\_int*) A reference count of requests for this particular membership.

### Interface Manipulation Functions

#### **if\_alloc()**

Allocate and initialize *struct ifnet*. Initialization includes the allocation of an interface index and may include the allocation of a *type* specific structure in *if\_l2com*.

#### **if\_alloc\_dev()**

Allocate and initialize *struct ifnet* as **if\_alloc()** does, with the addition that the *ifnet* can be tagged



with the appropriate NUMA domain derived from the *dev* argument passed by the caller.

#### **if\_alloc\_domain()**

Allocate and initialize *struct ifnet* as **if\_alloc()** does, with the addition that the *ifnet* will be tagged with the NUMA domain via the *numa\_domain* argument passed by the caller.

#### **if\_attach()**

Link the specified interface *ifp* into the list of network interfaces. Also initialize the list of addresses on that interface, and create a link-layer *ifaddr* structure to be the first element in that list. (A pointer to this address structure is saved in the *ifnet* structure and is accessed by the **ifaddr\_byindex()** function.) The *ifp* must have been allocated by **if\_alloc()**, **if\_alloc\_dev()** or **if\_alloc\_domain()**.

#### **if\_detach()**

Shut down and unlink the specified *ifp* from the interface list.

#### **if\_free()**

Free the given *ifp* back to the system. The interface must have been previously detached if it was ever attached.

#### **if\_free\_type()**

Identical to **if\_free()** except that the given *type* is used to free *if\_l2com* instead of the type in *if\_type*. This is intended for use with drivers that change their interface type.

#### **if\_down()**

Mark the interface *ifp* as down (i.e., **IFF\_UP** is not set), flush its output queue, notify protocols of the transition, and generate a message from the route(4) routing socket.

#### **if\_up()**

Mark the interface *ifp* as up, notify protocols of the transition, and generate a message from the route(4) routing socket.

#### **ifpromisc()**

Add or remove a promiscuous reference to *ifp*. If *pswitch* is true, add a reference; if it is false, remove a reference. On reference count transitions from zero to one and one to zero, set the **IFF\_PROMISC** flag appropriately and call **if\_ioctl()** to set up the interface in the desired mode.

#### **if\_allmulti()**

As **ifpromisc()**, but for the all-multicasts (**IFF\_ALLMULTI**) flag instead of the promiscuous flag.

**ifunit()**

Return an *ifnet* pointer for the interface named *name*.

**ifunit\_ref()**

Return a reference-counted (via **ifa\_ref()**) *ifnet* pointer for the interface named *name*. This is the preferred function over **ifunit()**. The caller is responsible for releasing the reference with **if\_rele()** when it is finished with the *ifnet*.

**ifioctl()**

Process the ioctl request *cmd*, issued on socket *so* by thread *td*, with data parameter *data*. This is the main routine for handling all interface configuration requests from user mode. It is ordinarily only called from the socket-layer ioctl(2) handler, and only for commands with class 'i'. Any unrecognized commands will be passed down to socket *so*'s protocol for further interpretation. The following commands are handled by **ifioctl()**:

SIOCGIFCONF	Get interface configuration. (No call-down to driver.)
SIOCSIFNAME	Set the interface name. RTM_IFANNOUNCE departure and arrival messages are sent so that routing code that relies on the interface name will update its interface list. Caller must have appropriate privilege. (No call-down to driver.)
SIOCGIFCAP	
SIOCGIFDATA	
SIOCGIFFIB	
SIOCGIFFLAGS	
SIOCGIFMETRIC	
SIOCGIFMTU	
SIOCGIFPHYS	Get interface capabilities, data, FIB, flags, metric, MTU, medium selection. (No call-down to driver.)
SIOCSIFCAP	Enable or disable interface capabilities. Caller must have appropriate privilege. Before a call to the driver-specific <b>if_ioctl()</b> routine, the requested mask for enabled capabilities is checked against the mask of capabilities supported by the interface, <i>if_capabilities</i> . Requesting to enable an unsupported capability is invalid. The rest is supposed to be done by the driver, which includes updating <i>if_capenable</i> and <i>if_data.if_hwassist</i> appropriately.
SIOCGIFCAPNV	NV(9) version of the SIOCGIFCAP ioctl. Caller must provide a

pointer to *struct ifreq\_cap\_nv* as *data*, where the member *buffer* points to some buffer containing *buf\_length* bytes. The serialized nvlist with description of the device capabilities is written to the buffer. If *buffer* is too short, the structure is updated with *buffer* member set to `NULL`, *length* set to the minimal required length, and error `EFBIG` is returned.

Elements of the returned nvlist for simple capabilities are boolean, identified by names. Presence of the boolean element means that corresponding capability is supported by the interface. Element's value describes the current configured state: `true` means that the capability is enabled, and `false` that it is disabled.

Driver indicates support for both `SIOCGIFCAPNV` and `SIOCSIFCAPNV` requests by setting `IFCAP_NV` non-modifiable capability bit in *if\_capabilities*.

<code>SIOCSIFCAPNV</code>	NV(9) version of the <code>SIOCSIFCAP</code> ioctl. Caller must provide the pointer to <i>struct ifreq_cap_nv</i> as <i>data</i> , where the member <i>buffer</i> points to serialized nvlist of length bytes. Each element of nvlist describes a requested update of one capability, identified by the element name. For simple capabilities, the element must be boolean. Its <code>true</code> value means that the caller asks to enable the capability, and <code>false</code> value to disable. Only capabilities listed in the nvlist are affected by the call.
<code>SIOCSIFFIB</code>	Sets interface FIB. Caller must have appropriate privilege. FIB values start at 0 and values greater or equals than <i>net.fibs</i> are considered invalid.
<code>SIOCSIFFLAGS</code>	Change interface flags. Caller must have appropriate privilege. If a change to the <code>IFF_UP</code> flag is requested, <code>if_up()</code> or <code>if_down()</code> is called as appropriate. Flags listed in <code>IFF_CANTCHANGE</code> are masked off, and the field <i>if_flags</i> in the interface structure is updated. Finally, the driver <code>if_ioctl()</code> routine is called to perform any setup requested.
<code>SIOCSIFMETRIC</code> <code>SIOCSIFPHYS</code>	Change interface metric or medium. Caller must have appropriate privilege.
<code>SIOCSIFMTU</code>	Change interface MTU. Caller must have appropriate privilege.

MTU values less than 72 or greater than 65535 are considered invalid. The driver **if\_ioctl()** routine is called to implement the change; it is responsible for any additional sanity checking and for actually modifying the MTU in the interface structure.

**SIOCADDMULTI**

**SIOCDELMULTI**

Add or delete permanent multicast group memberships on the interface. Caller must have appropriate privilege. The **if\_addmulti()** or **if\_delmulti()** function is called to perform the operation; qq.v.

**SIOCAIFADDR**

**SIOCDEFADDR**

The socket's protocol control routine is called to implement the requested action.

### Interface Address Functions

Several functions exist to look up an interface address structure given an address. **ifa\_ifwithaddr()** returns an interface address with either a local address or a broadcast address precisely matching the parameter *addr*. **ifa\_ifwithdstaddr()** returns an interface address for a point-to-point interface whose remote ("destination") address is *addr* and a fib is *fib*. If *fib* is RT\_ALL\_FIBS, then the first interface address matching *addr* will be returned.

**ifa\_ifwithnet()** returns the most specific interface address which matches the specified address, *addr*, subject to its configured netmask, or a point-to-point interface address whose remote address is *addr* if one is found. If *ignore\_ptp* is true, skip point-to-point interface addresses. The *fib* parameter is handled the same way as by **ifa\_ifwithdstaddr()**.

**ifaof\_ifpforaddr()** returns the most specific address configured on interface *ifp* which matches address *addr*, subject to its configured netmask. If the interface is point-to-point, only an interface address whose remote address is precisely *addr* will be returned.

**ifaddr\_byindex()** returns the link-level address of the interface with the given index *idx*.

All of these functions return a null pointer if no such address can be found.

### Interface Multicast Address Functions

The **if\_addmulti()**, **if\_delmulti()**, and **if\_findmulti()** functions provide support for requesting and relinquishing multicast group memberships, and for querying an interface's membership list, respectively. The **if\_addmulti()** function takes a pointer to an interface, *ifp*, and a generic address, *sa*. It also takes a pointer to a *struct ifmultiaddr \** which is filled in on successful return with the address of the group membership control block. The **if\_addmulti()** function performs the following four-step process:

1. Call the interface's **if\_resolvemulti()** entry point to determine the link-layer address, if any, corresponding to this membership request, and also to give the link layer an opportunity to veto this membership request should it so desire.
2. Check the interface's group membership list for a pre-existing membership for this group. If one is not found, allocate a new one; if one is, increment its reference count.
3. If the **if\_resolvemulti()** routine returned a link-layer address corresponding to the group, repeat the previous step for that address as well.
4. If the interface's multicast address filter needs to be changed because a new membership was added, call the interface's **if\_ioctl()** routine (with a *cmd* argument of SIOCADDMULTI) to request that it do so.

The **if\_delmulti()** function, given an interface *ifp* and an address, *sa*, reverses this process. Both functions return zero on success, or a standard error number on failure.

The **if\_findmulti()** function examines the membership list of interface *ifp* for an address matching *sa*, and returns a pointer to that *struct ifmultiaddr* if one is found, else it returns a null pointer.

## SEE ALSO

ioctl(2), link\_addr(3), queue(3), sysctl(3), bpf(4), ifmib(4), lo(4), netintro(4), polling(4), config(8), ppp(8), mbuf(9), rentry(9)

Gary R. Wright and W. Richard Stevens, *TCP/IP Illustrated*, Vol. 2, Addison-Wesley, ISBN 0-201-63354-X.

## AUTHORS

This manual page was written by Garrett A. Wollman.