

# **The TCP/IP stack in the FreeBSD kernel: an overview of the implementation**



**Kevin Lo**

**msi**

**The FreeBSD project**

# Examples of operating systems use FreeBSD-based network stack

DragonflyBSD, a fork of FreeBSD

OS X from Apple

Os<sup>v</sup> from Cloudeus Systems

RTEMS

# Peeking at the Linux kernel changelogs

- Increase the initial cwnd to 10 (RFC 6928)  
: 2.6.39
- Proportional Rate Reduction for TCP  
(RFC 6937) : 3.2
- Early Retransmit for TCP (RFC 5827)  
: 3.5
- TCP Fast Open : 3.6, 3.7

# Memory buffers (mbufs)

struct mbuf: the most important data structure in the FreeBSD networking subsystem, which is defined in `<sys/mbuf.h>`

Every packet sent/received is handled using the mbuf structure

Mbufs are fixed size data buffers (256 bytes each)

# mbuf structure

```
struct mbuf {
    struct m_hdr  m_hdr;
    union {
        struct {
            struct pkthdr  MH_pkthdr;      /* M_PKTHDR set */
            union {
                struct m_ext  MH_ext; /* M_EXT set */
                char          MH_databuf[MHLEN];
            } MH_dat;
        } MH;
        char M_databuf[MLEN];      /* !M_PKTHDR, !M_EXT */
    } M_dat;
};
```

# m\_hdr structure

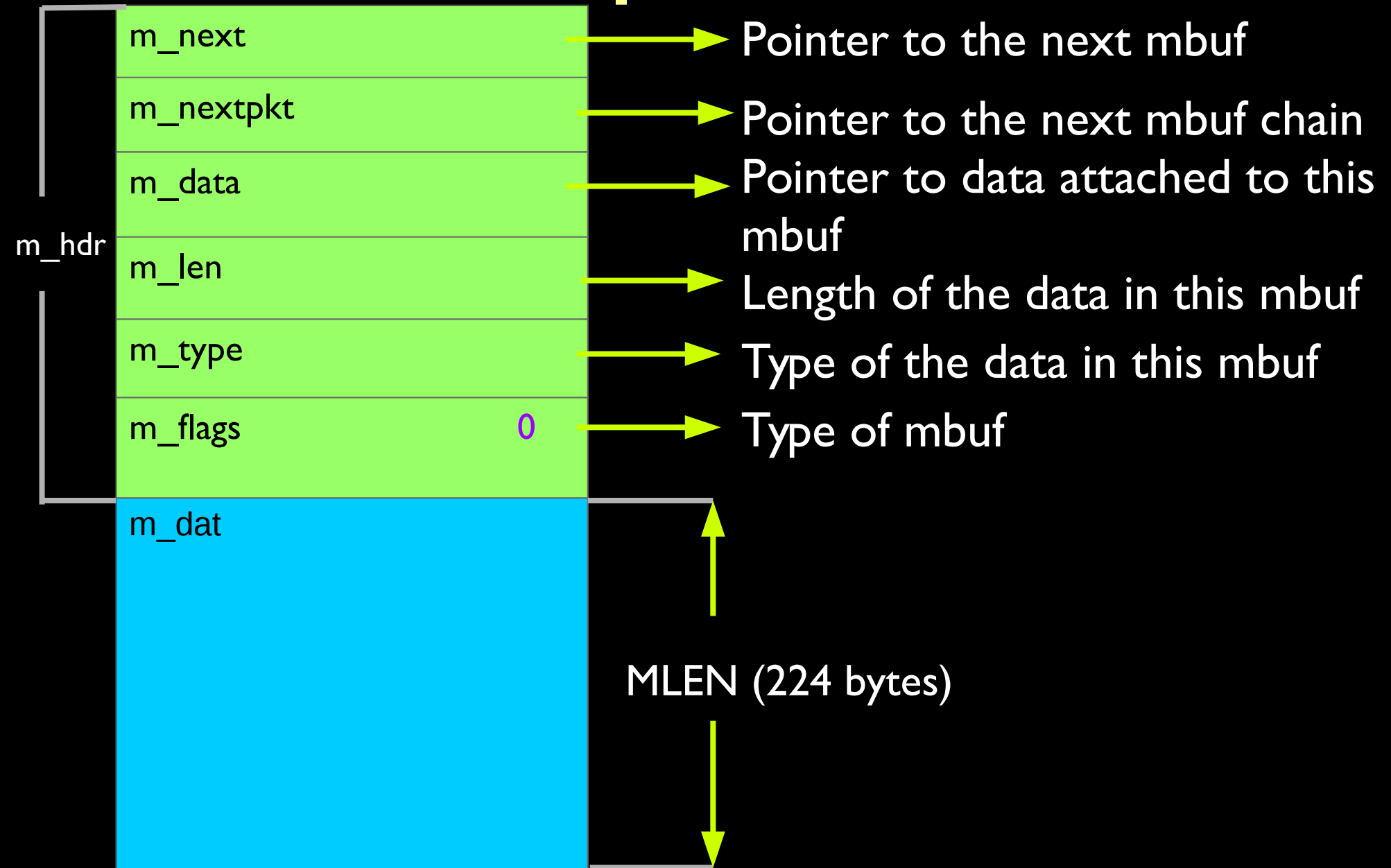
```
/*
 * Header present at the beginning of every mbuf.
 * Size ILP32: 24
 *      LP64: 32
 */
struct m_hdr {
    struct mbuf    *mh_next;      /* next buffer in chain */
    struct mbuf    *mh_nextpkt;  /* next chain in queue/record */
    caddr_t        mh_data;      /* location of data */
    int32_t        mh_len;       /* amount of data in this mbuf */
    uint32_t        mh_type:8,   /* type of data in this mbuf */
                  mh_flags:24;  /* flags; see below */
#ifdef __LP64__
    uint32_t        mh_pad;      /* pad for 64bit alignment */
#endif
};
```

On **64-bit** platforms, this results into  $3 * 8 \text{ bytes} + 2 * 4 \text{ bytes} = 32 \text{ bytes}$

# Simple mbuf

```
struct mbuf {
    struct m_hdr  m_hdr;
    union {
        struct {
            struct pkthdr  MH_pkthdr;      /* M_PKTHDR set */
            union {
                struct m_ext  MH_ext; /* M_EXT set */
                char          MH_databuf[MHLEN];
            } MH_dat;
        } MH;
        char M_databuf[MLEN];      /* !M_PKTHDR, !M_EXT */
    } M_dat;
};
```

# Simple mbuf

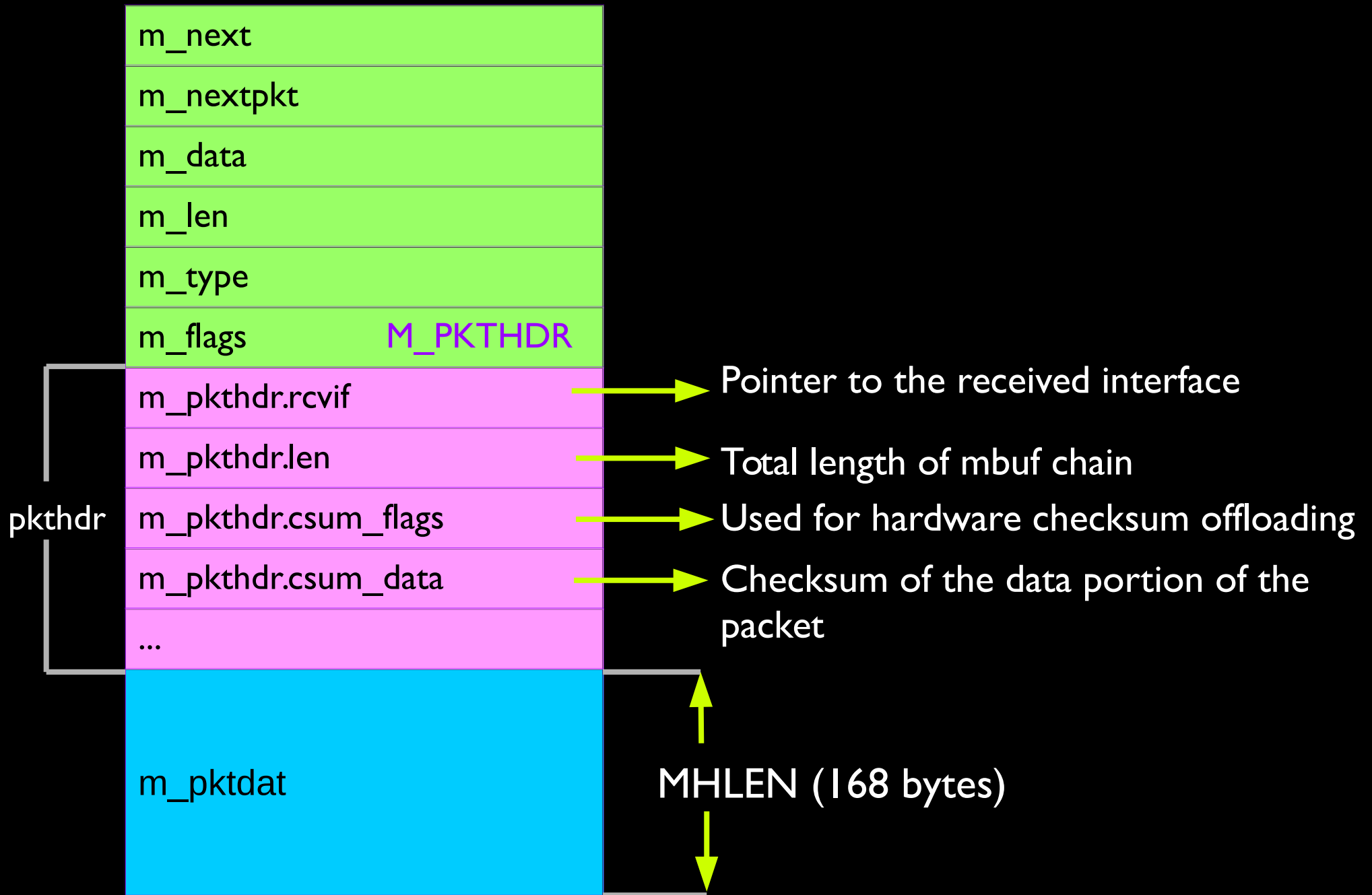




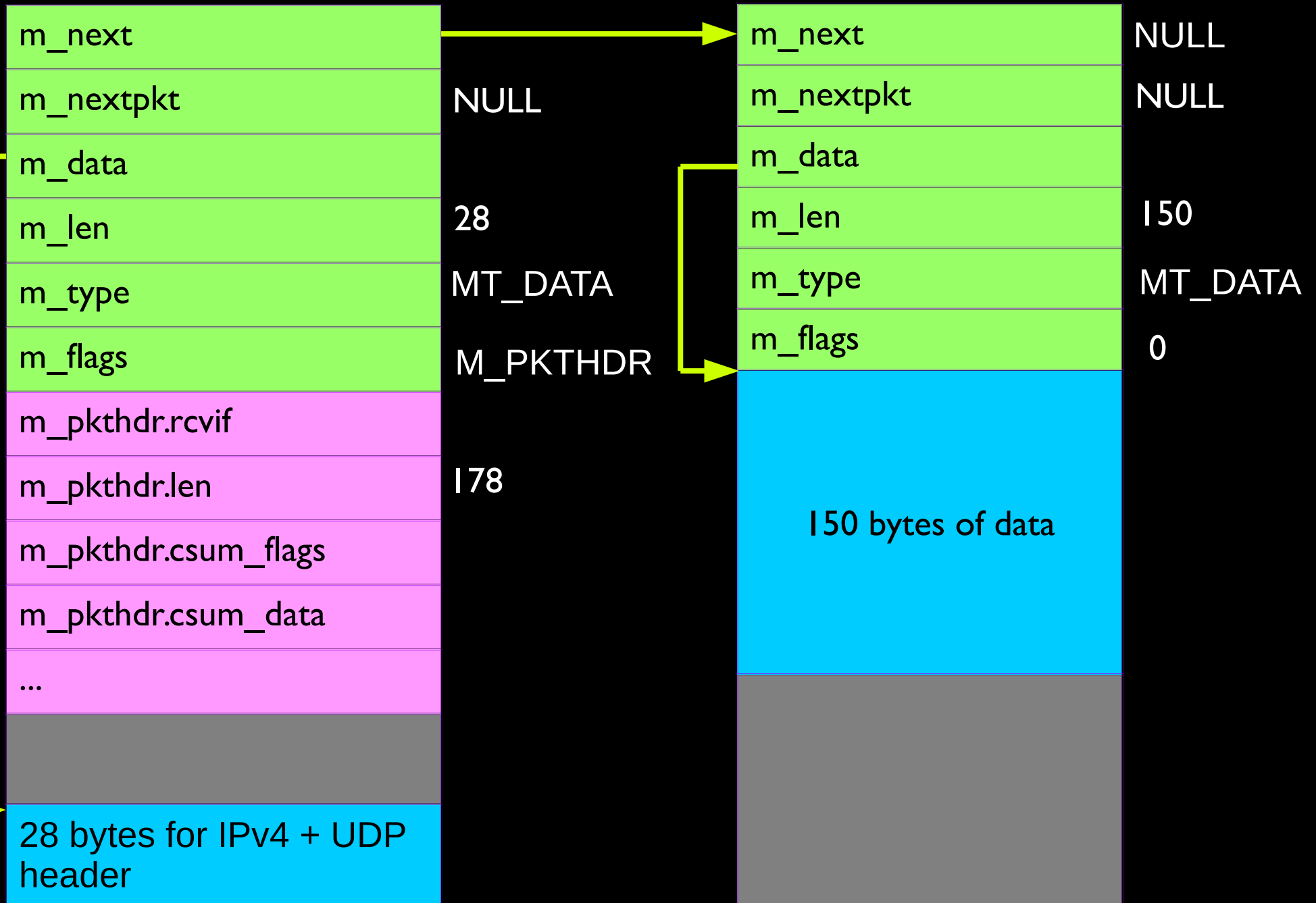
# Packet header mbuf

```
struct mbuf {
    struct m_hdr  m_hdr;
    union {
        struct {
            struct pkthdr  MH_pkthdr;      /* M_PKTHDR set */
            union {
                struct m_ext  MH_ext; /* M_EXT set */
                char          MH_databuf[MHLEN];
            } MH_dat;
        } MH;
        char M_databuf[MLEN];      /* !M_PKTHDR, !M_EXT */
    } M_dat;
};
```

# Packet header mbuf



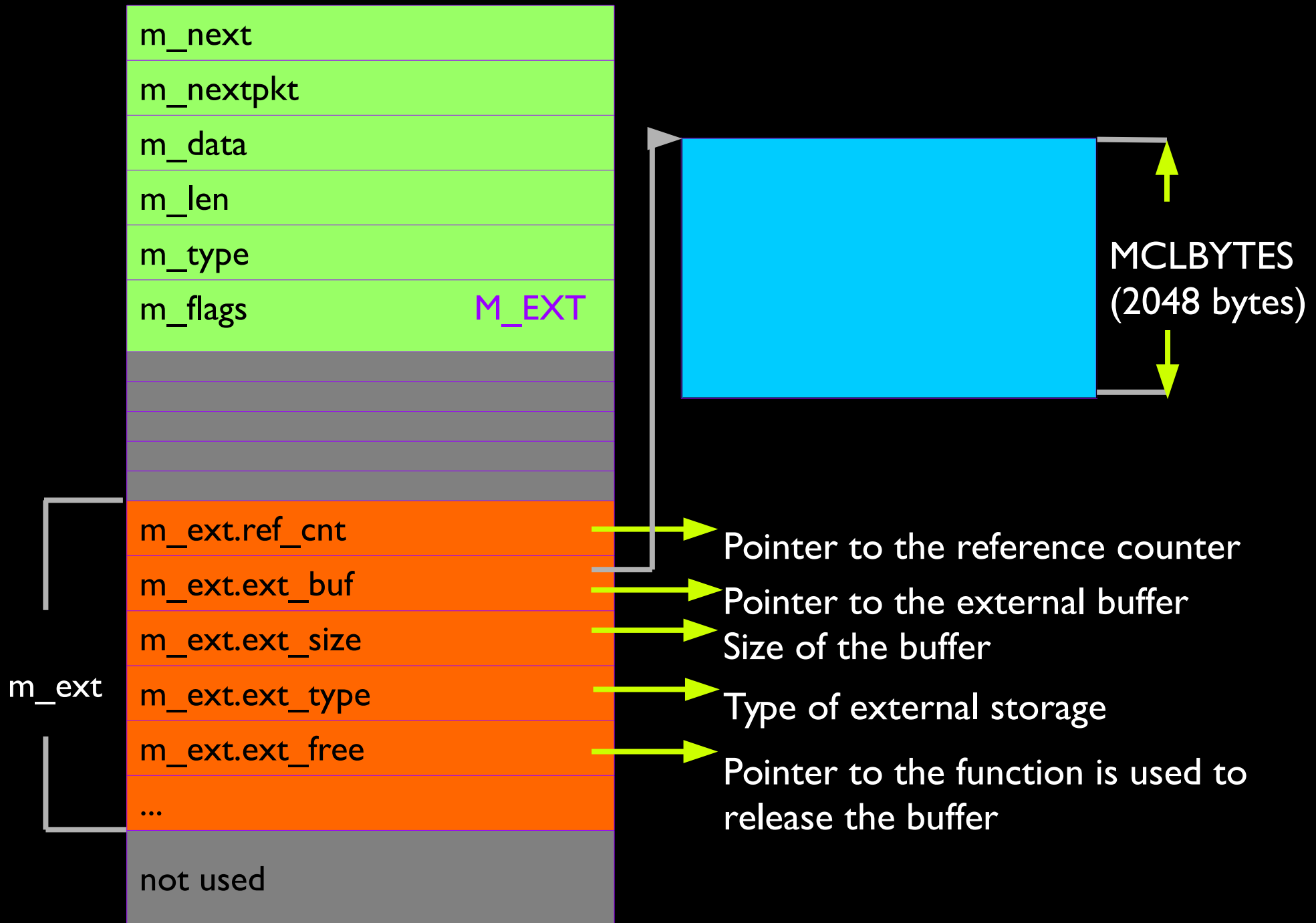
# A typical UDP packet



# Mbuf page cluster

```
struct mbuf {
    struct m_hdr  m_hdr;
    union {
        struct {
            struct pkthdr  MH_pkthdr;      /* M_PKTHDR set */
            union {
                struct m_ext  MH_ext; /* M_EXT set */
                char          MH_databuf[MHLEN];
            } MH_dat;
        } MH;
        char M_databuf[MLEN];      /* !M_PKTHDR, !M_EXT */
    } M_dat;
};
```

# Mbuf page cluster



# Packet header + page cluster

```
struct mbuf {
    struct m_hdr  m_hdr;
    union {
        struct {
            struct pkthdr  MH_pkthdr;      /* M_PKTHDR set */
            union {
                struct m_ext  MH_ext; /* M_EXT set */
                char          MH_databuf[MHLEN];
            } MH_dat;
        } MH;
        char M_databuf[MLEN];      /* !M_PKTHDR, !M_EXT */
    } M_dat;
};
```

# Packet header + page cluster



# Mbuf utility routines

MGET() / m\_get(): allocate an mbuf

MGETHDR() / m\_gethdr(): allocate an mbuf with a packet header

MCLGET() / m\_clget(): add an external cluster to an mbuf

m\_free(): free a single mbuf

m\_freem(): free a chain of mbufs

man mbuf



# Protocol data structures

The protocol layer uses three main types of structures:

- **domain** structure
- protocol switch structure (**protosw** & **ip6protosw**)
- protocol control block (**PCB**)

# Communication domains

Group of related protocols

Each has address family constant



# Supported address families

AF\_LOCAL / AF\_UNIX

Local communication

AF\_INET

Internet version 4

AF\_INET6

Internet version 6

AF\_ROUTE

Link layer interface

PF\_KEY

Internal key-management

AF\_NATM

Asynchronous transfer mode

AF\_NETGRAPH

Netgraph sockets

AF\_BLUETOOTH

Bluetooth protocols

AF\_INET\_SDP

OFED socket direct protocol

# RFC 2367 section 1.3

The PF\_KEY protocol family (PF\_KEY) symbol is defined in `<sys/socket.h>` in the same manner that other protocol families are defined. PF\_KEY does not use any socket addresses.

Applications using PF\_KEY **MUST NOT** depend on the availability of a symbol named AF\_KEY, but kernel implementations are encouraged to define that symbol for completeness.

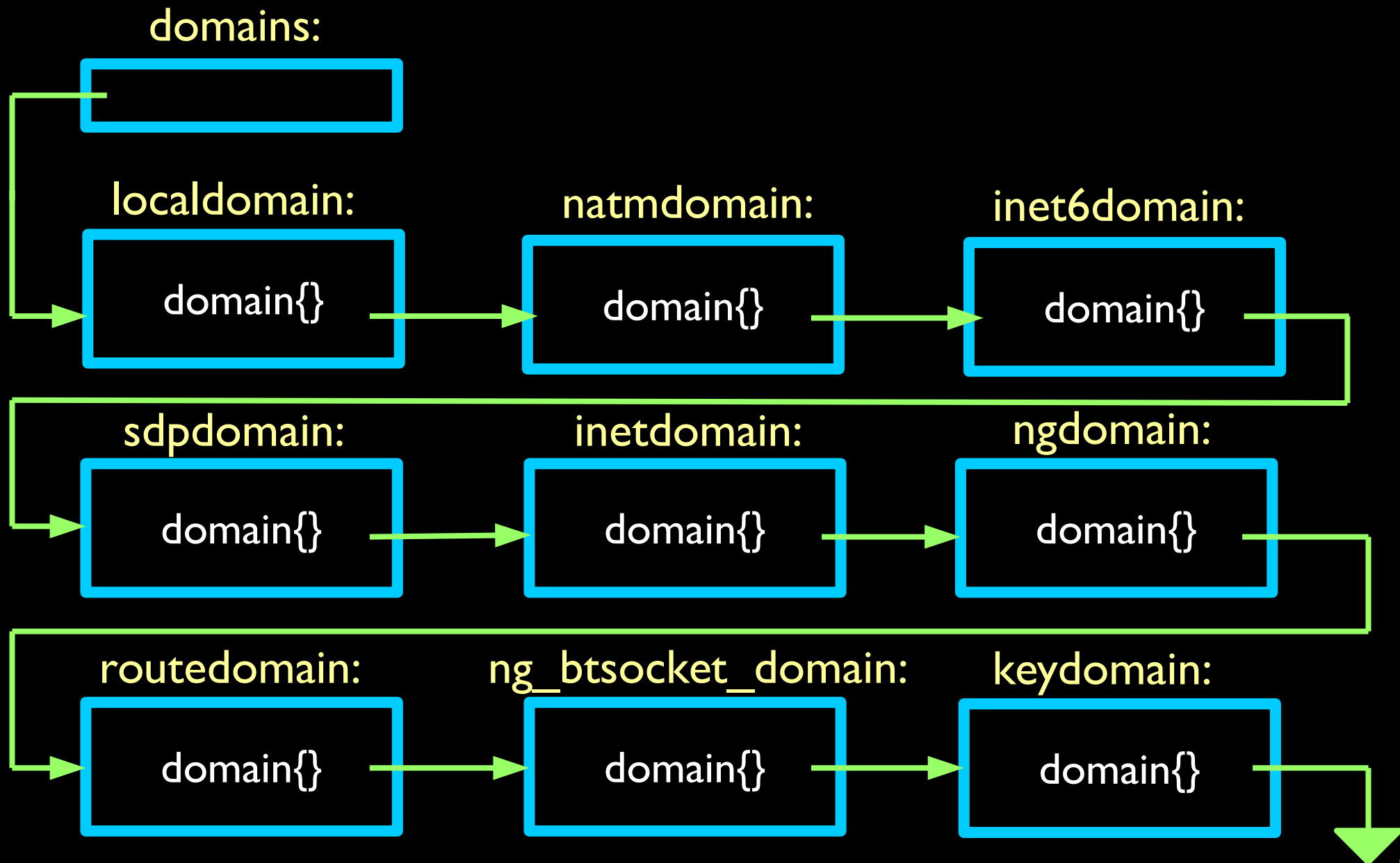
```
int s;
```

```
s = socket(AF_KEY, SOCK_RAW, PF_KEY_V2);
```

# inetdomain

```
struct domain inetdomain = {
    .dom_family =      AF_INET,
    .dom_name =       "internet",
    .dom_protosw =     inetsw,
    .dom_protoswNPROTOSW = &inetsw[sizeof(inetsw)/sizeof(inetsw[0])],
#ifdef RADIX_MPATH
    .dom_rtattach =    rn4_mpath_inithead,
#else
    .dom_rtattach =    in_inithead,
#endif
#ifdef VIMAGE
    .dom_rtdetach =    in_detachhead,
#endif
    .dom_rtoffset =    32,
    .dom_maxrtkey =    sizeof(struct sockaddr_in),
    .dom_ifattach =    in_domifattach,
    .dom_ifdetach =    in_domifdetach
};
VNET_DOMAIN_SET(inet);
```

# Domains list



# protosw structure

Defined in `<sys/protosw.h>`

```
/* USE THESE FOR YOUR PROTOTYPES ! */
```

```
typedef void pr_input_t (struct mbuf *, int);
```

```
typedef int pr_input6_t (struct mbuf **, int*, int); /* XXX FIX THIS */
```

```
typedef int pr_output_t (struct mbuf *, struct socket *);
```

```
typedef void pr_ctlinput_t (int, struct sockaddr *, void *);
```

```
typedef int pr_ctloutput_t (struct socket *, struct sockopt *);
```

```
typedef void pr_init_t (void);
```

```
typedef void pr_destroy_t (void);
```

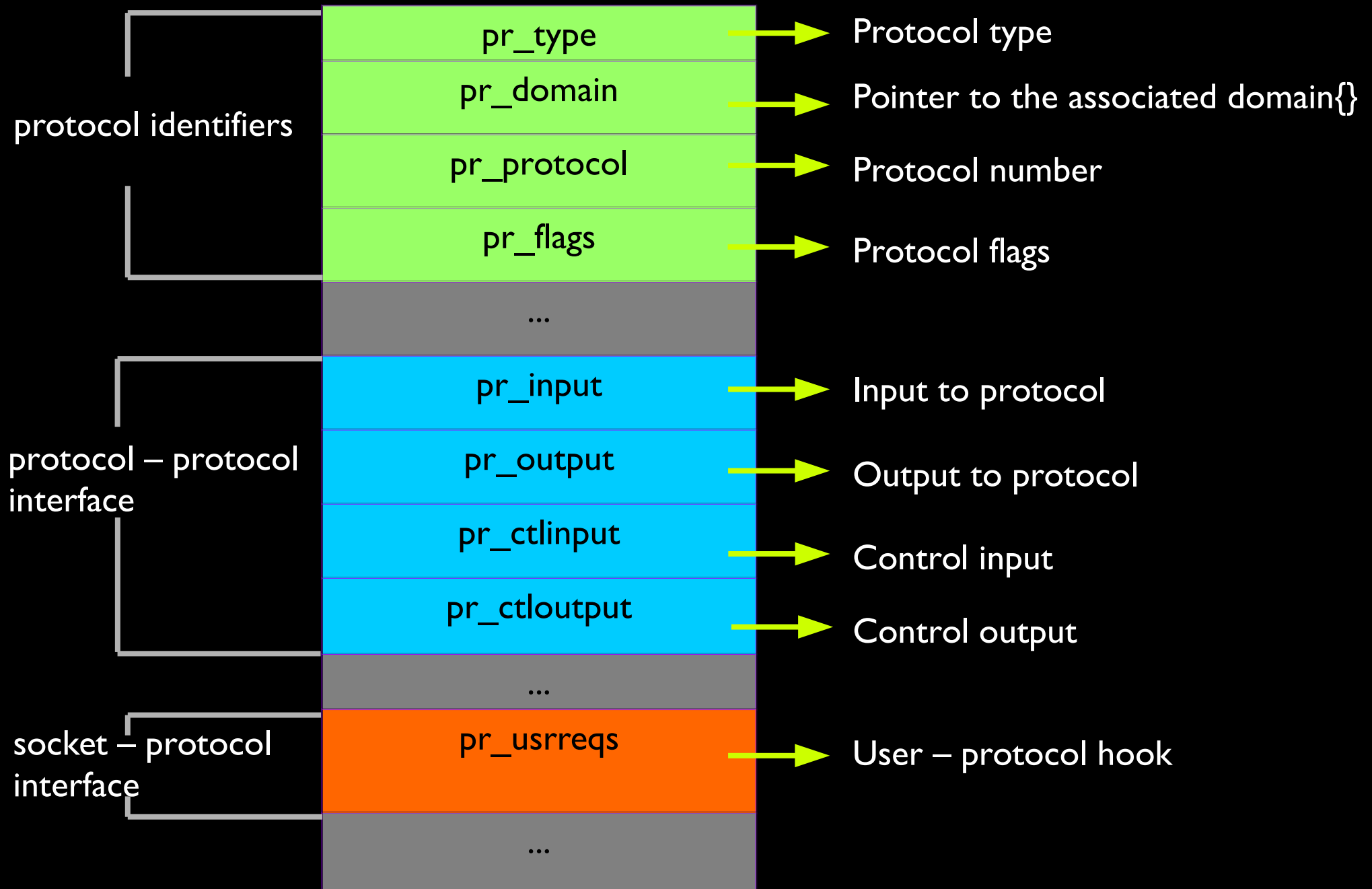
```
typedef void pr_fasttimo_t (void);
```

```
typedef void pr_slowtimo_t (void);
```

```
typedef void pr_drain_t (void);
```



# protosw structure (cont.)



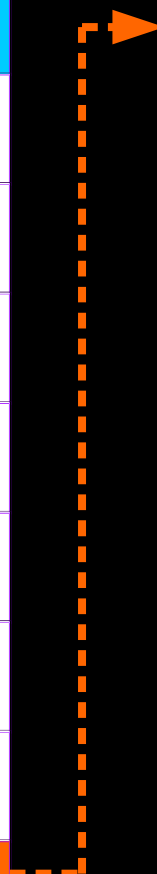
# inetsw[] switch table

```
struct protosw inetsw[] = {
{
    .pr_type =          SOCK_DGRAM,
    .pr_domain =       &inetdomain,
    .pr_protocol =     IPPROTO_UDP,
    .pr_flags =        PR_ATOMIC|PR_ADDR,
    .pr_input =        udp_input,
    .pr_ctlinput =     udp_ctlinput,
    .pr_ctloutput =    udp_ctloutput,
    .pr_init =         udp_init,
#ifdef VIMAGE
    .pr_destroy =      udp_destroy,
#endif
    .pr_usrreqs =      &udp_usrreqs
},
```

<b>socket{}</b>	
	so_count
	so_type
	so_options
	so_linger
	so_state
	so_qstate
protocol layer info	so_pcb
	<b>so_proto</b>
	...
socket buffers	so_rcv
	so_snd
	...

<b>protosw{}</b>	
	pr_type
	pr_domain
	pr_protocol
	pr_flags
	...
	pr_input
	pr_output
	pr_ctlinput
	pr_ctloutput
	...
	<b>pr_usrreqs</b>
	...

<b>pr_usrreqs{}</b>	
	pru_aboart
	pru_accept
	pru_attach
	pru_bind
	pru_connect
	...
	pru_detach
	pru_disconnect
	pru_listen
	pru_rcvd
	pru_send
	...



protocol  
layer info

socket  
buffers

# Protocol control block (pcb)

Hold protocol information

Stored as a doubly linked list

Internet protocol control block (inpcb)

- Foreign and local IP addresses

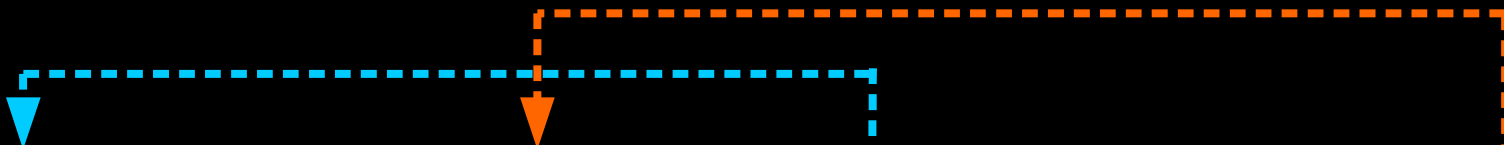
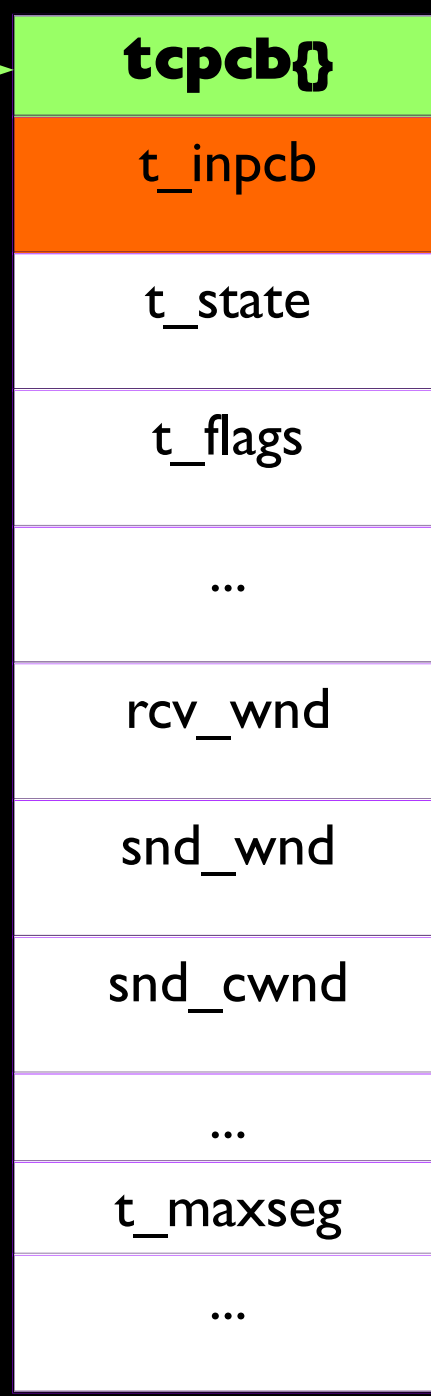
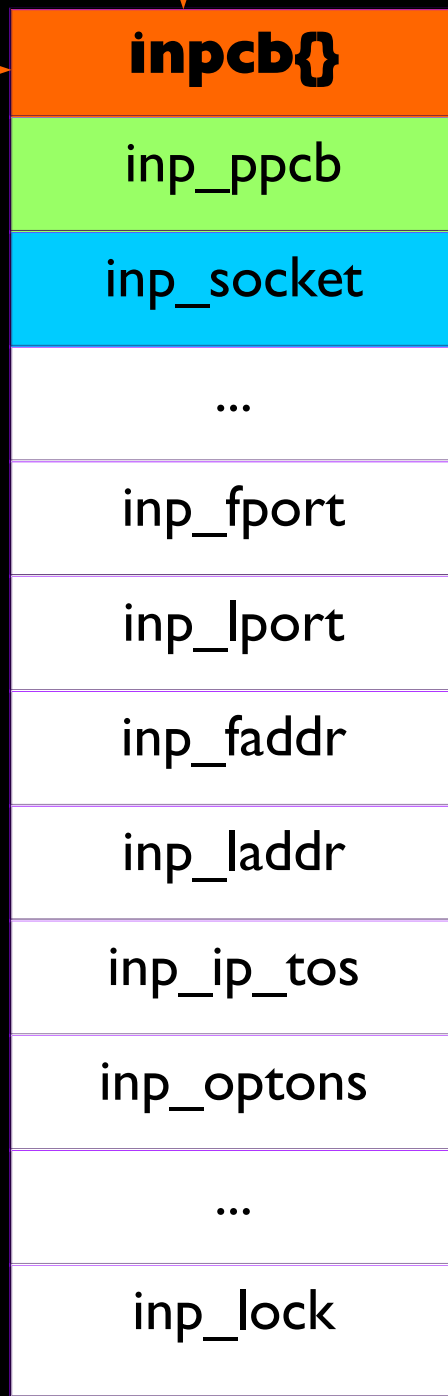
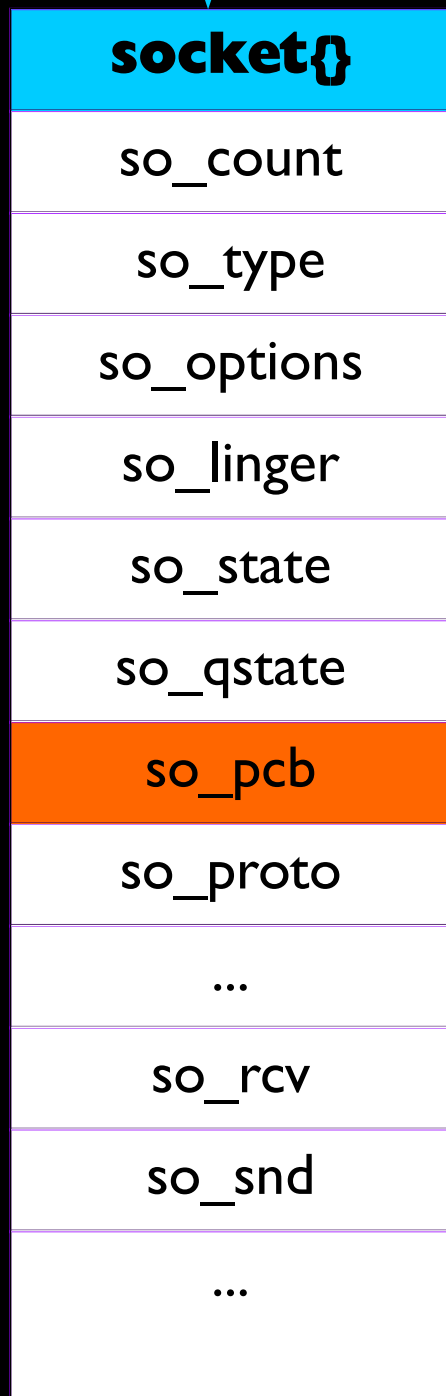
- Foreign and local port numbers

- Back pointer to socket

- Per-protocol pcb

TCP control block (tcpcb)

- Protocol state information



# Bibliography

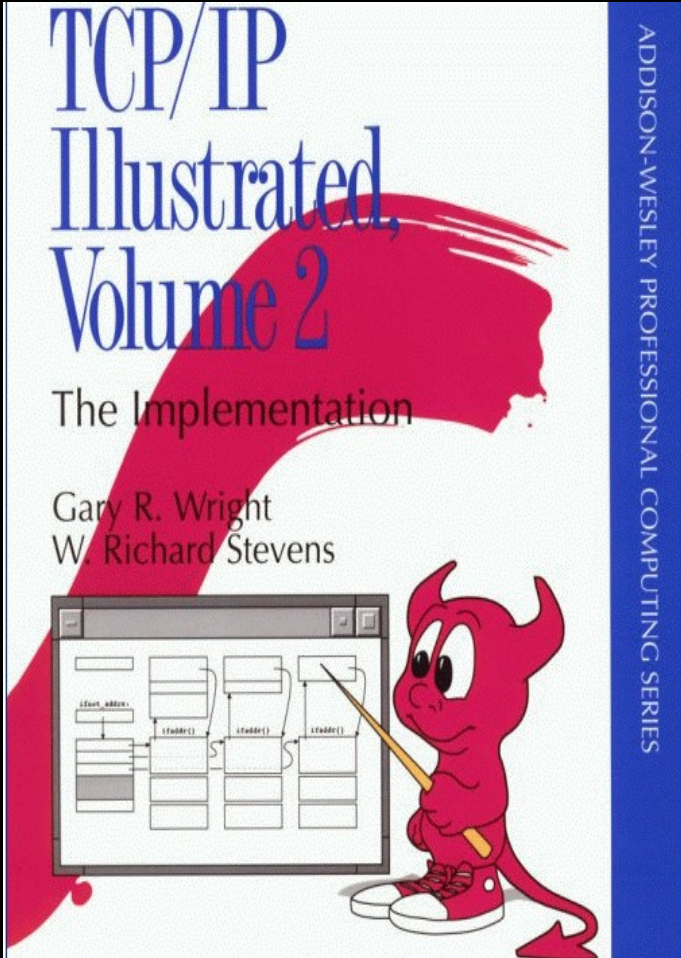


The book cover features a blue background with several green and brown cartoon turtles swimming in the water. The title is written in a large, white, serif font.

## IPv6 Core Protocols Implementation

Qing Li · Tatsuya Jinsei · Keiichi Shima

MIT



The book cover has a white background with a large, thick, red brushstroke that curves across the middle. A cartoon red devil character with horns and a tail is pointing a yellow pencil at a computer monitor. The monitor displays a network diagram with boxes and arrows, and labels like 'start\_addr()', 'start\_addr()', 'start\_addr()', and 'start\_addr()'. The title is in a blue serif font.

## TCP/IP Illustrated, Volume 2

The Implementation

Gary R. Wright  
W. Richard Stevens

ADDISON-WESLEY PROFESSIONAL COMPUTING SERIES



The book cover has a yellow background. A cartoon red devil character with horns and a tail is holding a pair of silver pliers. The title is in a mix of orange and white serif fonts. A yellow circle on the left contains the text 'FreeBSD version 5.2'. The authors' names are at the bottom. A small white star symbol is in the top right corner.

## The Design and Implementation of the FreeBSD Operating System

FreeBSD  
version  
5.2

Marshall Kirk McKusick  
George V. Neville-Neil