# Sheet 3

# sock structure

```
1   /*
2    * INET    An implementation of the TCP/IP protocol suite for the LINUX
3    *         operating system.  INET is implemented using the BSD Socket
4    *         interface as the means of communication with the user level.
5    *
6    *         Definitions for the AF_INET socket handler.
7    *
8    * Version:  @(#)sock.h    1.0.4    05/13/93
9    *
10   * Authors:  Ross Biro, <bir7@leland.Stanford.Edu>
11   *           Fred N. van Kempen, <waltje@uWalt.NL.Mugnet.ORG>
12   *           Corey Minyard <wf-rch!minyard@relay.EU.net>
13   *           Florian La Roche <flla@stud.uni-sb.de>
14   *
15   *           This program is free software; you can redistribute it and/or
16   *           modify it under the terms of the GNU General Public License
17   *           as published by the Free Software Foundation; either version
18   *           2 of the License, or (at your option) any later version.
19   */
20
```

21 **THIS IS A COMPLEX DATA STRUCTURE, WHICH CONTAINS STUFF THAT DOESN'T REALLY BELONG HERE, BUT WHICH IS**
22 **HERE FOR HISTORICAL REASONS. I HAVE CHANGED THE ORDER OF THIS SLIGHTLY SO THAT IT IS MORE LOGICAL AND I**
23 **HAVE DELETED QUITE A LOT OF IMPORTANT STUFF (.e.g all the locking code and much of the TCP related code), TO SHOW**
24 **THE BASIC STRUCTURE MORE CLEARLY.**
25

26   This structure is initialised in the following sequence:

27

28   At the end of net/ipv4/af_inet.c you will see a call to module_init(inet_init). As described in the definition of
29   include/linux/init.h::module_init, this is a marker for a driver initialisation point, which is called when the kernel boots or
30   when the module is loaded.

31

32   net/ipv4/af_inet.c::inet_init calls net/socket.c::sock_register. This latter routine is called by all protocol handlers that want to
33   advertise their address family. It creates one entry per address family in net/socket.c::net_families[family] of type
34   include/linux/net.h::net_proto_family. This has a field 'create' which is used to create a socket of that given family type. In
35   this case, this routine is set to point to net/ipv4/af_inet.c::inet_create

36

37   Socket creation: net/socket.c::sock_create -> calls create on the appropriate net_proto_family. In our case, this will call
38   through to the net/ipv4/af_inet.c::inet_create, as stored above. That initialises the sock datastructure, partly directly and
39   partly by calling net/core/sock.c::sock_init_data

40

41   `struct sock {`

42

43   The following are the source and destination information that must be entered into each IP packet. There appear to be two
44   sender addresses. rcv_saddr is the one used by hash lookups, and saddr is used for transmit. In the BSD API these are
45   almost always the same.

```
46       /* Socket demultiplex comparisons on incoming packets. */
47       __u32                   daddr;        /* Foreign IPv4 addr           */
48       __u32                   rcv_saddr;    /* Bound local IPv4 addr       */
```

```
49          __u32                   saddr;          /* Sending source                */
50          __u16                   dport;          /* Destination port              */
51          __u16                   sport;          /* Source port                   */
52      unsigned short              num;            /* Local port                    */
53
```

54  The next and prev components link sockets with the same hash value in the various socket hash tables. So, for example, in
55  net/ipv4/udp.c you find a definition of udp_hash, which is hashed on a port number. This is an open hash table of struct
56  socks which use linked lists, linked on the next and pprev values below.

```
57          /* Main hash linkage for various protocol lookup tables. */
58          struct sock             *next;
59          struct sock             **pprev;
60
```

61  TCP uses both the next and pprev fields above and the bind_next and bind_pprev and prev fields below for local binding
62  TCP hash as well as for fast bind/connect.

```
63          struct sock             *bind_next;
64          struct sock             **bind_pprev;
65          struct sock             *prev;
66
```

67  In our case this will be PF_INET

```
68      unsigned short              family;      /* Address family                   */
69
```

70  type is as for socket structure i.e. SOCK_STREAM, SOCK_DGRAM, SOCK_RAW

```
71      unsigned short              type;
72
```

73  Operation vector for the protocol with which this socket is associated. In this case, can be net/ipv4/tcp_ipv4.c::tcp_prot,
74  net/ipv4/udp.c::udp_prot, or net/ipv4/raw.c::raw_prot
75      struct proto          *prot;
76
77  In our case include/linux/in,.h::IPPROTO_TCP,  include/linux/in,.h::IPPROTO_UDP, or include/linux/in,.h::IPPROTO_IP
78      unsigned char          protocol;
79
80  State is dependent on protocol – main use is to drive TCP protocol state machine e.g. look for the enum with
81  TCP_ESTABLISHED in it in include/linux/tcp.h
82      volatile unsigned char state;       /* Connection state          */
83
84  Used when waiting for something to happen with this socket, e.g. waiting for connect in
85  net/ipv4/af_inet.c::inet_wait_for_connect,  net/ipv4/tcp.c::wait_for_tcp_connect and waiting for memory as in
86  net/ipv4/tcp.c::wait_for_tcp_memory
87      wait_queue_head_t     *sleep;       /* Sock wait queue           */
88
89      struct dst_entry      *dst_cache;  /* Destination cache         */
90
91  Packet queues. Note that there is also an error_queue, which I removed, but it's rarely used. See, for example,
92  net/ipv4/udp.c::udp_queue_rcv_skb in which a call is made to include/net/sock.h::sock_queue_rcv_skb. You can see the
93  write queue in use in net/ipv4/tcp_output.c::tcp_send_skb
94      struct sk_buff_head   receive_queue;   /* Incoming packets               */
95      struct sk_buff_head   write_queue;     /* Packet sending queue           */
96

97   **Space allocation variables.**
```
98       atomic_t                rmem_alloc;        /* Receive queue bytes committed   */
99       atomic_t                wmem_alloc;        /* Transmit queue bytes committed  */
100      atomic_t                omem_alloc;        /* "o" is "option" or "other" */
101      int                     wmem_queued;       /* Persistent queue size           */
102      int                     forward_alloc;     /* Space allocated forward.        */
```
103  **Allocation is the priority with which memory is requested for this socket**
```
104      unsigned int    allocation;                /* Allocation mode                 */
105
```
106  **Maximum amount of memory that can be requested for this socket when sending or receiving packets**
```
107      int             rcvbuf;                    /* Size of receive buffer in bytes */
108      int             sndbuf;                    /* Size of send buffer in bytes    */
109
110
```
111  **A non zero value means that we are allowed to reuse port numbers for ports that are in the TIME_WAIT state.**
```
112      unsigned char           reuse;             /* SO_REUSEADDR setting            */
113
```
114  **This says something about the way we are shutting down.**
```
115      unsigned char           shutdown;
116
```
117  **The volatile keyword is used when we have something that might change as a result of an external event, and where the**
118  **compiler will reuse the physical address rather than optimising access. E.g.  if my code looks like**

119      **A = sk->dead;**

120      **B = sk->dead;**

121 then the compiler will do both dereferences. If dead was not volatile, the compiler would normally optimise this to
122 A = B = sk->dead i.e. it would only do one dereference of sk. This is not helpful if its value is changes by an external agency
123 in between A's access and B's. In any case, these are various options that can be set for a socket.

```
124     volatile char          dead, done, urginline, keepopen, linger, destroy,
125                            no_check, broadcast, bsdism;
126     unsigned long          lingertime;
127
```

128 SO_TIMESTAMP option – if enabled then recvmsg returns a timestamp corresponding to when datagram was received.

```
129     unsigned char          rcvtstamp;
130
```

131 Says something about the features of the network device, like whether it can do the checksumming of TCP/UDP packets,
132 and whether it can DMA. Look for NETIF_F_* in include/linux/netdevice.h:net_device

```
133     int                    route_caps;
134
```

135 The proc variable is used to contain a process or process group which will be sent a signal on receipt of out-of-band data

```
136     int                    proc;
137
```

138 Used when we have peered sockets, such as with unix (local) sockets. See e.g. net/unix/af_unix.c

```
139     struct sock            *pair;
140
```

141 A process may 'lock' socket state so that it can't be changed. In particular this means that it can't be changed by bottom
142 half (interrupt driven) handlers i.e. arriving packets are blocked so we don't get any new data or changes to the state here.
143 Whilst locked, bottom half processing can add packets to the backlog queue.

```
144          /* The backlog queue is special, it is always used with
145           * the per-socket spinlock held and requires low latency
146           * access.  Therefore we special case its implementation.
147           */
148          struct {
149              struct sk_buff *head;
150              struct sk_buff *tail;
151          } backlog;
152
```

153 **tcp stuff – there's more stuff that I've deleted and some of the options described above only really apply to TCP**

```
154          union {
155              struct tcp_opt    af_tcp;
156 #if defined(CONFIG_INET) || defined (CONFIG_INET_MODULE)
157              struct raw_opt    tp_raw4;
158 #endif
159          } tp_pinfo;
160
161      int                     hashent;
162
```

163 **Error conditions**

```
164      int                     err, err_soft; /* Soft holds errors that don't
165                                                 cause failure but are the cause
166                                                 of a persistent failure not just
167                                                 'timed out' */
168
```

169 **backlog is the second parameter to the listen routine. It represents the maximum number of pending connections there can**
170 **be. Here, max_ack_backlog is this number and ack_backlog is a count of the number of connections pending at any given**
171 **time. The latter is manipulated using helper routines in include/net/tcp.h**

```
172        unsigned short         max_ack_backlog;
173        unsigned short         ack_backlog;
174
```

175 **Used to set the TOS field. Packets with a higher priority may be processed first, depending on the device's queueing**
176 **discipline. See SO_PRIORITY**

```
177        __u32                  priority;
```

178 **Route locally only if set – set by SO_DONTROUTE option.**

```
179        unsigned char          localroute;        /* Route locally only                 */
```

180 **From SO_PEERCRED option**

```
181        struct ucred           peercred;
```

182 **From SO_RCVLOWAT**

```
183        int                    rcvlowat;
```

184 **From SO_RCVTIMEO**

```
185        long                   rcvtimeo;
```

186 **From SO_SNDTIMEO**

```
187        long                   sndtimeo;
188
```

189 **Private data for each address family (truncated)**

```
190        /* This is where all the private (optional) areas that don't
191         * overlap will eventually live.
```

```
192          */
193         union {
194               void                  *destruct_hook;
195               struct unix_opt   af_unix;
196 #if defined(CONFIG_INET) || defined (CONFIG_INET_MODULE)
197               struct inet_opt   af_inet;
198 #endif
199         } protinfo;
200
```

201 Timer functions. You'll find a lot of useful timer stuff in include/linux/timer.h and kernel/timer.c In this case, the timer is used
202 for SO_KEEPALIVE (i.e. sending occasional keepalive probes to a remote site – by default, set to 2 hours in
203 include/net/tcp.h). stamp is simply the time that the last packet was received.

```
204         /* This part is used for the timeout functions. */
205         struct timer_list     timer;        /* This is the sock cleanup timer. */
206         struct timeval        stamp;
207
```

208 A backpointer to the enclosing include/linux/net.h::socket structure.

```
209         /* Identd and reporting IO signals */
210         struct socket         *socket;
211
```

212 The `state_change` operation is called whenever the status of the socket is changed. Similarly, `data_ready` is called
213 when data have been received, `write_space` when free memory available for writing has increased and
214 `error_report` when an error occurs.

```
215         /* Callbacks */
```

```
216        void                       (*state_change)(struct sock *sk);
217        void                       (*data_ready)(struct sock *sk,int bytes);
218        void                       (*write_space)(struct sock *sk);
219        void                       (*error_report)(struct sock *sk);
220
221        int                        (*backlog_rcv) (struct sock *sk, struct sk_buff *skb);
222
223    Get rid of the socket.
224        void                       (*destruct)(struct sock *sk);
225    };
```