

# Ch.7: Socket Options

- *getsockopt* and *setsockopt* functions
- Check options and obtain default values
- Generic socket options
- IPv4 socket options
- IPv6 socket options
- TCP socket options
- *fcntl* function
- Summary

# *getsockopt* and *setsockopt* Functions

```
#include <sys/socket.h>
int getsockopt (int sockfd, int level, int optname, void *optval,
                socklen_t *optlen);
int setsockopt (int sockfd, int level, int optname, const void *optval,
                socklen_t optlen);
```

Both return: 0 if OK, -1 if error

Types of options: flag and value

Levels:

Generic --	SOL_SOCKET
IP --	IPPROTO_IP
ICMPv6 --	IPPROTO_ICMPV6
IPv6 --	IPPROTO_IPV6
TCP --	IPPROTO_TCP

# Check Options and Obtain Default Values

- Program declaration:
  - declare union of possible option values
  - define printing function prototypes
  - define `sock_opts` structure, initialize `sock_opt[ ]` array
- Check and print options:
  - create TCP socket, go through all options
  - call *getsockopt*
  - print option's default value

Declaration for Socket Options  
(see Figure 7.2)

Program to Check and Print Socket Options  
(see Figures 7.3 and 7.4)

# Generic Socket Options

- **SO\_BROADCAST**: permit sending of broadcast datagram (UDP), only on broadcast links
- **SO\_DEBUG**: enable the kernel to track details about packets sent or received by TCP socket, allowing *trpt* program to examine the kernel circular buffer
- **SO\_DONTROUTE**: bypass routing table lookup, used by routing daemons (*routed* and *gated*) in case the routing table is incorrect

# SO\_ERROR Socket Option (1/2)

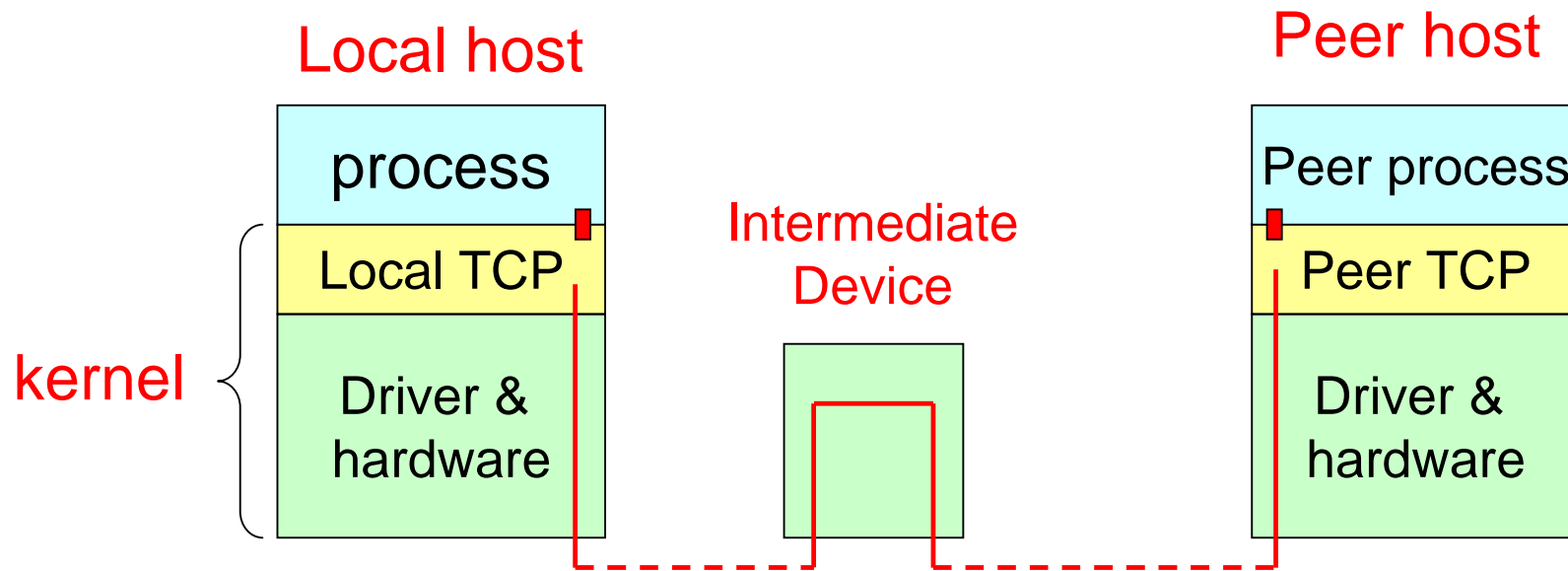
- **SO\_ERROR**: get pending error and clear
- 當 socket 發生錯誤時，kernel 會將錯誤代碼放入變數 `so_error` 中，稱為該 socket 的 pending error
  - 如 process 正 block 在 `select` 中測試該 socket 是否 readable 或 writable，則 `select` 會 return
  - 如該 process 使用 signal-driven I/O (五種 I/O model 中的第四種)，則 kernel 會產生 SIGIO signal 給 process
- 無論是哪一種情形，process 都可以透過 `SO_ERROR` option 取得 `so_error` 的值

# SO\_ERROR Socket Option (2/2)

- `so_error`的值被取回後即會被kernel歸零
- 當process呼叫`read`時`so_error`的值不為0且沒有data可讀，則`read` return -1
  - 有data可讀時則傳回0或讀取data數，不會傳回-1。
- 當process呼叫`write`時`so_error`的值不為0，則`write` return -1 (無論是否writable)
- 此時`so_error`的值會被設給全域變數`errno`後歸零。Process只須檢查`errno`的值即可。

# SO\_KEEPALIVE Socket Option (1/2)

- **SO\_KEEPALIVE**: when enabled, test if TCP connection still alive periodically (default 2 hours, can be changed by TCP\_KEEPALIVE)





# SO\_KEEPALIVE Socket Option (2/2)

- 當任一方超過兩小時未交換資料，kernel TCP會自動發*keepalive probe* (a TCP segment)給Peer TCP
  - 如果Peer TCP回應ack，應用程式不受任何影響。
  - 如果Peer TCP回應RST (Peer process已不存在)，則local TCP close此socket且設so\_error = ECONNRESET
  - 如果沒收到peer的回應，so\_error = ETIMEOUT
  - 如果收到ICMP unreachable的錯誤，so\_error = EHOSTUNREACH

# Peer Host Vs. Peer Process

- The purpose of keepalive is to detect if the peer *host* (not the peer *process*) crashes
  - If peer process crashes, peer TCP sends an FIN
  - If local TCP sends another segment, peer TCP responds with a RST
  - If process sends yet another segment, local TCP sends local process a `SIGPIPE` signal
  - For other cases see Figure 7.5
- This option is normally used by servers to detect if client host crash (avoids *half-open*)

# SO\_LINGER Socket Option

- Specifies how `close` operates for TCP
- By default, `close` returns immediately
  - While the kernel will try to deliver any data remaining in the socket send buffer
- Application process may set this option to
  - enable TCP to abort the connection by discarding any pending data, or
  - put the process to sleep until (1) all the data is sent and acknowledged by the peer or (2) time out

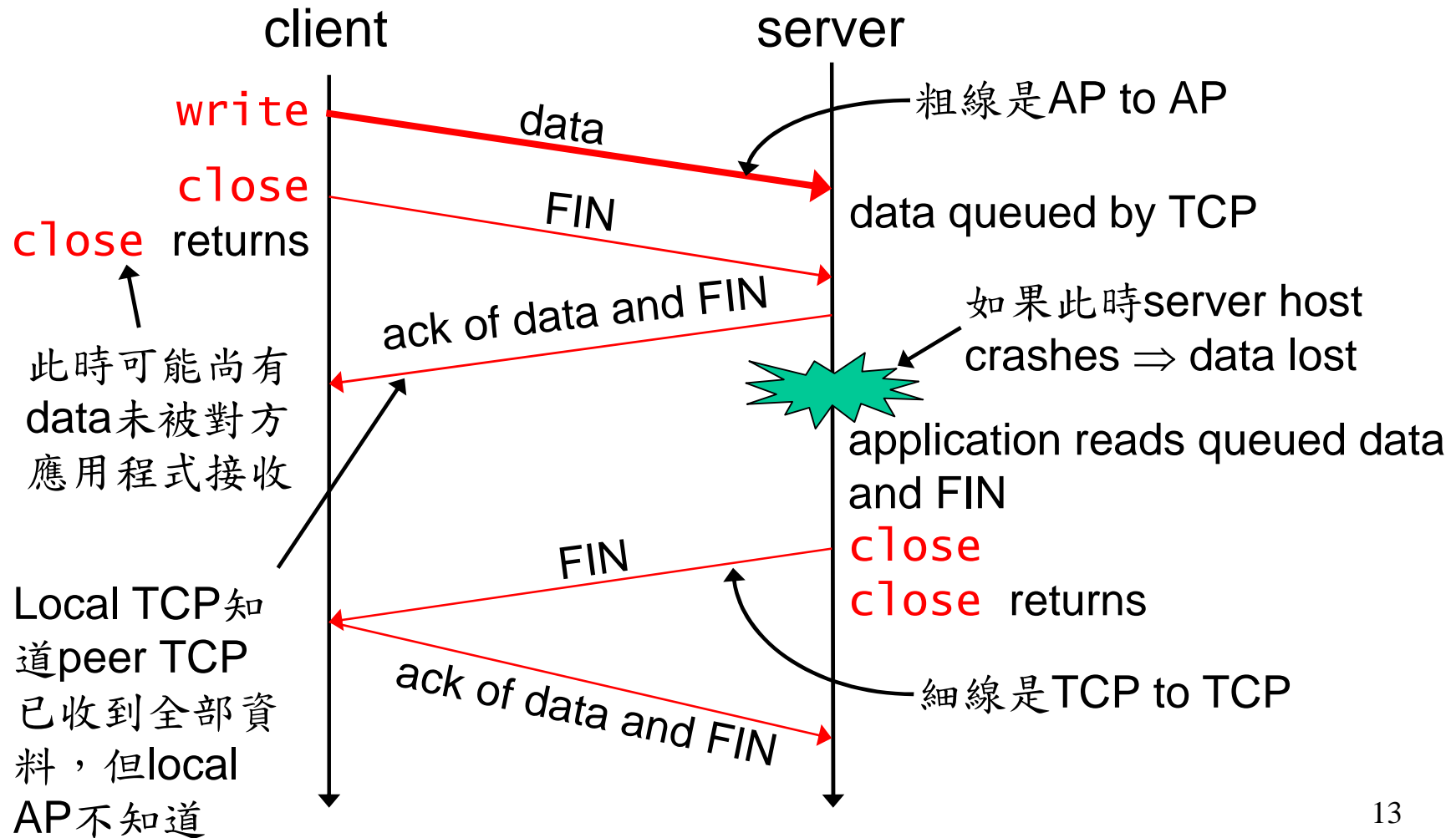
# SO\_LINGER: Data Structure

```
#include <sys/socket.h>

struct linger {
    int l_onoff; /* 0=off, nonzero=on */
    int l_linger; /* linger time */
};
```

- ① `l_onoff = 0`: `close` returns immediately (default)
- ② `l_onoff = 1` and `l_linger = 0`: aborts connection
- ③ `l_onoff = 1` and `l_linger ≠ 0`: linger on `close`

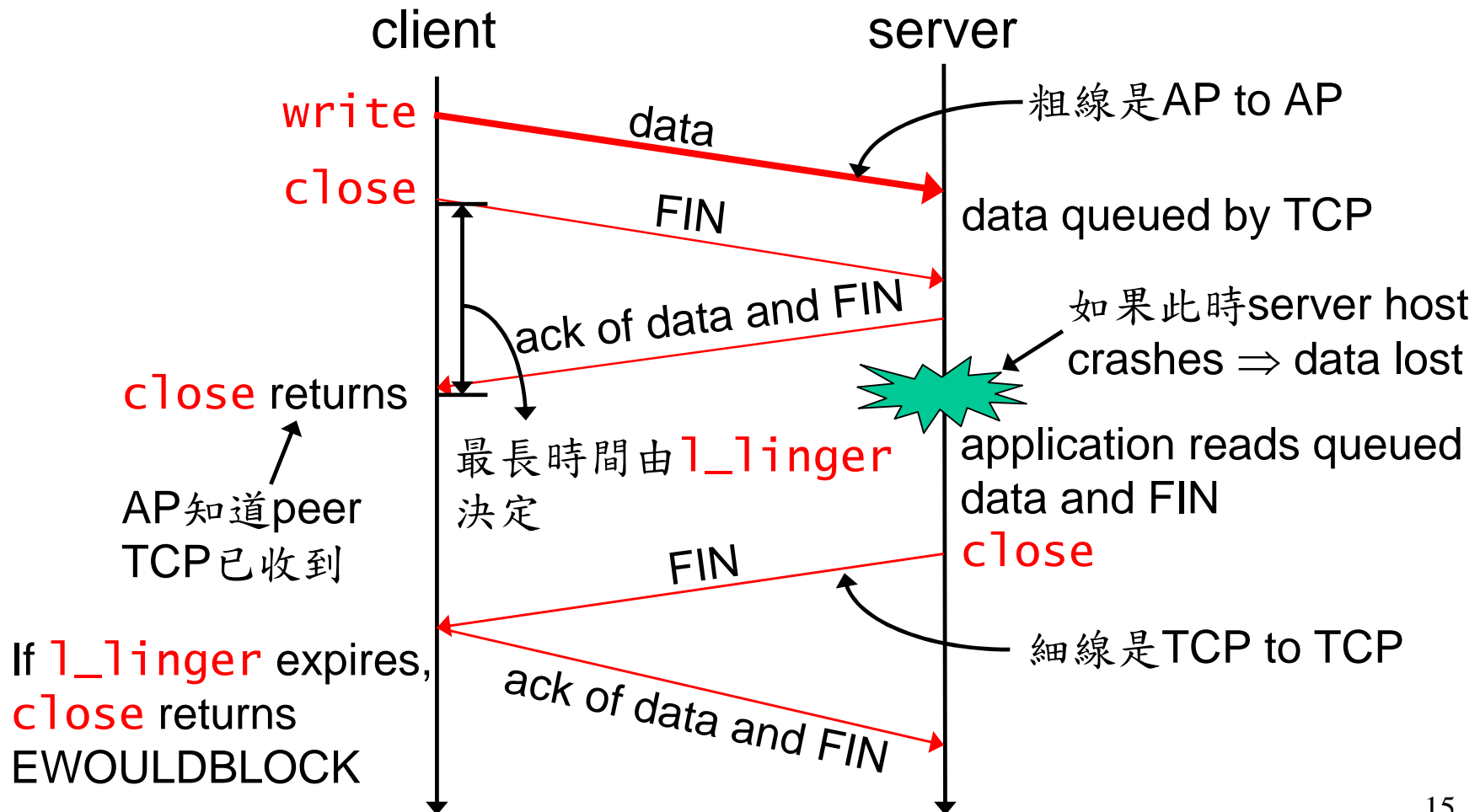
# Option 1: $\_onoff = 0$



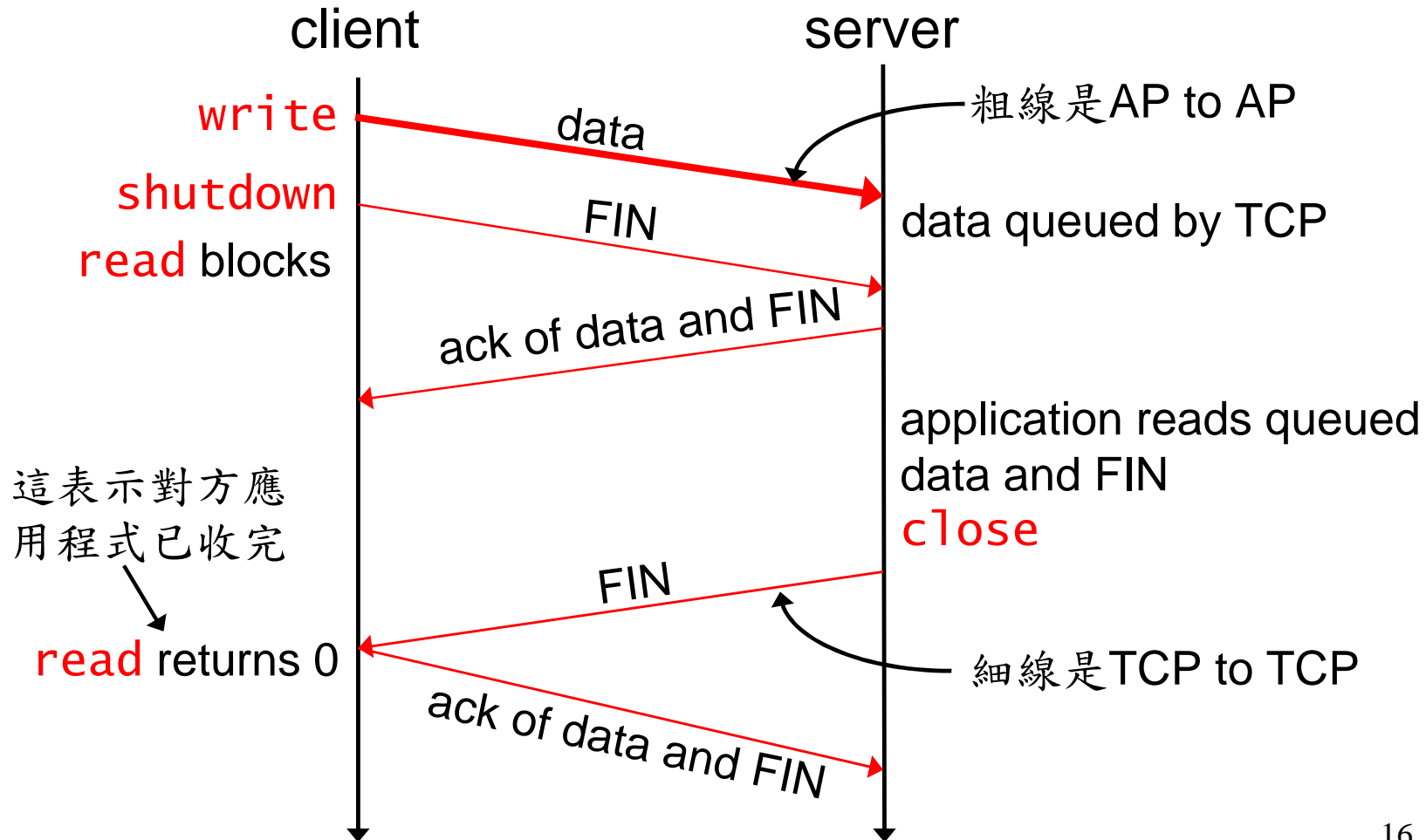
Option 2:  $\text{TCP\_linger} = 1$  And  
 $\text{TCP\_linger\_onoff} = 0$

- TCP discards any data still remaining in the socket send buffer (這些data對應的ack都尚未收到), and sends an RST to the peer TCP
- This avoids TCP's TIME\_WAIT state, which may cause problems.

# Option 3: $\text{TCP\_linger} = 1$ And $\text{TCP\_linger} \neq 0$



# To Know That Peer *Process* Has Read Our Data: Using **shutdown**





# When We Close Our End of the Connection...

- The return can be occurred at three different times
  - `close` returns immediately, without waiting at all (AP knows nothing)
  - `close` lingers until the ACK of our FIN is received (peer TCP has received all data)
  - `shutdown` followed by a `read` waits until we receive the peer's FIN (peer process has received all data)

## Another Alternative: Using Application-Level Acknowledgement

```
char ack;
```

```
write(sockfd, data, nbytes);
```

```
n = Read(sockfd, &ack, 1);
```

```
nbytes = Read(sockfd, buff, sizeof(buff));
```

```
/* server 確定已收妥全部資料 */
```

```
write(sockfd, "", 1);
```

Server's ACK back to client

# Generic Socket Options (Cont.)

- **SO\_OOBINLINE**: leave received out-of-band data inline in the normal input queue
- **SO\_RCVBUF/SO\_SNDBUF**: socket receive / send buffer size, TCP default: 8192-61440, UDP default: 40000/9000
- **SO\_RCVLOWAT/ SO\_SNDLOWAT**: receive / send buffer low water mark for *select* to return

# Generic Socket Options (Cont.)

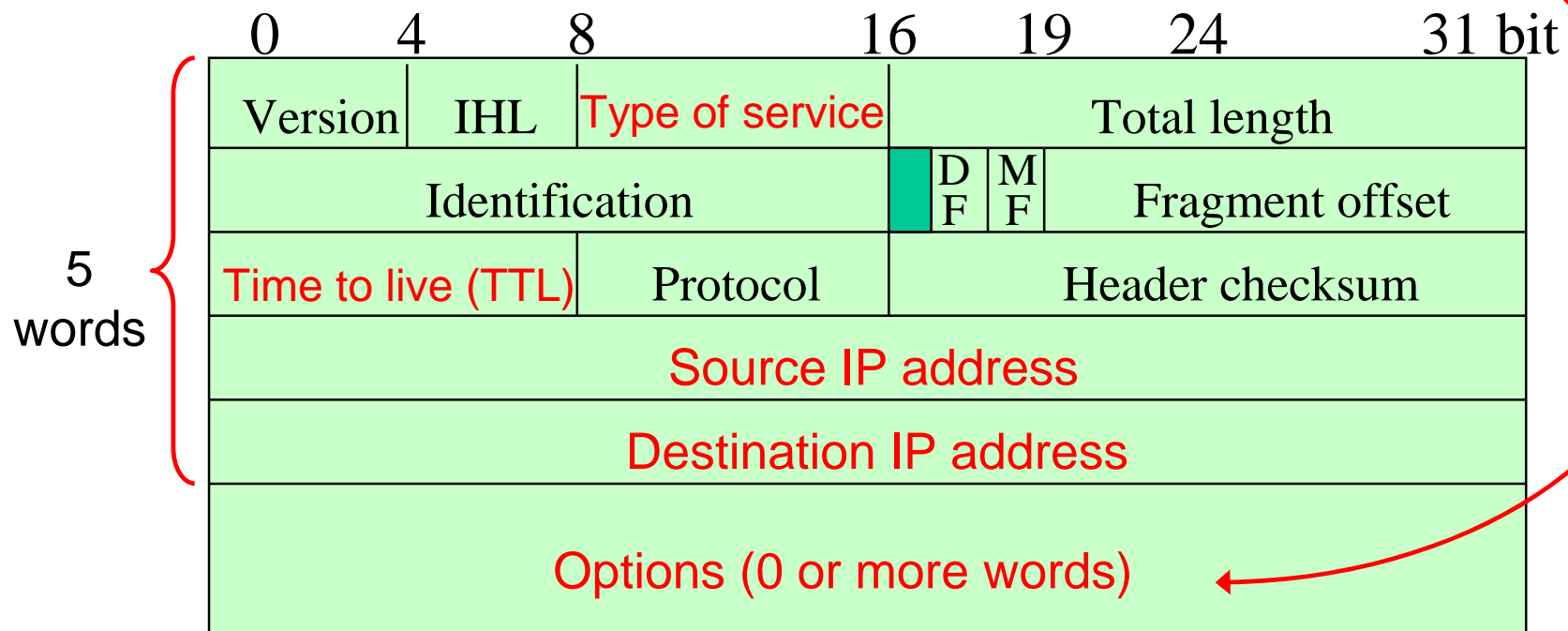
- **SO\_RCVTIMEO/SO\_SNDTIMEO**: receive / send timeout for socket read/write
- **SO\_REUSEADDR/SO\_REUSEPORT**: allow local address reuse for TCP server restart, IP alias, UDP duplicate binding for multicasting
- **SO\_TYPE**: get socket type (SOCK\_STREAM or SOCK\_DGRAM)
- **SO\_USELOOPBACK**: applies only to routing socket; gets copy of what it sends

# IPv4 Socket Options

- **IP\_HDRINCL**: If set for a raw IP socket, we must build our own IP header for all the datagrams that we send on the raw socket.
  - Normally the kernel builds the IP header
- e.g. *traceroute* builds its own IP header on a raw socket

# IPv4 Socket Options (Cont.)

- **IP\_OPTIONS**: specify IP options in the IPv4 header (source route, timestamp, etc.)



# IPv4 Socket Options (Cont.)

- **IP\_RECVSTADDR**: return destination IP address of a received UDP datagram by **recvmsg** (Ch. 8)
- **IP\_RECVIF**: return received interface index for a received UDP datagram by **recvmsg**
- **IP\_TOS**: set IP Type of Service (TOS) field of outgoing packets for TCP/UDP socket.
  - options: IPTOS\_LOWDELAY, IPTOS\_THROUGHPUT, IPTOS\_RELIABILITY, IPTOS\_LOWCOST

# IP Socket Options (Cont.)

- **IP\_TTL**: set and fetch the default TTL for outgoing packets, 64 for TCP/UDP sockets, 255 for raw sockets, used in *traceroute*
- **IP\_MULTICAST\_IF**, **IP\_MULTICAST\_TTL**, **IP\_MULTICAST\_LOOP**, **IP\_ADD\_MEMBERSHIP**, **IP\_DROP\_MEMBERSHIP** (Sec. 19.5)



# IPv6 Socket Options

- **ICMP6\_FILTER**: fetch and set `icmp6_filter` structure specifying message types to pass
- **IPV6\_ADDFORM**: change address format of socket between IPv4 and IPv6
- **IPV6\_CHECKSUM**: offset of checksum field for raw socket
- **IPV6\_DSTOPTS**: return destination options of received datagram by `recvmsg`
- **IPV6\_HOPLIMIT**: return hop limit of received datagrams by `recvmsg`

# IPv6 Socket Options (Cont.)

- **IPV6\_HOPOPS**: return hop-by-hop options of received datagrams by **recvmsg**
- **IPV6\_NEXTHOP**: specify next hop address as a socket address structure for a datagram
- **IPV6\_PKTINFO**: return packet info, dest IPv6 address and arriving interface, of received datagrams by **recvmsg**

# IPv6 Socket Options (Cont.)

- **IPV6\_PKTOPTIONS**: specify socket options of TCP socket (UDP uses **recvmsg** and **sendmsg**)
- **IPV6\_RTHDR**: receive routing header (source route)
- **IPV6\_UNICAST\_HOPS**: ~ IP\_TTL
- IPV6\_MULTICAST\_IF/HOPS/LOOP, IPV6\_ADD/DROP\_MEMBERSHIP (Sec. 19.5)

# TCP Socket Options

- **TCP\_KEEPALIVE**: seconds between probes
- **TCP\_MAXRT**: TCP max retx time
- **TCP\_MAXSEG**: TCP max segment size
- **TCP\_NODELAY**: disable Nagle algorithm, to reduce the number of small packets
- **TCP\_STDURG**: interpretation of TCP's urgent pointer, used with out-of-band data

# fcntl Function

- Calling `fcntl` is a preferred way to
  - Set socket for non-blocking I/O (第二種I/O)
  - Set socket for signal-driven I/O (第四種I/O)
  - Set socket owner
    - Socket owner 是接收 SIGIO 和 SIGURG 信號的 process
  - Get socket owner

# fcntl Function Prototyping

```
#include <fcntl.h>
int fcntl (int fd, int cmd, ... /* int arg */ );
           returns: depends on cmd if OK, -1 on error
```

Two flags that affect a socket:

(fetch by **F\_GETFL**, set by **F\_SETFL**)

**O\_NONBLOCK** (nonblocking I/O) 第15章

**O\_ASYNC** (signal-driven I/O notification)

To set a flag : ①fetch ②**OR/AND~** ③set

```
int flags;
if ( (flags = fcntl (fd, F_GETFL, 0) ) < 0) error_sys ("F_GETFL error");
flags |= O_NONBLOCK; /* turn on */
( or flags &= ~O_NONBLOCK; /* turn off */ )
if (fcntl(fd, F_SETFL, flags) < 0) error_sys ("F_SETFL error");
```

# Set and Get Socket Owner

- **F\_SETOWN** command lets us set socket owner (process ID or process group ID) to receive **SIGIO** and **SIGURG** signals
  - **SIGIO** is generated if signal-driven I/O is enabled for a socket
  - **SIGURG** is generated when a out-of-band data arrives for a socket.
- **F\_GETOWN** command gets socket owner

# Summary

- Commonly used socket options: **SO\_KEEPALIVE**, **SO\_RCVBUF**, **SO\_SNDBUF**, **SO\_REUSEADDR**
- **SO\_REUSEADDR** set in TCP server before calling `bind`
- **SO\_LINGER** gives more control over when close returns and forces RST to be sent
- **SO\_RCVBUF/SO\_SNDBUF** for bulk data transfer across long fat pipes