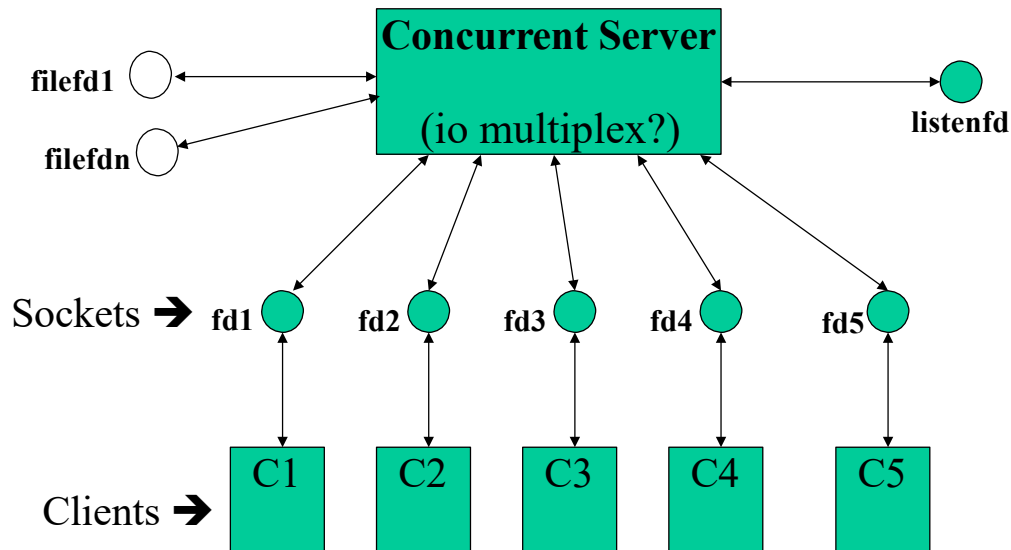# I/O Multiplexing

# What is I/O Multiplexing?

- When an application needs to handle multiple I/O descriptors at the same time
  - E.g., file and socket descriptors, or multiple socket descriptors
- Application processes both *interactive* input and *network* socket
- Server handles multiple ports and protocols
  - Server handles both TCP listen and connected sockets
  - Server handles both TCP and UDP
- When I/O on any one descriptor can result in *blocking*
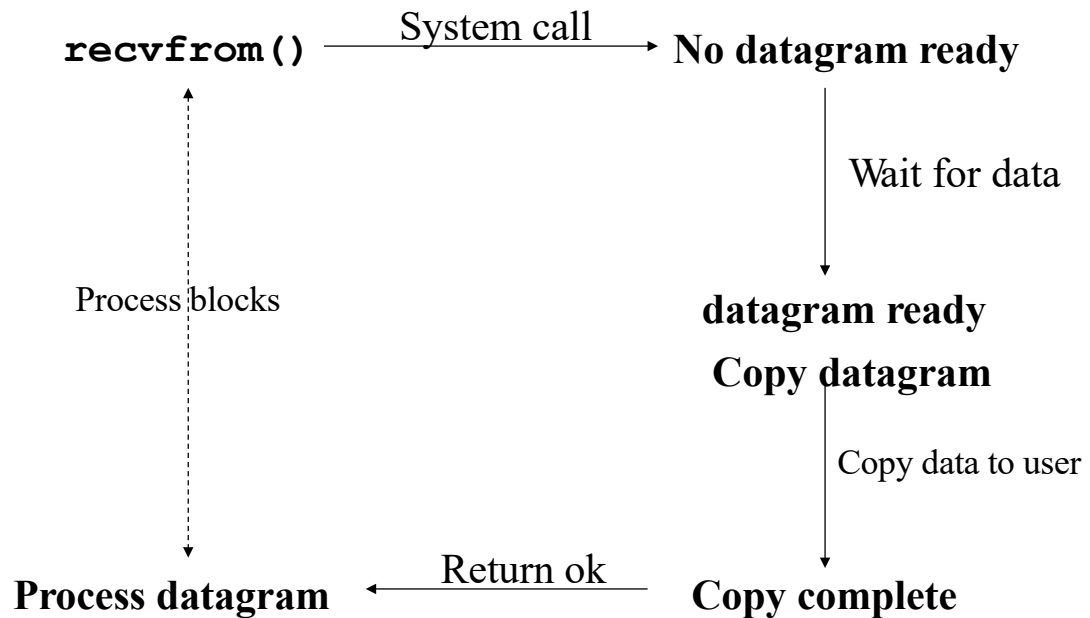
# Non-Forking Concurrent Server

**filefd1** ◯ ⟷ **Concurrent Server**

**Concurrent Server**

**(io multiplex?)** ⟷ ● **listenfd**

**filefdn** ◯

Sockets ➜ **fd1** ●    **fd2** ●    **fd3** ●    **fd4** ●    **fd5** ●

Clients ➜   C1    C2    C3    C4    C5

---

# I/O models

1. Blocking I/O
2. Non-blocking I/O
3. I/O multiplexing – select(), poll()
4. Signal driven I/O
5. Asynchronous I/O
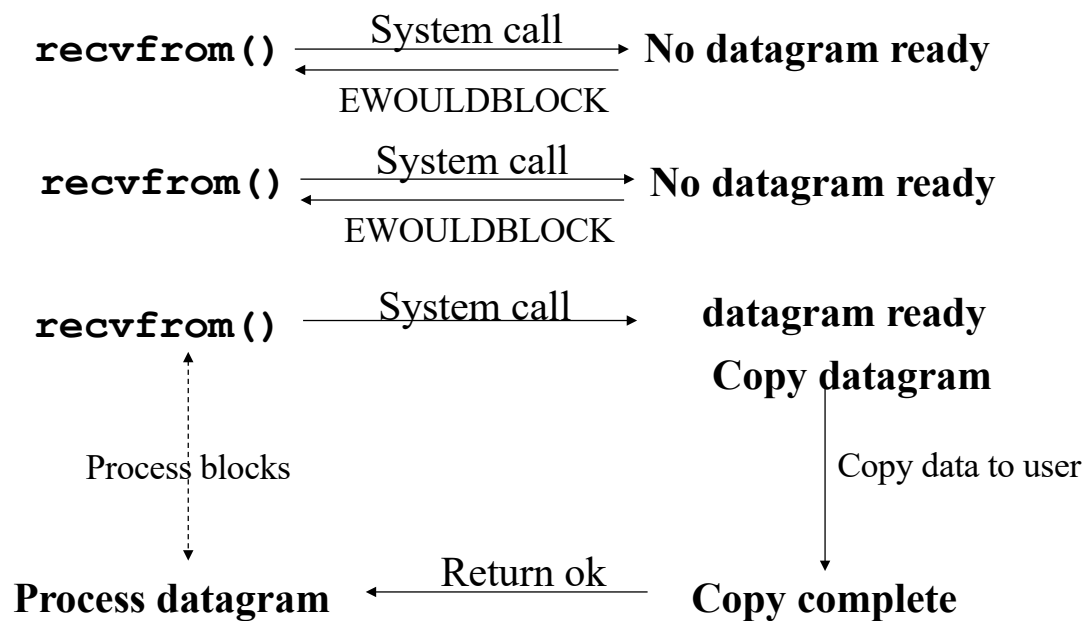
# Blocking I/O

Application                           Operating system

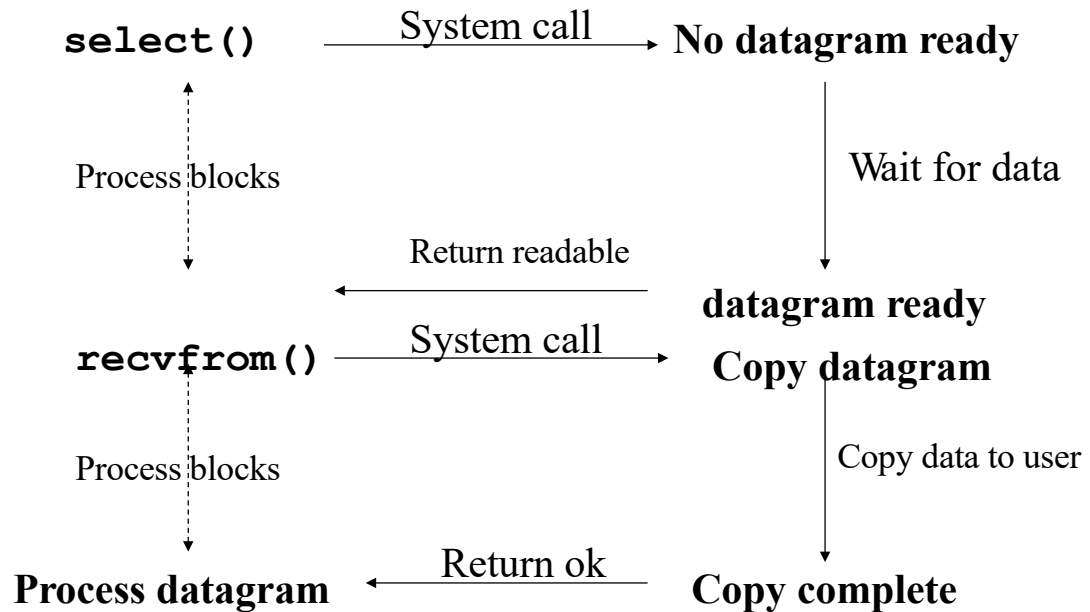**recvfrom()** — System call → **No datagram ready**

Wait for data

**datagram ready**

**Copy datagram**

Copy data to user

Process blocks

**Process datagram** ← Return ok — **Copy complete**

# Non-Blocking I/O

Application                       Operating system

**recvfrom()** — System call → **No datagram ready**
← EWOULDBLOCK

**recvfrom()** — System call → **No datagram ready**
← EWOULDBLOCK

**recvfrom()** — System call → **datagram ready**

**Copy datagram**

Copy data to user

Process blocks

**Process datagram** ← Return ok — **Copy complete**

# I/O Multiplexing

Application            Operating system

**select()** — System call → **No datagram ready**

Process blocks

Wait for data

← Return readable

**datagram ready**

**recvfrom()** — System call → **Copy datagram**

Process blocks

Copy data to user

**Process datagram** ← Return ok — **Copy complete**

# Signal driven I/O

Application            Operating system

Establish SIGIO
Signal handler — System call → No datagram ready

Process continues

Wait for data

Signal Handler ← Deliver SIGIO

datagram ready

**recvfrom()** — System call → Copy datagram

Process blocks

Copy data to user

Process datagram ← Return ok — Copy complete

# Asynchronous I/O

Application

Operating system

**aio_read()**  ⟶ System call ⟶  **No datagram ready**
⟵ return

Process continues

Wait for data

**datagram ready**

**Copy datagram**

Copy data to user

Signal handler
Process datagram  ⟵ Return ok  **Copy complete**

---

# select() call

- Allows a process to *wait* for an event to occur on any one of its descriptors:
  - Wait forever
  - Wait for a fixed amount of time
  - Do not wait at all! This is called *polling*

- Types of *events*
  - Ready for read
  - Ready for write
  - Exception condition

# select() call

```
int  select(int maxfdp1,        /* max. fd + 1 */
            fd_set *readfds,    /* read ready? */
            fd_set *writefds,   /* write ready? */
            fd_set *exceptfds,  /* exception? */
            struct timeval *timeout);
```

Returns positive count of ready descriptors, 0 on timeout, –1 on error

```
struct  timeval {
  long tv_sec;   /* seconds */
  long tv_usec; /* microseconds */
};
```

# fd_set

- Describes a set of descriptors that we want to wait on for events
- Traditionally held 1024 descriptors, but should be configurable for newer kernels
- Defines manipulation macros

```
void FD_ZERO(fd_set *fds)
void FD_SET(int fd, fd_set *fds)
void FD_CLR(int fd, fd_set *fds)
int FD_ISSET(int fd, fd_set *fds)
```

# Value-result arguments

- Select modifies descriptor sets pointed to by *readset*, *writeset*, and *exceptset* pointers
- On function call
  - Specify value of descriptors that we are interested in
- On function return
  - Result indicates which descriptors are ready
- Use FD_ISSET macro on return to test a specific descriptor in an fd_set structure
  - Any descriptor not ready will have its bit cleared
  - You need to turn on all the bits in which you are interested on all the descriptor sets each time you call select

# maxfdp1 argument

- Specifies the maximum number of descriptors to be tested
- Its value is the maximum descriptor plus one
  - Descriptors 0, 1, 2, up through and including *maxfdp1*-1 are tested
  - Calculate the maxfdp1 value based on the current descriptor set *for efficiency* of processing fd_set
- Constant FD_SETSIZE  defined by including header file <sys/select.h>
  - It is the number of descriptors in the fd_set datatype (traditionally 1024)

# Socket conditions
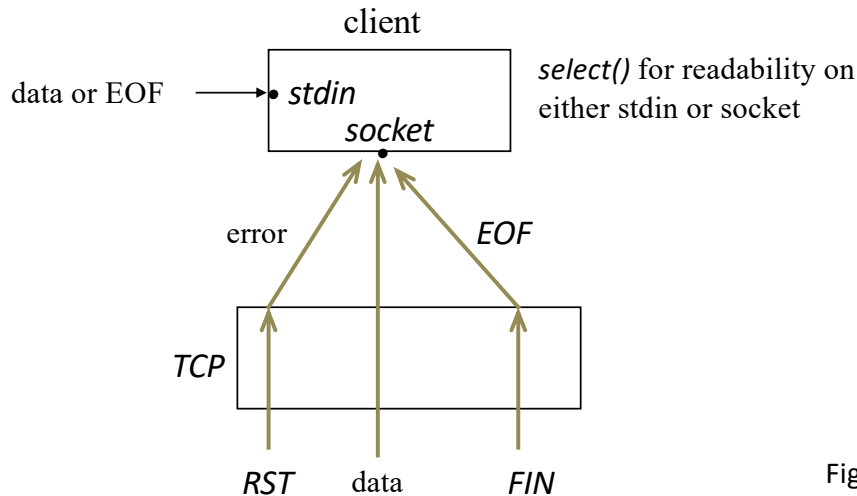
- Handling conditions where peer can send data, RST or FIN

```
                        client
                   ┌──────────────┐      select() for readability on
data or EOF ──────▶│• stdin       │      either stdin or socket
                   │     socket   │
                   │        •     │
                   └──────────────┘
                      ▲   ▲   ▲
               error  │   │   │  EOF
                   ┌──┼───┼───┼──┐
              TCP  │  │   │   │  │
                   └──┼───┼───┼──┘
                      │   │   │
                     RST data FIN          Fig. 6.8
```

# Socket conditions (2)

- Can handle conditions where peer can send data, RST or FIN
- Peer TCP sends data
  - The socket becomes readable and read returns greater than 0 (number of bytes of data)
- Peer TCP sends a FIN (terminate or close)
  - The socket becomes readable and read returns 0 (i.e., EOF)
- Peer TCP sends a RST (crashed and rebooted)
  - The socket becomes readable and returns -1
  - **errno** contains the specific error code

# Ready conditions

Summary of conditions that cause a socket to be ready for select()

| Condition | Readable? | Writable? | Exception? |
|---|:---:|:---:|:---:|
| Data to read<br>Read-half of the connection closed<br>New connection ready for listening socket | •<br>•<br>• | | |
| Space available for writing<br>Write-half of the connection closed | | •<br>• | |
| Pending error | • | • | |
| TCP out-of-band data | | | • |

# **Non-forking concurrent server**

```
fdset rdset, wrset;
int listenfd, connfd1, connfd2;
int maxfdp1;

/* do connection establishment etc */
    ……
/* initialize */
FD_ZERO(&rdset);
FD_ZERO(&wrset);
```

```
for ( ;; ) {
   FD_SET(connfd1, &rdset);
   FD_SET(connfd2, &wrset);
   FD_SET(listenfd, &rdset);

   maxfdp1 = max(connfd1, connfd2, listenfd) + 1;
   /* wait for some event */
   Select(maxfdp1, &rdset, &wrset, NULL, NULL);

   if (FD_ISSET(connfd1, &rdset)) {
       /* read data from connfd */
   }
   if (FD_ISSET(connfd2, &wrset)) {
       /* write data to connfd */
   }
   if (FD_ISSET(listenfd, &rdset)) {
       /* process a new connection */
   }
}
```

# High resolution timer!

- Select can be used as a millisecond resolution timer.

```
select(0, NULL, NULL, NULL, &timeval);
```

- Usual **sleep**() call has resolution of seconds.

# Fig 6.13 (strcliselect02.c)

```
1 #include   "unp.h"
2 void
3 str_cli(FILE *fp, int sockfd)
4 {
5   int   maxfdp1, stdineof;
6   fd_set  rset;
7   char   buf[MAXLINE];
8   int   n;
9   stdineof = 0;
10  FD_ZERO(&rset);
11  for ( ; ; ) {
12    if (stdineof == 0)
13      FD_SET(fileno(fp), &rset);
14    FD_SET(sockfd, &rset);
15    maxfdp1 = max(fileno(fp), sockfd) + 1;
16    Select(maxfdp1, &rset, NULL, NULL, NULL);

17    if (FD_ISSET(sockfd, &rset)) {  /* socket is readable */
18      if ( (n = Read(sockfd, buf, MAXLINE)) == 0) {
19        if (stdineof == 1)
20          return;     /* normal termination */
21        else
22          err_quit("str_cli: server terminated prematurely");
23      }
24      Write(fileno(stdout), buf, n);
25    }
26    if (FD_ISSET(fileno(fp), &rset)) {  /* input is readable */
27      if ( (n = Read(fileno(fp), buf, MAXLINE)) == 0) {
28        stdineof = 1;
29        Shutdown(sockfd, SHUT_WR);  /* send FIN */
30        FD_CLR(fileno(fp), &rset);
31        continue;
32      }
33      Writen(sockfd, buf, n);
34    }
35  }
36 }
```

# Multiple client descriptors (e.g.)

- Instead of fork() per client, use select()
- After the second client connection is established (assuming connected descriptor returned by accept is 5)
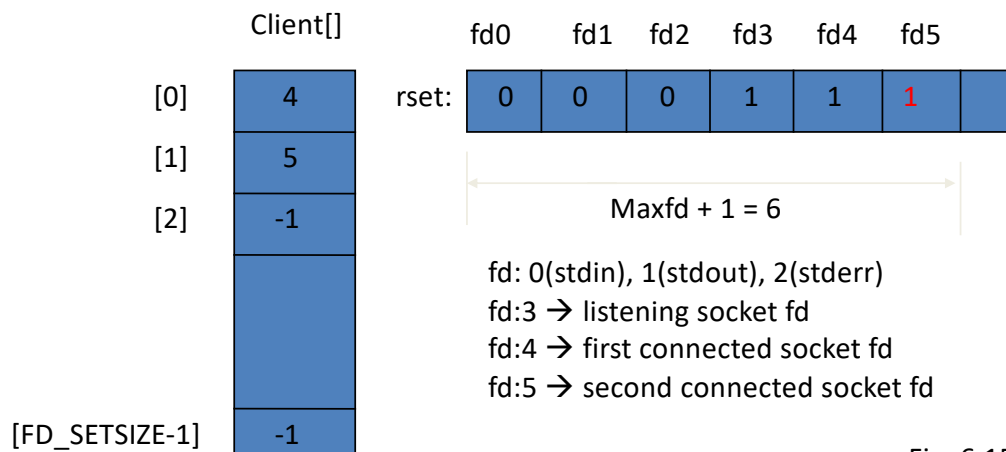


fd: 0(stdin), 1(stdout), 2(stderr)
fd:3 → listening socket fd
fd:4 → first connected socket fd
fd:5 → second connected socket fd

Fig. 6.15

# Fig. 6.22 (tcpservselect01.c)

```
25  for ( ; ; ) {
26      rset = allset;                      /* structure assignment */
27      nready = Select(maxfd + 1, &rset, NULL, NULL, NULL);

28      if (FD_ISSET(listenfd, &rset)) {    /* new client connection */
29          clilen = sizeof(cliaddr);
30          connfd = Accept(listenfd, (SA *) &cliaddr, &clilen);

31          for (i = 0; i < FD_SETSIZE; i++)
32              if (client[i] < 0) {
33                  client[i] = connfd; /* save descriptor */
34                  break;
35              }
36          if (i == FD_SETSIZE)
37              err_quit("too many clients");
38          FD_SET(connfd, &allset);        /* add new descriptor to set */
39          if (connfd > maxfd)
40              maxfd = connfd; /* for select */
41          if (i > maxi)
42              maxi = i;       /* max index in client[] array */

43          if (--nready <= 0)
44              continue;       /* no more readable descriptors */
45      }
46      for (i = 0; i <= maxi; i++) {        /* check all clients for data */
47          if ( (sockfd = client[i]) < 0)
48              continue;
49          if (FD_ISSET(sockfd, &rset)) {
50              if ( (n = Read(sockfd, buf, MAXLINE)) == 0) {
51                  /* connection closed by client */
52                  Close(sockfd);
53                  FD_CLR(sockfd, &allset);
54                  client[i] = -1;
55              } else
56                  Writen(sockfd, buf, n);

57              if (--nready <= 0)
58                  break;      /* no more readable descriptors */
59          }
60      }
61  }
62 }
```

# pselect()

- Const nanosec timespec for timeout parameter
- Set a blocking mask while waiting, restore original mask when returning

```
int   pselect(int maxfdp1,
              fd_set *readset,
              fd_set *writeset,
              fd_set *exceptset,
              const struct timespec *timout,
              const sigset_t *sigmask);


struct    timespec    {
    time_t   tv_sec;    /* seconds (long) */
    long     tv_nsec;   /* nanoseconds */
};
```

# pselect()

Race condition can exist between the call to select() and signal delivery

```
if (intr_flag)
   handle_intr();     /* handle the signal */

                                                      ←        ??
if ( (nready = select( ... )) < 0) {
   if (errno == EINTR) {
      if (intr_flag)
         handle_intr();
   }
   ...
}
```

# pselect()

Signal mask argument overrides the existing signal mask

```
sigset_t newmask, oldmask, zeromask;

sigemptyset(&zeromask);
sigemptyset(&newmask);
sigaddset(&newmask, SIGINT);

sigprocmask(SIG_BLOCK, &newmask, &oldmask);        /* block SIGINT */
if (intr_flag)
   handle_intr();    /* handle the signal */
if ( (nready = pselect ( ... , &zeromask)) < 0 ) {
   if (errno == EINTR)  {
      if (intr_flag)
         handle_intr ();
   }
   ...
}
```

# poll()

- Provides similar functionality to select()
- Events encoded in bits specifying certain conditions
- Wait forever INFTIM (-1 or < 0); or don't block 0; or wait for > 0

```
#include <poll.h>

int  poll(struct pollfd *fdarray,
          unsigned long nfds, /* array size */
          int timeout);        /* msec to wait */

struct          pollfd          {
   int          fd;                /* descriptor to check */
   short        events;           /* events of interest */
   short        revents;          /* events occurred */
};
```

# Poll() Parameters

- Have two variables per descriptor
  - Input value events member **event**
  - Returns status for a descriptor in the corresponding **revent**s member

|  | Constant | Input to events ? | Result from revents ? | Description |
|---|---|---|---|---|
| inputs | POLLIN | • | • | Normal or priority band data can be read |
|  | POLLRDNORM | • | • | Normal data can be read |
|  | POLLRDBAND | • | • | Priority band data can be read |
|  | POLLPRI | • | • | High-priority data can be read |
| outputs | POLLOUT | • | • | Normal data can be written |
|  | POLLWRNORM | • | • | Normal data can be written |
|  | POLLWRBAND | • | • | Priority band data can be written |
| errors | POLLERR |  | • | Error has occurred |
|  | POLLHUP |  | • | Hangup has occurred |
|  | POLLNVAL |  | • | Descriptor is not an open file |

- Set **fd** member < 0 to ignore it
- Note fdarray allocated by caller; no fixed-size data type similar to fd_set

# Poll Conditions (e.g.)

- Three classes of data identified:
  - Normal, priority, high-priority
- Some events for backward compatibility, e.g., POLLIN and POLLOUT
  - POLLIN: data other than high-priority data may be read without blocking
  - POLLRDNORM: normal data (priority band equals 0) may be read without blocking
  - POLLWRNORM: same as POLLOUT
- Regular TCP/UDP considered normal
- TCP out-of-band is priority band
- TCP read half-close is normal
- TCP error can be either normal or error (POLLERR)

***