



BITS Pilani
Pilani Campus

Network Programming

K Hari Babu
Department of Computer Science & Information Systems



BITS Pilani
Pilani Campus



Outline

Outline



- Advanced I/O
 - `recv()`, `send()`
 - `readv()`, `writew()`
 - `recvmsg()`, `sendmsg()`



Advanced I/O Functions (T1: ch 14)

recv() and send()



- The recv() and send() system calls perform I/O on connected sockets (TCP or connected UDP sockets).
- Socket-Specific I/O System Calls:
 - They provide socket-specific functionality not available with read() and write().

```
1  #include <sys/socket.h>
2  ssize_t recv(int sockfd , void * buffer , size_t length , int flags );
3  //Returns number of bytes received, 0 on EOF, or -1 on error
4  ssize_t send(int sockfd , const void * buffer , size_t length , int flags );
5  //Returns number of bytes sent, or -1 on error
```

- Same as read() and write() except for *flags*.
- Return values are same as read() and write().

- **MSG_DONTWAIT:**
 - perform a non-blocking recv().
 - can be done using fcntl() call but that will make sock fd non-blocking. Here only this operation is non-blocking.
- **MSG_OOB:**
 - receive out-of-band data on the socket.
- **MSG_PEEK:**
 - retrieve a copy of the requested bytes from the socket buffer.
 - Data is not removed from the socket buffer.
 - Used for knowing the no of bytes available on the buffer.
- **MSG_WAITALL**
 - Blocks until *length* bytes are read from socket buffer.
 - May get interrupted by signals.

send() flags



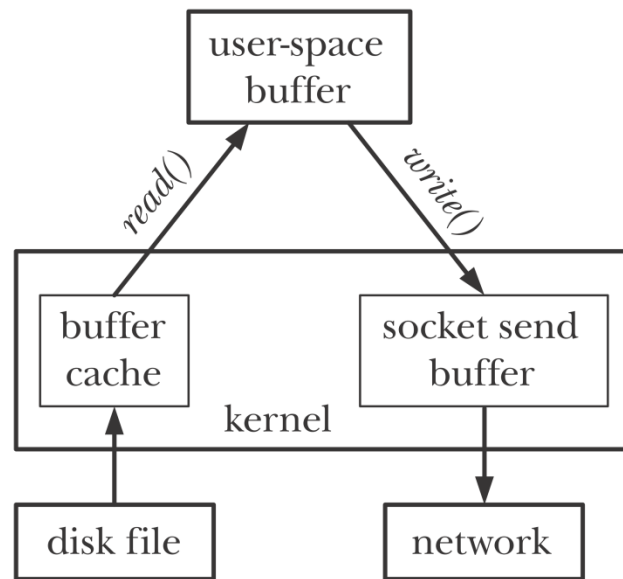
- **MSG_DONTWAIT**
 - Perform a non-blocking send.
- **MSG_MORE**
 - Data written using send() or sendto() calls with this flag is packaged into a single datagram until a send() without this flag.
- **MSG_NOSIGNAL**
 - Do not generate SIGPIPE signal. Return only EPIPE error.
- **MSG_OOB**
 - Write out of band data on TCP.

sendfile() sys call (R1:61.4)

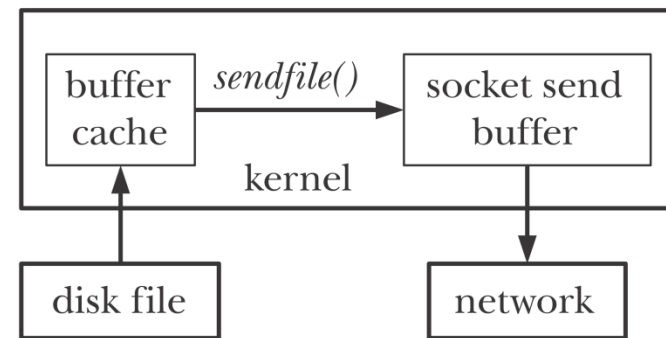


- Transferring large file in web servers requires repeated calls to `read()` and `write()`.
 - This is inefficient.

```
1 while ((n = read(diskfilefd, buf, BUZ_SIZE)) > 0)
2   write(sockfd, buf, n);
```



a) `read()` + `write()`



b) `sendfile()`

sendfile() sys call (R1:61.4)



- The sendfile() sys call is designed to eliminate copying file data into user space.
 - File contents are directly transferred to the socket without going through user space.
 - This is referred as a *zero-copy transfer*.

```
1  #include <sys/sendfile.h>
2  ssize_t sendfile(int out_fd, int in_fd, off_t * offset, size_t count );
3  //Returns number of bytes transferred, or -1 on error
```

- *out_fd*: is the socket fd.
- *In_fd*: is regular file fd.
- *off_t*: is the offset. This is a value-result argument.
- *count* is the number of bytes to be transferred.
- *sendfile* doesn't change the file offset for *in_fd*.

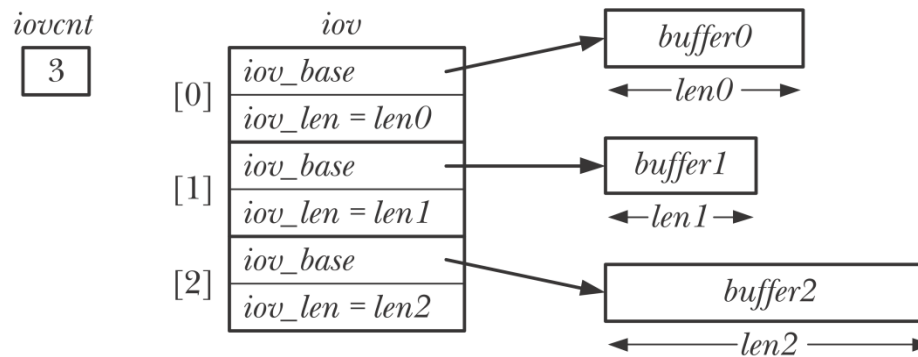
readv() and writev()



- The readv() and writev() system calls perform scatter-gather I/O.
- *iov* points to an array of buffers, each in *iovec* structure.

```
1  #include <sys/uio.h>
2  ssize_t readv(int fd , const struct iovec * iov , int iovcnt );
3  //Returns number of bytes read, 0 on EOF, or -1 on error
4  ssize_t writev(int fd , const struct iovec * iov , int iovcnt );
5  //Returns number of bytes written, or -1 on error
```

```
8  struct iovec {
9      void *iov_base; /* Start address of buffer */
10     size_t iov_len;  /* Number of bytes to transfer to/from buffer */
11 };
```



readv() and writev()



- The readv() system call performs scatter input:
 - Reads from the file and puts the data into the buffer starting at iov[0]. Once the first buffer is full, it goes to another.
- readv() completes atomically:
 - Kernel performs a single data transfer.
 - Assured that all the bytes read are contiguous in the file. File offset can't be changed by other process.
- The writev() call performs gather output:
 - Starting from the first buffer, writes the data contiguously into the file.
 - Partial write is possible.
- writev() completes atomically.
- readv() and writev() are used for convenience and speed.
 - Reduce number of sys calls.

sendmsg() & recvmsg() sys calls



- The sendmsg() and recvmsg() system calls are the most general purpose of the socket I/O system calls.
 - The sendmsg() system call can do everything that is done by write(), send(), and sendto();
 - the recvmsg() system call can do everything that is done by read(), recv(), and recvfrom().

```
1  #include <sys/socket.h>
2  ssize_t recvmsg(int sockfd, struct msghdr *msg, int flags);
3  ssize_t sendmsg(int sockfd, struct msghdr *msg, int flags);
4  //Both return: number of bytes read or written if OK, -1 on error
```

```
1  struct msghdr {
2      void            *msg_name;           /* protocol address */
3      socklen_t       msg_namelen;        /* size of protocol address */
4      struct iovec    *msg_iov;           /* scatter/gather array */
5      int             msg_iovlen;         /* # elements in msg_iov */
6      void            *msg_control;        /* ancillary data (cmsghdr struct) */
7      socklen_t       msg_controllen;     /* length of ancillary data */
8      int             msg_flags;          /* flags returned by recvmsg() */
9  };
```

sendmsg() & recvmsg() sys calls



- Can be used to send or receive ancillary data (control information).
- Flags are

Flag	Examined by: Send flags Sendto flags Sendmsg flags	Examined by: recv flags recvfrom flags recvmsg flags	Returned by: Recvmsg msg_flags
MSG_DONTROUTE MSG_DONTWAIT MSG_PEEK MSG_WAITALL	• •	• • •	
MSG_EOR MSG_OOB	• •	•	• •
MSG_BCAST MSG_MCAST MSG_TRUNC MSG_CTRUNC			• • • •

Flags returned by rcvmsg()

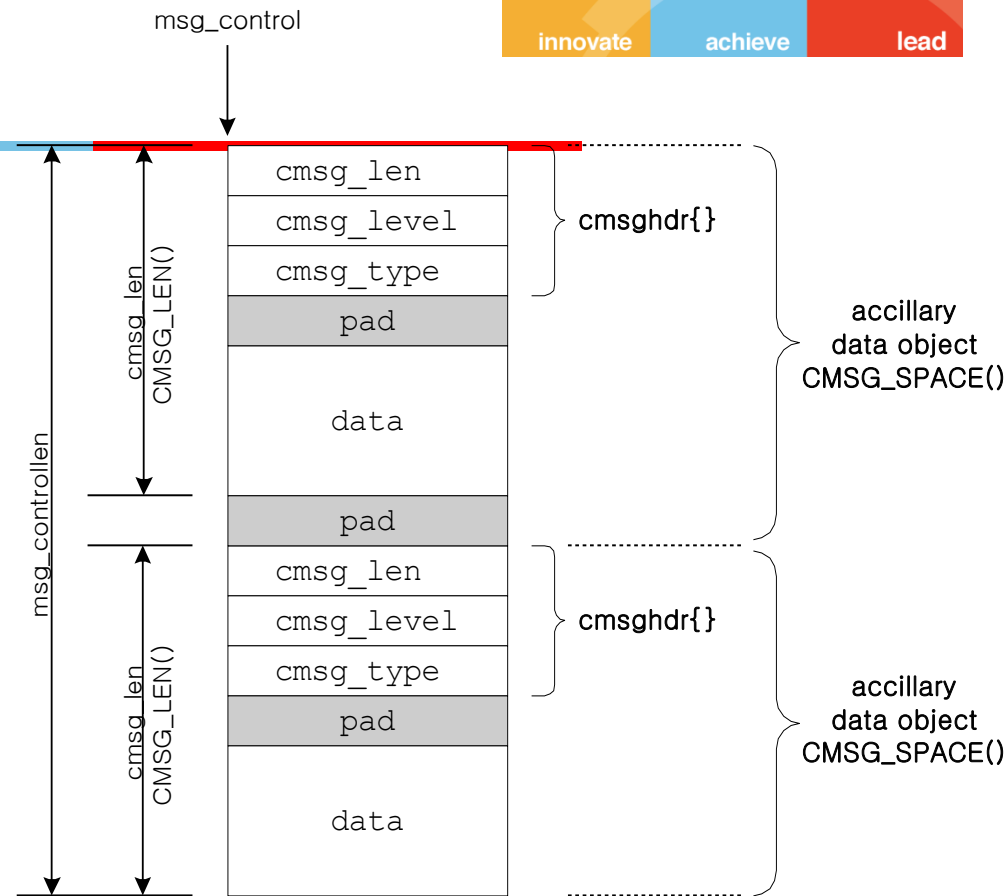


- **MSG_BCAST**
 - is returned if the datagram was received as as a broadcast.
- **MSG_MCAST**
 - is returned if the datagram was received as a link-layer multicast.
- **MSG_TRUNC**
 - is returned if the datagram was truncated
- **MSG_CTRUNC**
 - is returned if the ancillary data was truncated
- **MSG_EOR**
 - is turned on if the returned data ends a logical record.
- **MSG_OOB**
 - This flag is never returned for TCP out-of-band data. This flag is returned by other protocol suites (e.g., the OSI protocols).
- **MSG_NOTIFICATION**
 - This flag is returned for SCTP receivers to indicate that the message read is an event notification, not a data message.

Ancillary Data



- Ancillary data can be sent and received using the `msg_control` and `msg_controllen` members of the `msghdr` structure.
 - Another term for ancillary data is control information.



```
1 struct cmsghdr {  
2     socklen_t  cmsg_len;    /* length in bytes, including this structure */  
3     int        cmsg_level; /* originating protocol */  
4     int        cmsg_type;  /* protocol-specific type */  
5     /* followed by unsigned char cmsg_data[] */  
6 };
```

Ancillary Data



- Ancillary data is domain specific.

Protocol	cmsg_level	Cmsg_type	Description
IPv4	IPPROTO_IP	IP_RECVDSTADDR	receive destination address with UDP datagram
		IP_RECVIF	receive interface index with UDP datagram
IPv6	IPPROTO_IPV6	IPV6_DSTOPTS	specify / receive destination options
		IPV6_HOPLIMIT	specify / receive hop limit
		IPV6_HOPOPTS	specify / receive hop-by-hop options
		IPV6_NEXTHOP	specify next-hop address
		IPV6_PKTINFO	specify / receive packet information
		IPV6_RTHDR	specify / receive routing header
Unix domain	SOL_SOCKET	SCM_RIGHTS	send / receive descriptors
		SCM_CREDS	send / receive user credentials

Ancillary Data



- File descriptors and process credentials can be passed between unrelated processes using ancillary data.

cmsghdr{}	
msg_len	16
msg_level	SOL_SOCKET
msg_type	SCM_RIGHTS
discriptor	

cmsghdr{}	
msg_len	16
msg_level	SOL_SOCKET
msg_type	SCM_CREDS
fcred{}	

How much data is Queued?



- Use *recv()* with MSG_PEEK flag.

```
2 int numbytes = recv(fd, buf, bufsize, MSG_PEEK);
```

- For TCP, this will give the number of bytes available in socket *recv* buffer.
 - This value could change in between two reads.
 - For a connected UDP socket, this return the number of bytes in the next available datagram.
 - Between two reads this value remains same.
- Use *ioctl()* call with FIONREAD command.

```
2 ioctl(fd, FIONREAD, &numbytes)
```

- In UDP case, the size of datagram can be zero. This makes it difficult to distinguish between nodata or data.
 - Safer to use *select()* first and then call I/O.

Sockets and Standard I/O



- TCP and UDP sockets are full duplex. File streams can also be full duplex.
- We can open a file stream on a socket using *fdopen()*.
- If we open with mode r+, then
 - An input function can't be followed by output function without calling *fseek()*.
 - An output function can't be followed by input function without calling *fseek()*.
 - But we can't call *lseek()* on sockets.
- So open two separate streams on a socket: one for reading, one for writing.

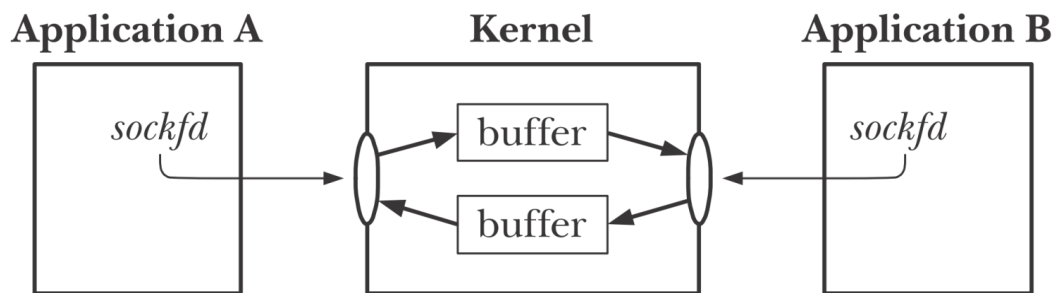


Unix Domain Sockets

Internet Domain vs Unix Domain



- Internet Domain: *AF_INET* or *AF_INET6*
 - Used in network communication.
 - Can be used between two processes in the same host also.
- Unix Domain: *AF_UNIX* or *AF_LOCAL*
 - Used for communication between processes on the same host. Same API as sockets API.
 - No TCP/IP protocol stack. A socket is made of two buffers in the kernel.
 - No header processing, no checksums. Reliable communication. Unix domain sockets are twice faster.



Unix Domain Sockets - Usage



- Unix domain sockets are used for three reasons:
 - Unix domain sockets are often twice as fast as a TCP socket when both peers are on the same host.
 - X Windows
 - used when passing file descriptors between processes on the same host.
 - unix domain sockets provide the client's process credentials (user ID and group IDs) to the server, which can provide additional security checking

Unix Domain Sockets



- Two types of sockets are provided.
 - Stream sockets
 - Similar to TCP
 - Datagram Sockets
 - Similar to UDP sockets.
 - Message boundaries are preserved.
 - Communication is reliable unlike UDP.

Unix Domain Socket Address



- End Point Address
 - pathnames within the normal file system
 - The pathname associated with a Unix domain socket should be an absolute pathname.

```
1 struct sockaddr_un {  
2     sa_family_t sun_family; /* Always AF_UNIX */  
3     char sun_path[108]; /* Null-terminated socket pathname */  
4 };
```


Binding End Point to a Socket



- When used to bind a UNIX domain socket, bind() creates an entry in the file system.

```
1  const char *SOCKNAME = "/tmp/mysock";
2  int sfd;
3  struct sockaddr_un addr;
4  sfd = socket(AF_UNIX, SOCK_STREAM, 0);    /* Create socket */
5  if (sfd == -1)
6      errExit("socket");
7  memset(&addr, 0, sizeof(struct sockaddr_un)); /* Clear structure */
8  addr.sun_family = AF_UNIX;    /* UNIX domain address */
9  strncpy(addr.sun_path, SOCKNAME, sizeof(addr.sun_path) - 1);
10 if (bind(sfd, (struct sockaddr *) &addr, sizeof(struct sockaddr_un)) == -1)
11     errExit("bind");
```

- We can't bind a socket to an existing pathname (bind() fails with the error EADDRINUSE).
- A socket may be bound to only one pathname; conversely, a pathname can be bound to only one socket.
 - When the socket is no longer required, its pathname entry should be removed using unlink().

Unix Domain Stream Sockets



- **Server**
 - Absolute pathname is required.
 - Pathname specified in connect() should be existing, and bound to a socket
 - If the listening socket's queue is full, ECONNREFUSED is immediately returned.
- **Client**
 - Create socket
 - Connect to the server.

Unix Domain Stream Server



```
1  /*sockets/us_xfr_sv.c*/
2  main(int argc, char *argv[])
3  {   struct sockaddr_un addr;
4      int sfd, cfd;
5      ssize_t numRead;
6      char buf[BUF_SIZE];
7      sfd = socket(AF_UNIX, SOCK_STREAM, 0);
8      if (sfd == -1)  errExit("socket");
9      if (remove(SV_SOCK_PATH) == -1 && errno != ENOENT)
10         errExit("remove-%s", SV_SOCK_PATH);
11     memset(&addr, 0, sizeof(struct sockaddr_un));
12     addr.sun_family = AF_UNIX;
13     strncpy(addr.sun_path, SV_SOCK_PATH, sizeof(addr.sun_path) - 1);
14     if (bind(sfd, (struct sockaddr *) &addr, sizeof(struct sockaddr_un)) == -1)
15         errExit("bind");
16     if (listen(sfd, BACKLOG) == -1)  errExit("listen");
17     for (;;) {
18         cfd = accept(sfd, NULL, NULL);
19         if (cfd == -1)  errExit("accept");
20         while ((numRead = read(cfd, buf, BUF_SIZE)) > 0)
21             if (write(STDOUT_FILENO, buf, numRead) != numRead)
22                 fatal("partial/failed write");
23         if (numRead == -1) errExit("read");
24         if (close(cfd) == -1) errMsg("close");
25     }}
```

Unix Domain Stream Client



```
1  /*sockets/us_xfr_cl.c*/
2  main(int argc, char *argv[])
3  {
4      struct sockaddr_un addr;
5      int sfd;
6      ssize_t numRead;
7      char buf[BUF_SIZE];
8      sfd = socket(AF_UNIX, SOCK_STREAM, 0); /* Create client socket */
9      if (sfd == -1) errExit("socket");
10     /* Construct server address, and make the connection */
11     memset(&addr, 0, sizeof(struct sockaddr_un));
12     addr.sun_family = AF_UNIX;
13     strncpy(addr.sun_path, SV_SOCKET_PATH, sizeof(addr.sun_path) - 1);
14     if (connect(sfd, (struct sockaddr *) &addr,
15               sizeof(struct sockaddr_un)) == -1)
16         errExit("connect");
17     /* Copy stdin to socket */
18     while ((numRead = read(STDIN_FILENO, buf, BUF_SIZE)) > 0)
19         if (write(sfd, buf, numRead) != numRead)
20             fatal("partial/failed write");
21     if (numRead == -1)
22         errExit("read");
23     exit(EXIT_SUCCESS); /* Closes our socket; server sees EOF */
24 }
```

Unix Domain Datagram Sockets



- Datagram sockets are reliable unlike UDP sockets.
 - Datagrams are not lost.
 - Datagrams are delivered in order and without duplicates.
- Server
 - Creates a socket
 - binds to well-known path.
- Client
 - Creates a socket
 - **binds the socket to an address, so that the server can send its reply.**
 - The client address is made unique by including the client's process ID in the pathname.

Unix Domain Datagram Server



```
1  /*sockets/ud_ucose_sv.c*/
2  main(int argc, char *argv[])
3  {
4      struct sockaddr_un svaddr, claddr;
5      sfd = socket(AF_UNIX, SOCK_DGRAM, 0);          /* Create server socket */
6      if (remove(SV_SOCK_PATH) == -1 && errno != ENOENT)
7          errExit("remove-%s", SV_SOCK_PATH);
8      memset(&svaddr, 0, sizeof(struct sockaddr_un));
9      svaddr.sun_family = AF_UNIX;
10     strncpy(svaddr.sun_path, SV_SOCK_PATH, sizeof(svaddr.sun_path) - 1);
11     if (bind(sfd, (struct sockaddr *) &svaddr, sizeof(struct sockaddr_un)) == -1)
12         errExit("bind");
13     for (;;) {
14         len = sizeof(struct sockaddr_un);
15         numBytes = recvfrom(sfd, buf, BUF_SIZE, 0,
16                             (struct sockaddr *) &claddr, &len);
17         if (numBytes == -1) errExit("recvfrom");
18         printf("Server received %ld bytes from %s\n", (long) numBytes,
19               claddr.sun_path);
20         for (j = 0; j < numBytes; j++)
21             buf[j] = toupper((unsigned char) buf[j]);
22         if (sendto(sfd, buf, numBytes, 0, (struct sockaddr *) &claddr, len) !=
23             numBytes)
24             fatal("sendto");
25     }
26 }
```

Unix Domain Datagram Client



```
1 ▾ /* sockets/ud_ucase_cl.c */
2  main(int argc, char *argv[])
3  {   struct sockaddr_un svaddr, claddr;
4      sfd = socket(AF_UNIX, SOCK_DGRAM, 0);
5      memset(&claddr, 0, sizeof(struct sockaddr_un));
6      claddr.sun_family = AF_UNIX;
7      snprintf(claddr.sun_path, sizeof(claddr.sun_path),
8               "/tmp/ud_ucase_cl.%ld", (long) getpid());
9  if (bind(sfd, (struct sockaddr *) &claddr, sizeof(struct sockaddr_un)) == -1)
10     errExit("bind");
11 ▾  /* Construct address of server */
12  memset(&svaddr, 0, sizeof(struct sockaddr_un));
13  svaddr.sun_family = AF_UNIX;
14  strncpy(svaddr.sun_path, SV_SOCKET_PATH, sizeof(svaddr.sun_path) - 1);
15 ▾  /* Send messages to server; echo responses on stdout */
16 ▾  for (j = 1; j < argc; j++) {
17      msgLen = strlen(argv[j]);          /* May be longer than BUF_SIZE */
18      if (sendto(sfd, argv[j], msgLen, 0, (struct sockaddr *) &svaddr,
19                sizeof(struct sockaddr_un)) != msgLen)
20          fatal("sendto");
21      numBytes = recvfrom(sfd, resp, BUF_SIZE, 0, NULL, NULL);
22      if (numBytes == -1) errExit("recvfrom");
23      printf("Response %d: %.*s\n", j, (int) numBytes, resp); }
24  remove(claddr.sun_path); /* Remove client socket pathname */
25 }
```

socketpair()

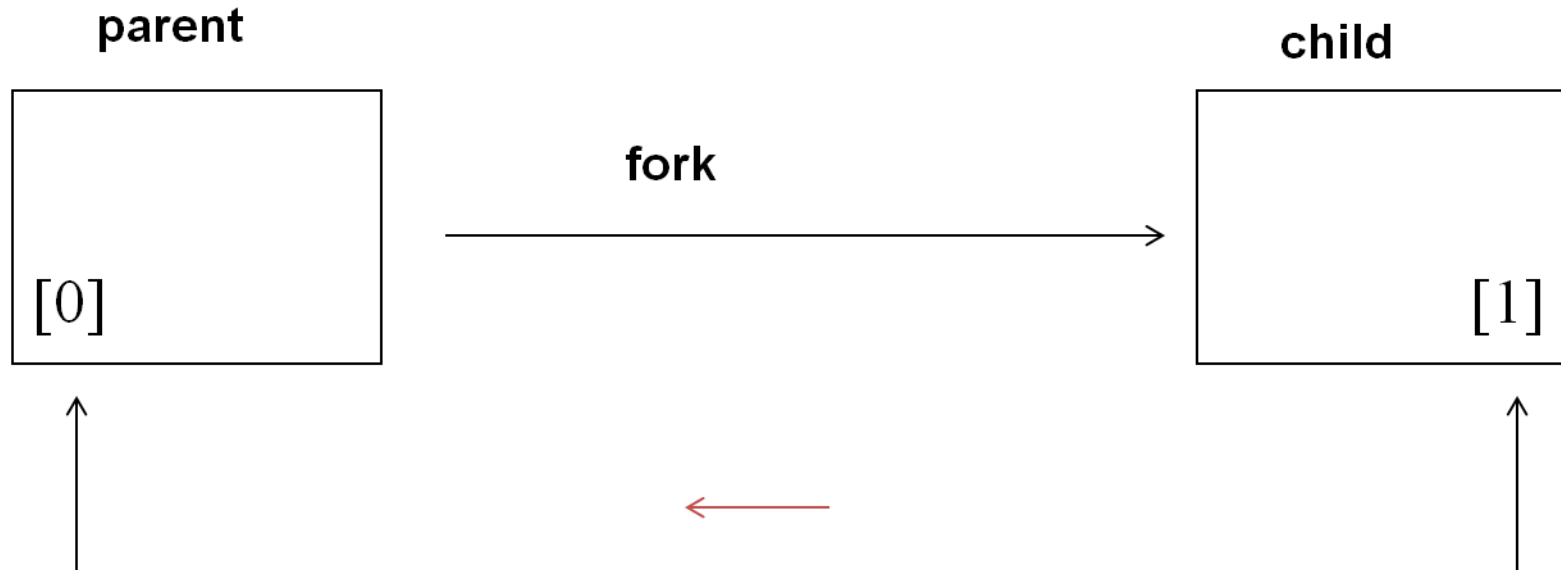


- creates a pair of sockets and connect them together.

```
1  #include <sys/socket.h>
2  int socketpair(int domain, int type ,int protocol ,int sockfd [2]);
3  //Returns 0 on success, or -1 on error
```

- Returns two socket fds.
- No path names bound for sockets. Not visible outside the process.
- Type SOCK_STREAM creates the equivalent of a bidirectional pipe (also known as a stream pipe).
- Each socket can be used for both reading and writing.
- Just like in pipe, after creating calling socketpair(), fork() is called.
- Parent and child can communicate using sockets.

socketpair()



socketpair()



```
1  #include <sys/types.h>
2  #include <unistd.h>
3  #include <sys/socket.h>
4  main ()
5  {
6      int i;
7      int p[2];
8      pid_t ret;
9      socketpair(AF_UNIX, SOCK_STREAM, 0, p);
10     ret = fork ();
11     if (ret == 0)
12     {
13         close (1);
14         dup (p[1]);
15         close (p[0]);
16         execlp ("ls", "ls", "-l", (char *) 0);
17     }
18     if (ret > 0)
19     {
20         close (0);
21         dup (p[0]);
22         close (p[1]);
23         execlp ("wc", "wc", "-l", (char *) 0);
24     }
25 }
```

Passing File Descriptors



- Unix system provide a way to pass any open descriptor from one process to any other process.(using *sendmsg()*)
- It allows one process (typically a server) to do the privileged execution
 - dialing a modem, negotiating locks for the file or deal with database
 - simply pass back to the calling process a descriptor that can be used with all the I/O functions.
- All the details involved in opening the file or device are hidden from the client.

Passing File Descriptors

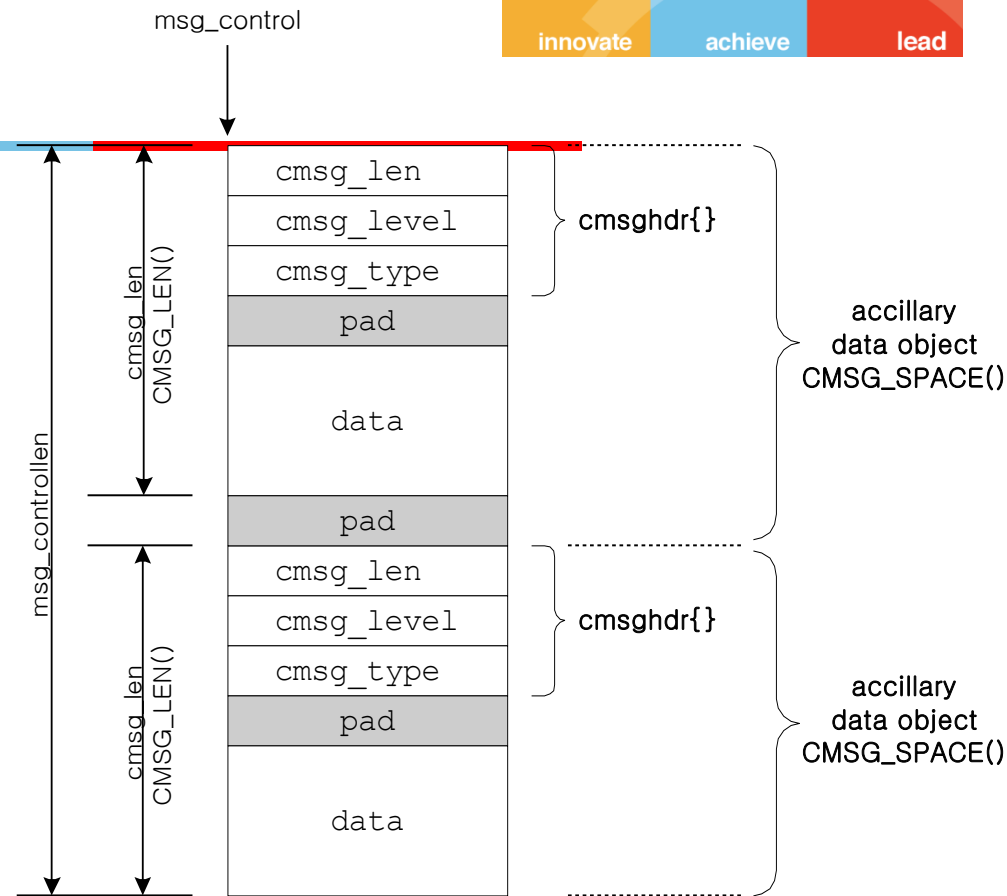


- Steps involved
 - Create a unix domain socket(stream or datagram)
 - one process opens a descriptor by calling any of the unix function that returns a descriptor
 - the sending process build a ***msghdr*** structure containing the descriptor to be passed
 - Sending process sends ancillary data using ***sendmsg()*** with SCM_RIGHTS
 - the receiving process calls ***recvmsg()*** to receive the descriptor on the unix domain socket
- Passing a descriptor is not same as passing descriptor number
 - involves creating a new descriptor in the receiving process that refers to the same file table entry within the kernel.

Ancillary Data



- Ancillary data can be sent and received using the `msg_control` and `msg_controllen` members of the `msghdr` structure.
 - Another term for ancillary data is control information.



```
1 struct cmsghdr {  
2     socklen_t  cmsg_len;    /* length in bytes, including this structure */  
3     int        cmsg_level; /* originating protocol */  
4     int        cmsg_type;  /* protocol-specific type */  
5     /* followed by unsigned char cmsg_data[] */  
6 };
```

Ancillary Data



- Ancillary data is domain specific.

Protocol	cmsg_level	Cmsg_type	Description
IPv4	IPPROTO_IP	IP_RECVDSTADDR IP_RECVIF	receive destination address with UDP datagram receive interface index with UDP datagram
IPv6	IPPROTO_IPV6	IPV6_DSTOPTS IPV6_HOPLIMIT IPV6_HOPOPTS IPV6_NEXTHOP IPV6_PKTINFO IPV6_RTHDR	specify / receive destination options specify / receive hop limit specify / receive hop-by-hop options specify next-hop address specify / receive packet information specify / receive routing header
Unix domain	SOL_SOCKET	SCM_RIGHTS SCM_CREDS	send / receive descriptors send / receive user credentials

Ancillary Data



- File descriptors and process credentials can be passed between unrelated processes using ancillary data.

cmsghdr{}	
msg_len	16
msg_level	SOL_SOCKET
msg_type	SCM_RIGHTS
discriptor	

cmsghdr{}	
msg_len	16
msg_level	SOL_SOCKET
msg_type	SCM_CREDS
fcred{}	

Passing File Descriptors Example



- Protocol
 - If `buf[1]<0` then there is error.
 - If `buf[1]=0` then it is success.
- Receiver
 - Create unix domain stream socket
 - Send file name
 - Call `recv_fd`
- Sender
 - Create unix domain stream socket
 - Open the file descriptor
 - Call `send_fd`

- Macros associated with ancillary data

```
1  #include <sys/socket.h>
2  #include <sys/param.h> /* for ALIGN macro on many implementations */
3  struct cmsghdr *CMSG_FIRSTHDR(struct msghdr *mhdrptr) ;
4  //Returns: pointer to first cmsghdr structure or NULL if no ancillary data
5  struct cmsghdr *CMSG_NXTHDR(struct msghdr *mhdrptr, struct cmsghdr *cmsgptr) ;
6  //Returns: pointer to next cmsghdr structure or NULL if no more ancillary data
7  unsigned char *CMSG_DATA(struct cmsghdr *cmsgptr) ;
8  //Returns: pointer to first byte of data associated with cmsghdr structure
9  unsigned int CMSG_LEN(unsigned int length) ;
10 //Returns: value to store in cmsg_len given the amount of data
11 unsigned int CMSG_SPACE(unsigned int length) ;
12 //Returns: total size of an ancillary data object given the amount of data
```

Send fd

innovate

achieve

lead

```
1  #include <sys/socket.h>
2  #define CONTROLLEN  CMSG_LEN(sizeof(int))
3  static struct cmsghdr  *cmptr = NULL;
4  int send_fd(int fd, int fd_to_send)
5  {
6      struct iovec      iov[1];
7      struct msghdr      msg;
8      char              buf[2]; /* send_fd()/recv_fd() 2-byte protocol */
9      iov[0].iov_base = buf; iov[0].iov_len = 2;
10     msg.msg_iov      = iov; msg.msg_iovlen = 1;
11     msg.msg_name      = NULL; msg.msg_namelen = 0;
12     if (fd_to_send < 0) {
13         msg.msg_control      = NULL;
14         msg.msg_controllen = 0;
15         buf[1] = -fd_to_send; /* nonzero status means error */
16     }
17     else { if (cmptr == NULL && (cmptr = malloc(CONTROLLEN)) == NULL)
18         return(-1);
19         cmptr->cmsg_level = SOL_SOCKET;
20         cmptr->cmsg_type = SCM_RIGHTS;
21         cmptr->cmsg_len = CONTROLLEN;
22         msg.msg_control = cmptr;
23         msg.msg_controllen = CONTROLLEN;
24         *(int *)CMSG_DATA(cmptr) = fd_to_send; /* the fd to pass */
25         buf[1] = 0; /* zero status means OK */
26     }
27     buf[0] = 0; /* null byte flag to recv_fd() */
28     if (sendmsg(fd, &msg, 0) != 2)
29         return(-1);
30     return(0);
31 }
```

Receive fd

innovate

achieve

lead

```
1 int recv_fd(int sockfd )
2 {
3     #define CONTROLLEN  CMSG_LEN(sizeof(int))
4     static struct cmsghdr  *cmptr = NULL;
5
6     struct iovec    iov[1];
7     struct msghdr    msg;
8     char            buf[2]; /* send_fd()/recv_fd() 2-byte protocol */
9     memset(&msg, 0, sizeof(msg));
10    iov.iov_base     = buf;
11    iov.iov_len       = sizeof(data)-1;
12    msg.msg_iov       = &iov;
13    msg.msg_iovlen    = 1;
14    if (cmptr == NULL && (cmptr = malloc(CONTROLLEN)) == NULL)
15        return(-1);
16    msg.msg_control    = cmptr;
17    msg.msg_controllen = CONTROLLEN;
18    recvmsg(sockfd, &msg, 0)
19    if (buf[1]<0)) {
20        printf("failed to open %s: %s\n", name, data);
21        return -1;
22    }
23    /* Loop over all control messages */
24    cmsg = CMSG_FIRSTHDR(&msg);
25    while (cmsg != NULL) {
26        if (cmsg->cmsg_level == SOL_SOCKET
27            && cmsg->cmsg_type == SCM_RIGHTS)
28            return *(int *) CMSG_DATA(cmsg);
29        cmsg = CMSG_NXTHDR(&msg, cmsg);
30    }
31 }
```

Passing Credentials



- A process can pass its credentials as ancillary data using SCM_CREDENTIALS option.
- The structure for credentials

```
1 struct ucred {  
2     pid_t pid;      /* process ID of the sending process */  
3     uid_t uid;      /* user ID of the sending process */  
4     gid_t gid;      /* group ID of the sending process */  
5 };
```

- This structure is filled by the kernel and passed onto the receiver process.
- Example:
 - Sender process sends file name and access mode.
 - Server verifies the credentials and passes on the fd.

```
1 sockfd = socket(AF_LOCAL, SOCK_STREAM, 0);
2 bzero(&servaddr, sizeof(servaddr));
3 servaddr.sun_family = AF_LOCAL;
4 strcpy(servaddr.sun_path, "PATH");
5 connect(sockfd, (struct sockaddr *) &servaddr, sizeof(servaddr));
6 msgh.msg_iov = &iov;
7 msgh.msg_iovlen = 1;
8 /* Send Filename and Access Mode to server */
9 strcat(data, argv[1]);strcat(data, "#");
10 strcat(data, argv[2]);strcat(data, "#");
11 iov.iov_base = data;
12 iov.iov_len = MAX_DATA;
13 msgh.msg_name = NULL;
14 msgh.msg_namelen = 0;
15 msgh.msg_control = NULL;
16 msgh.msg_controllen = 0;
17
18 if(sendmsg(sockfd, &msggh, 0) < 0)
19 {
20     perror("Error sending message");
21     exit(1);
22 }
```

Receiver

innovate

achieve

lead

```
1  optval = 1;
2  /* Set SO_PASSCRED socket option for receiving credentials of other processes */
3  setsockopt(*(int *)arg, SOL_SOCKET, SO_PASSCRED, &optval, sizeof(optval));
4  /* Set 'control_un' to describe ancillary data that we want to receive */
5  control_un.cmh.cmsg_len = CMSG_LEN(sizeof(struct ucred));
6  control_un.cmh.cmsg_level = SOL_SOCKET;
7  control_un.cmh.cmsg_type = SCM_CREDENTIALS;
8  /* Set 'msg' fields to describe 'control_un' */
9  msg.cmsg_control = control_un.control;
10 msg.cmsg_controllen = sizeof(control_un.control);
11 msg.cmsg_iov = &iov;    msg.cmsg_iovlen = 1;
12 iov.iov_base = data;
13 iov.iov_len = MAX_DATA;
14 msg.cmsg_name = NULL;
15 msg.cmsg_namelen = 0;
16 /* Receive real plus ancillary data */
17 nr = recvmsg(*(int *)arg, &msg, 0);
18 /* Extract credentials information from received ancillary data */
19 cmhp = CMSG_FIRSTHDR(&msg);
20 ucredp = (struct ucred *) CMSG_DATA(cmhp);
21 printf("Received Credentials pid: %ld, uid: %ld, gid: %ld\n",
22 (long) ucredp->pid, (long) ucredp->uid, (long) ucredp->gid);
```

The Linux Abstract Socket Namespace



- Is a Linux-specific feature that allows us to bind a UNIX domain socket to a name without that name being created in the file system.
- It is not necessary to unlink the socket pathname when we have finished using the socket. The abstract name is automatically removed when the socket is closed.
 - To create an abstract binding, we specify the first byte of the `sun_path` field as a null byte (`\0`).
 - The remaining bytes of the `sun_path` field then define the abstract name for the socket. These bytes are interpreted in their entirety, rather than as a null-terminated string.

The Linux Abstract Socket Namespace Example



```
1 struct sockaddr_un addr;
2 memset(&addr, 0, sizeof(struct sockaddr_un)); /* Clear address structure */
3 addr.sun_family = AF_UNIX; /* UNIX domain address */
4 /* addr.sun_path[0] has already been set to 0 by memset() */
5 strncpy(&addr.sun_path[1], "xyz", sizeof(addr.sun_path) - 2);
6 /* Abstract name is "xyz" followed by null bytes */
7 sockfd = socket(AF_UNIX, SOCK_STREAM, 0);
8 if (sockfd == -1)
9     errExit("socket");
10 if (bind(sockfd, (struct sockaddr *) &addr,
11         sizeof(struct sockaddr_un)) == -1)
12     errExit("bind");
```




Unix I/O Models

- While doing I/O there are two phases
 - Waiting for the data
 - Copying the data
- Each I/O model differs how it deals with these two phases.
- There are five I/O models
 - blocking I/O
 - nonblocking I/O
 - I/O multiplexing (select and poll)
 - signal driven I/O (SIGIO)
 - asynchronous I/O (the POSIX aio_functions)

Blocking I/O Model

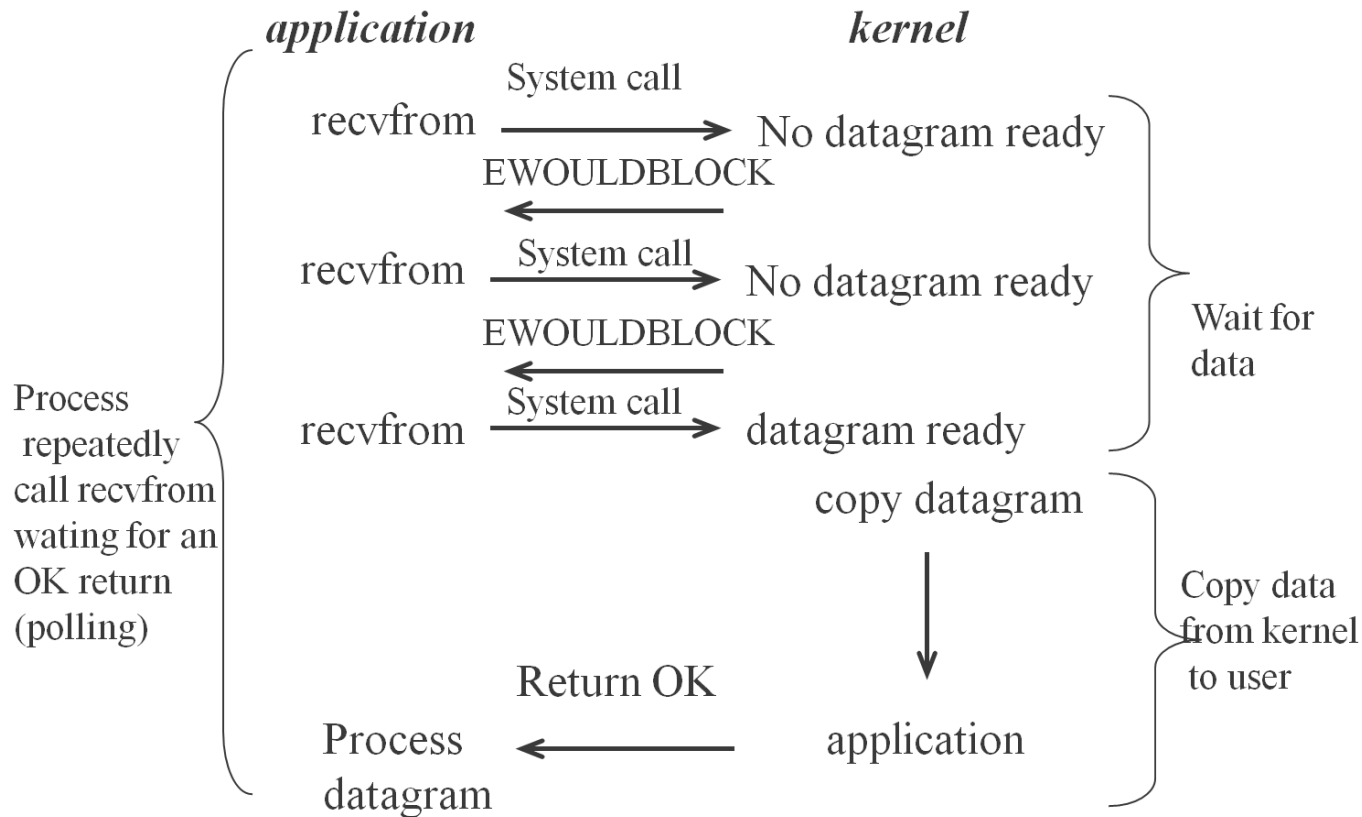


- Most prevalent model

Nonblocking I/O Model



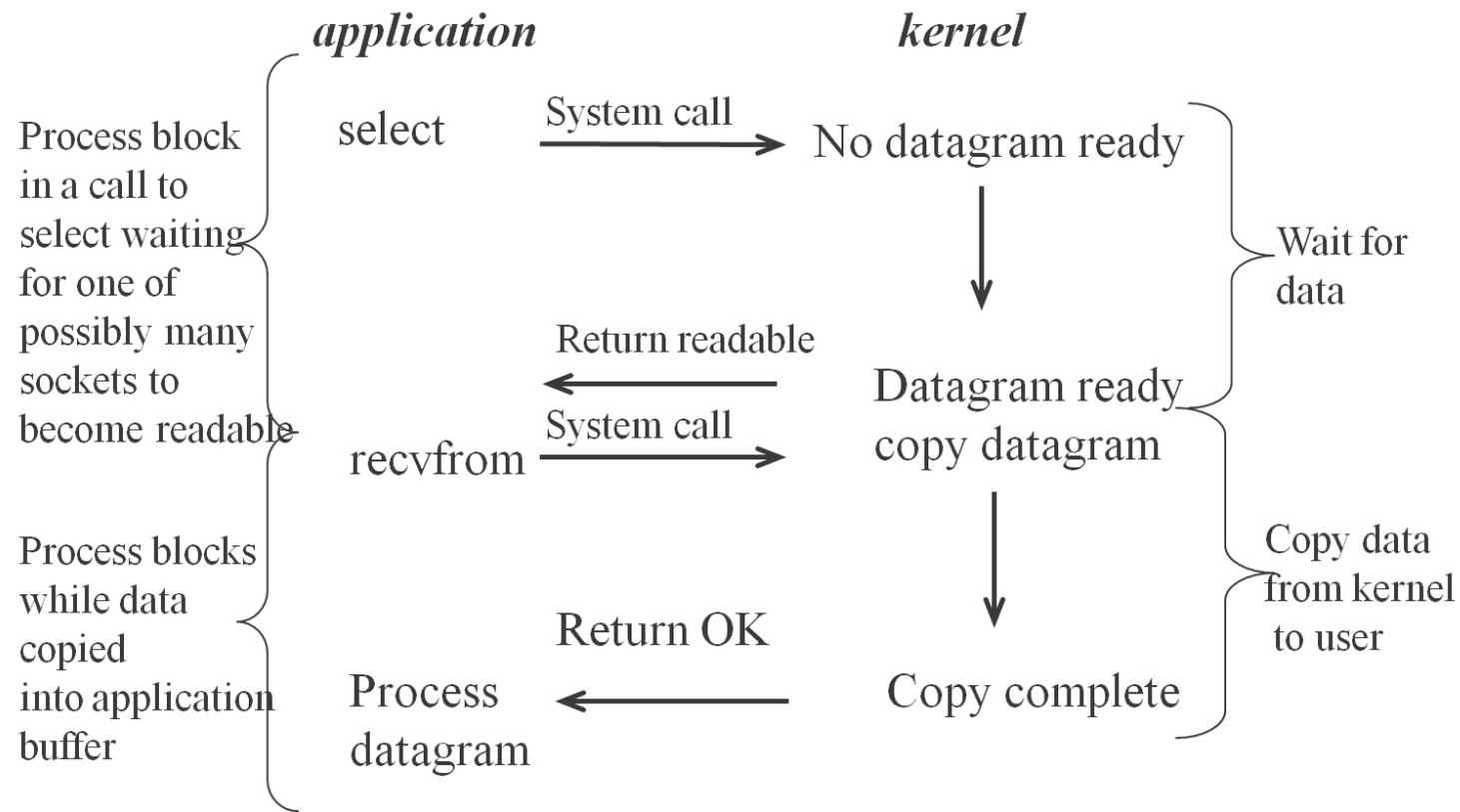
- When the socket is set to be non-blocking,
 - We tell the kernel that do not put the process to sleep if IO can't be completed.



I/O Multiplexing Model



- Block in select() or poll() instead of blocking in actual I/O system call.



I/O Multiplexing Model

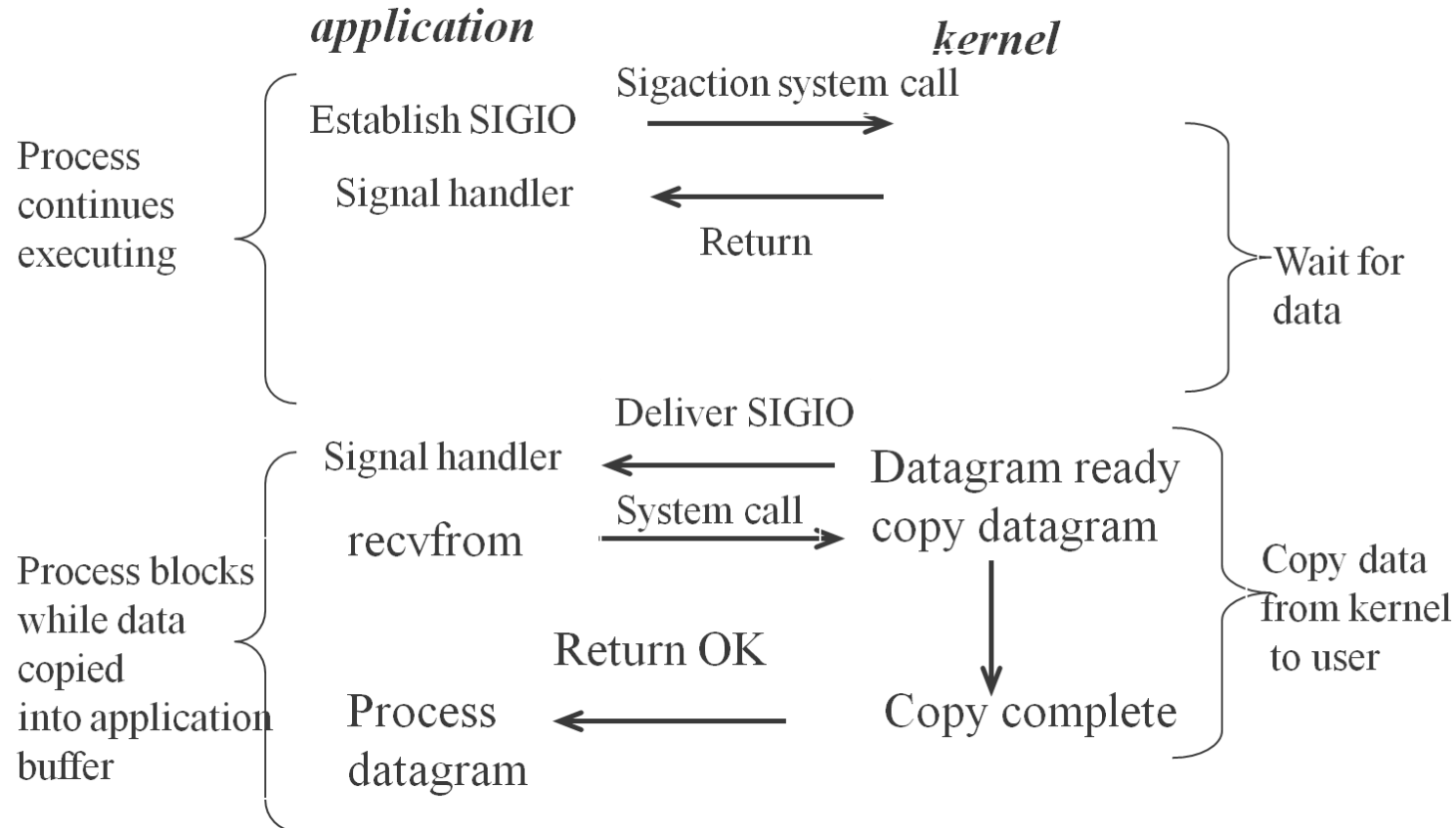


- Blocking and I/O multiplexing seem to be non-different and actually calling two sys calls in I/O multiplexing.
- Advantage with I/O multiplexing is that it can wait for I/O on multiple fds.

Signal-Driven I/O Model



- Tell the kernel to notify us with the SIGIO signal when the descriptor is ready.



Signal-Driven I/O Model

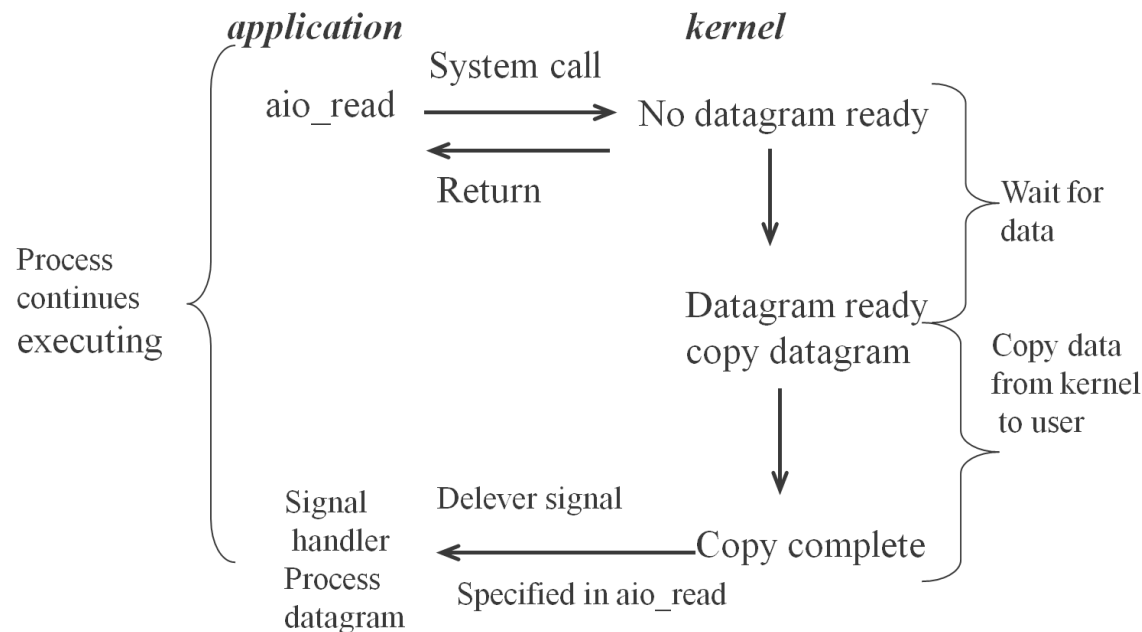


- To use signal-driven I/O with a socket (SIGIO) requires the process to perform the following three steps:
 - A signal handler must be established for the SIGIO signal.
 - The socket owner must be set, normally with the F_SETOWN command of fcntl.
 - Signal-driven I/O must be enabled for the socket, normally with the F_SETFL command of fcntl to turn on the O_ASYNC flag.

Asynchronous I/O Model



- The main difference between this model and the signal-driven I/O models that
 - with signal-driven I/O, the kernel tells us when an I/O operation can be initiated,
 - but with asynchronous I/O, the kernel tells us when an I/O operation is complete.



Asynchronous I/O Model



- POSIX API for asynchronous IO is implemented in a very few systems.
- *aio*cb structure

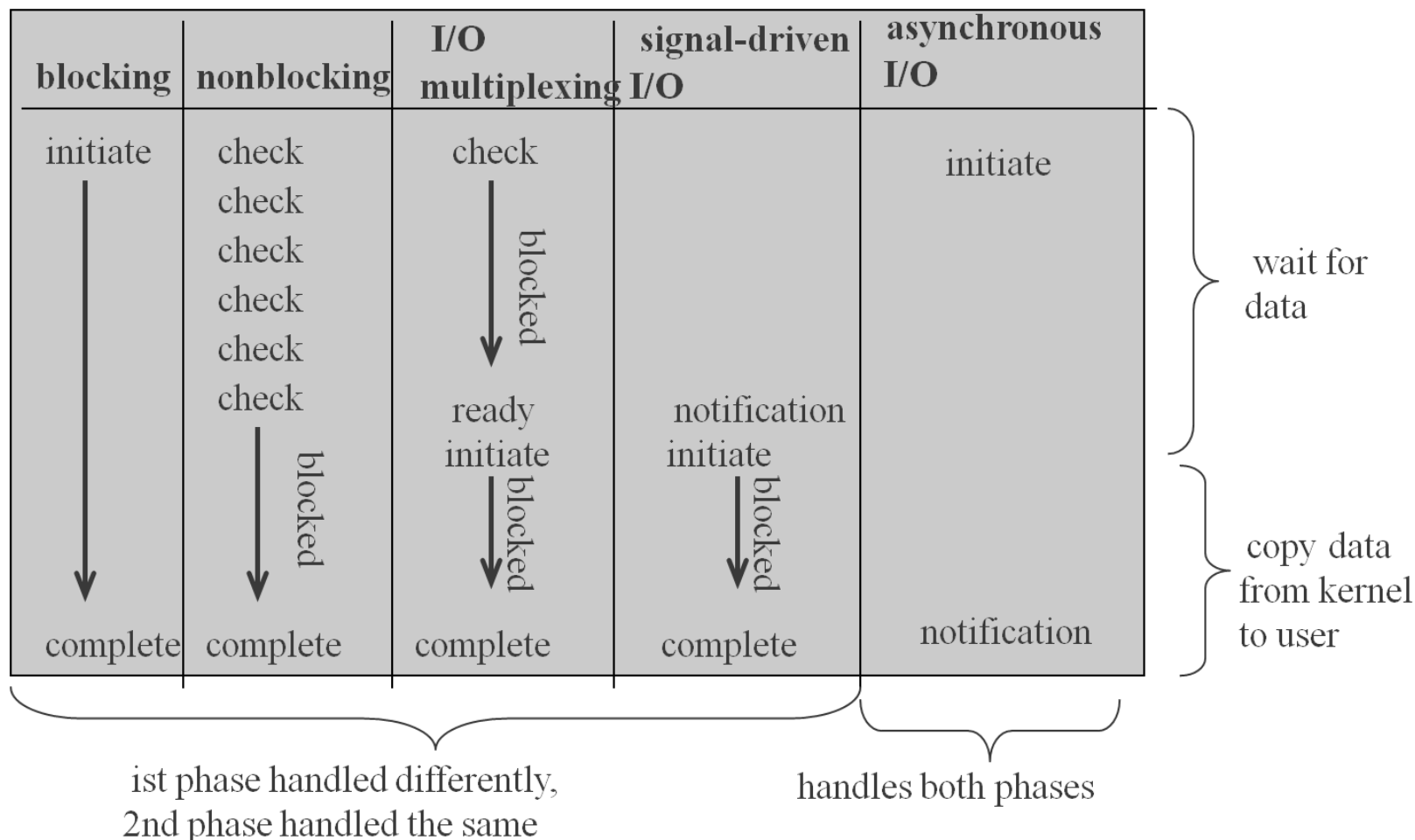
```
1 struct aioctx{
2     int          aio_fdses //    file descriptor
3     off_t        aio_offset //  file offset
4     volatile void* aio_buf //    location of buffer
5     size_t       aio_nbytes //   length of transfer
6     int          aio_reqprio //  request priority offset
7     struct sigevent aio_sigevent // signal number and value
8     int          aio_lio_opcode //operation to be performed
9 };
```

```
1 #include <aio.h>
2 int aio_read(struct aioctx *aioctx);
```

Comparison of I/O Models



- First four models differ only in first phase.



Synchronous I/O vs Asynchronous I/O



- POSIX defines these two terms as follows:
 - A synchronous I/O operation causes the requesting process to be blocked until that I/O operation completes.
 - An asynchronous I/O operation does not cause the requesting process to be blocked.
- Using these definitions,
 - the first four I/O models
 - blocking,
 - nonblocking,
 - I/O multiplexing,
 - and signal-driven I/O
 - are all synchronous because the actual I/O operation (recvfrom) blocks the process.
 - Only the asynchronous I/O model matches the asynchronous I/O definition.



I/O Multiplexing

T1: Ch6

I/O Multiplexing



- I/O multiplexing allows us to simultaneously monitor multiple file descriptors to see if I/O is possible on any of them.
- `select()`, appeared along with the sockets API in BSD. This was historically the more widespread of the two system calls. The other system call, `poll()`, appeared in System V.
- We can use `select()` and `poll()` to monitor file descriptors for regular files, terminals, pseudoterminals, pipes, FIFOs, sockets, and some types of character devices.
- Both system calls allow a process either to block indefinitely waiting for file descriptors to become ready or to specify a timeout on the call.

select()



- The select() system call blocks until one or more of a set of file descriptors becomes ready.

```
1  #include <sys/time.h>    /* For portability */
2  #include <sys/select.h>
3  int select(int  nfds , fd_set * readfds , fd_set * writefds,
4             fd_set * exceptfds, struct timeval * timeout );
5  //Returns number of ready file descriptors, 0 on timeout, or -1 on error
```

- *nfds*: highest number assigned to a descriptor +1.
- *readfds*: set of descriptors we want to read from.
- *writefds*: set of descriptors we want to write to.
- *exceptfds*: set of descriptors to watch for exceptions.
- *timeout*: maximum time select should wait

```
7  struct timeval {
8      long tv_usec;    /* seconds */
9      long tv_usec;    /* microseconds */
10 }
```

select()



- `timeval==NULL`
 - Wait forever : return only when descriptor is ready
- `timeval != NULL`: wait up to a fixed amount of time
 - `timeval = 0`
 - Do not wait at all : return immediately after checking the descriptors
 - `Timeval>0`
 - Return only if descriptor is ready or timeval expires.

File descriptor sets



- The `readfds`, `writfds`, and `exceptfds` arguments are pointers to file descriptor sets, represented using the data type `fd_set`.
- the `fd_set` data type is implemented as a bit mask.

```
1  #include <sys/select.h>
2  void FD_ZERO(fd_set *fdset);
3  /* clear all bits in fdset */
4  void FD_SET(int fd, fd_set *fdset);
5  /* turn on the bit for fd in fdset */
6  void FD_CLR(int fd, fd_set *fdset);
7  /* turn off the bit for fd in fdset */
8  int FD_ISSET(int fd, fd_set *fdset);
9  /* is the bit for fd on in fdset ? */
10 //Returns true (1) if fd is in fdset, or false (0) otherwise
```

- A file descriptor set has a maximum size, defined by the constant `FD_SETSIZE`. On Linux, this constant has the value 1024.

select()



- *nfds*
 - Its value is the maximum descriptor to be tested, plus one
 - example: fds 1,2,5 => nfds: 6
- *readset*
 - descriptor set for checking readable
- *writeset*
 - descriptor set for checking writable
- *exceptset*
 - descriptor set for checking two exception conditions
 - arrival of out of band data for a socket
 - the presence of control status information to be read from the master side of a pseudo terminal
- When select returns value > 1, these sets have been modified by kernel. Now they contain the fds which are ready.

When is the descriptor ready for reading?



- The number of bytes of data in the socket receive buffer is greater than or equal to the current size of the low-water mark for the socket receive buffer. `SO_RCVLOWAT` socket option. It defaults to 1 for TCP and UDP sockets
- The read half of the connection is closed (i.e., a TCP connection that has received a FIN)
- The socket is a listening socket and the number of completed connections is nonzero.
- A socket error is pending. A read operation on the socket will not block and will return an error (`-1`) with `errno` set to the specific error condition.

When the socket is ready for writing?

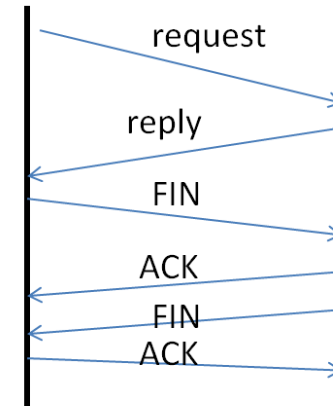


- The number of bytes of available space in the socket send buffer is greater than or equal to the current size of the low-water mark for the socket send buffer. 2048 bytes.
- The write half of the connection is closed. A write operation on the socket will generate SIGPIPE.
- A socket using a non-blocking connect has completed the connection, or the connect has failed
- A socket error is pending. A write operation on the socket will not block and will return an error (−1) with errno set to the specific error condition.
- These pending errors can also be fetched and cleared by calling getsockopt with the SO_ERROR socket option.

Client Handling Multiple Descriptors



- A client is handling two descriptors.
 - *stdin*
 - *socket*
- Sequential handling:
 - First wait on *stdin*.
 - Write to *socket*
 - Read from *socket*.
 - Write to *stdout*.



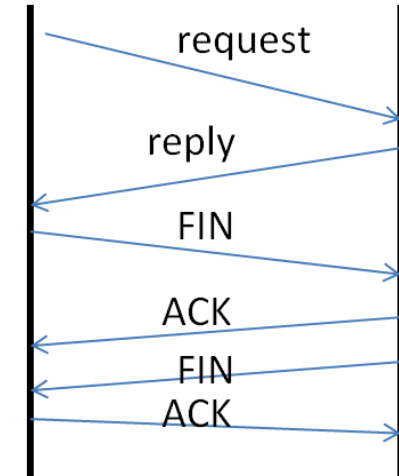
Normal course of actions

```
1 void str_cli(FILE *fp, int sockfd)
2 {
3     char sendline[MAXLINE], recvline[MAXLINE];
4     while (Fgets(sendline, MAXLINE, fp) != NULL) {
5         Writen(sockfd, sendline, strlen(sendline));
6         if (Readline(sockfd, recvline, MAXLINE) == 0)
7             err_quit("str_cli: server terminated prematurely");
8         Fputs(recvline, stdout);
9     }
10 }
```

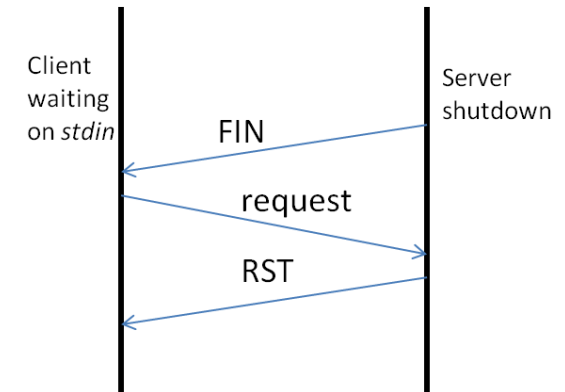
Client Handling Multiple Descriptors



- read() call on both stdin and socket will block until data is available.
- **Consider a case:**
 - If client is blocked in waiting for user to enter data, meanwhile TCP receives FIN from server.
 - Server is down. So sending request is meaningless.
- How to handle uncertainty of availability of data on descriptors?



Normal course of actions

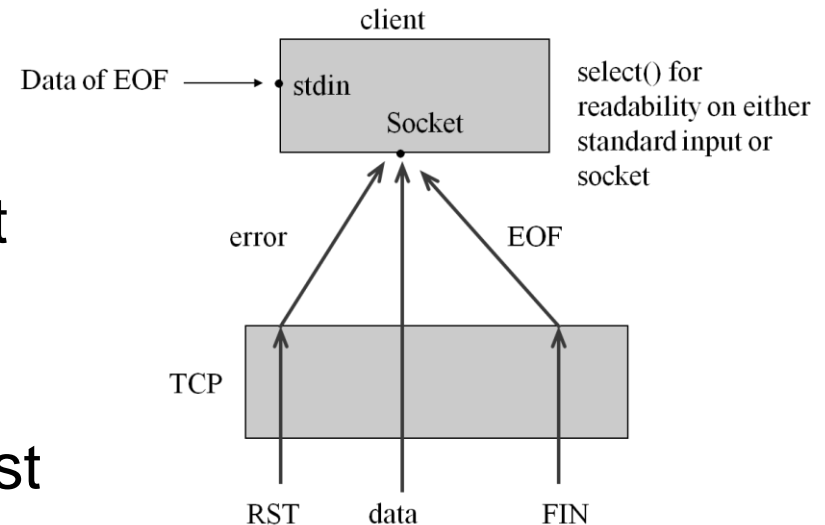


Unexpected Server shutdown

read() on socket



- Peer TCP sends data, the socket becomes readable and *read* returns greater than 0.
- Peer TCP send a FIN(peer process terminates), the socket become readable and *read* returns 0(end-of-file)
- Peer TCP send a RST(peer host has crashed and rebooted), the socket become readable and returns -1 and *errno* contains the specific error code



Client Handling Multiple Descriptors



- To avoid a situation where data has arrived from socket but client is unable to take note of it, use I/O Multiplexing.
- Client can wait on select().
 - Add stdin, socket to fd_set.
 - Call select () with fd_set for readability
 - When select() returns, find out which descriptor is ready with data.
 - Call read() on that fd.
- This will enable client to give timely response and avoid error situations.

Client Handling Multiple Descriptors

innovate

achieve

lead

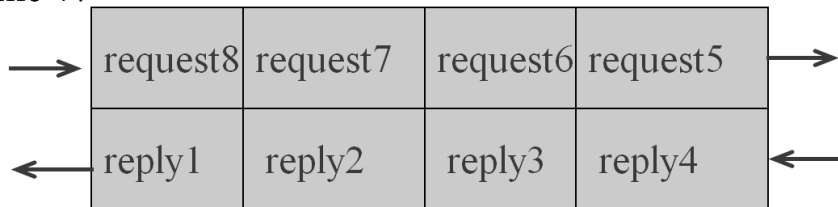
```
1 void str_cli(FILE *fp, int sockfd)
2 {
3     int maxfdp1;
4     fd_set rset;
5     char sendline[MAXLINE], recvline[MAXLINE];
6     FD_ZERO(&rset);
7     for ( ; ; ) {
8         FD_SET(fileno(fp), &rset);
9         FD_SET(sockfd, &rset);
10        maxfdp1 = max(fileno(fp), sockfd) + 1;
11        select(maxfdp1, &rset, NULL, NULL, NULL);
12        if (FD_ISSET(sockfd, &rset)) { /* socket is readable */
13            if (Readline(sockfd, recvline, MAXLINE) == 0)
14                err_quit("str_cli: server terminated prematurely");
15            Fputs(recvline, stdout);
16        }
17        if (FD_ISSET(fileno(fp), &rset)) { /* input is readable */
18            if (Fgets(sendline, MAXLINE, fp) == NULL)
19                return; /* all done */
20            Writen(sockfd, sendline, strlen(sendline));
21        }
22    } //for
23 } //str_cli
```

Batch Mode Client

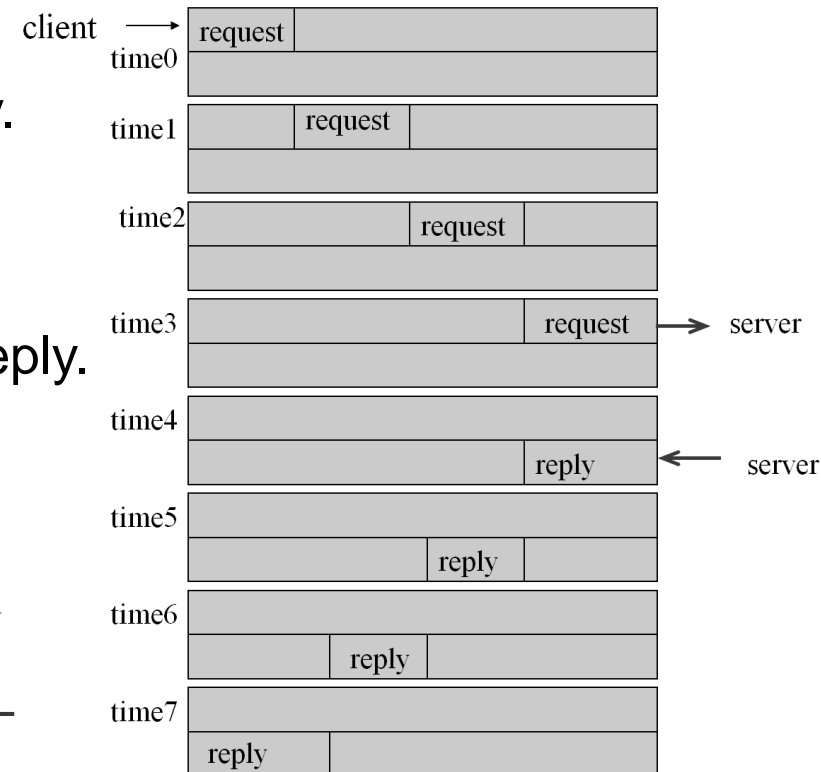
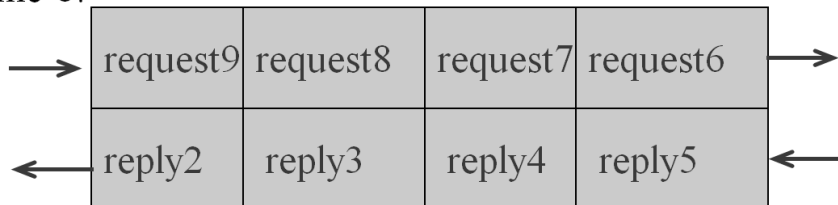


- Stop and Wait client
 - Send one request and wait for reply.
 - Usual in interactive mode.
- Batch Mode clients
 - Send requests without waiting for reply.
 - Better bandwidth utilization.

Time 7:



Time 8:



- Need for closing a socket partially:
 - We tell server that we have sent all requests by closing socket.
 - It will send FIN to server.
 - But in batch mode, by closing socket, we send FIN but we can't read replies which are yet to reach the client.
- `close()` vs `shutdown()` sys calls
 - closes the socket partially (either read end or write end) unlike `close()` sys call.
 - `close()` closes completely.
 - Irrespective of reference count it closes the socket.
 - `close()` will initiate FIN only if reference count for the fd reaches 0.

shutdown() sys call



- Sometimes, it is useful to close one half of the connection, so that data can be transmitted in just one direction through the socket.

```
1  #include <sys/socket.h>
2  int shutdown(int sockfd , int how );
3  //Returns 0 on success, or -1 on error
```

- SHUT_RD : read-half of the connection closed. Subsequent reads will return end-of-file (0).
 - SHUT_RD can't be used meaningfully for TCP sockets.
- SHUT_WR : write-half of the connection closed. Also called *socket half-close*. Buffered data will be sent followed by termination sequence.
 - Common use of shutdown()
 - Subsequent writes to the local socket yield the SIGPIPE signal and an EPIPE error.
- SHUT_RDWR : both closed
 - Note that shutdown() doesn't close the file descriptor, even if how is specified as SHUT_RDWR . To close the file descriptor, we must additionally call close().

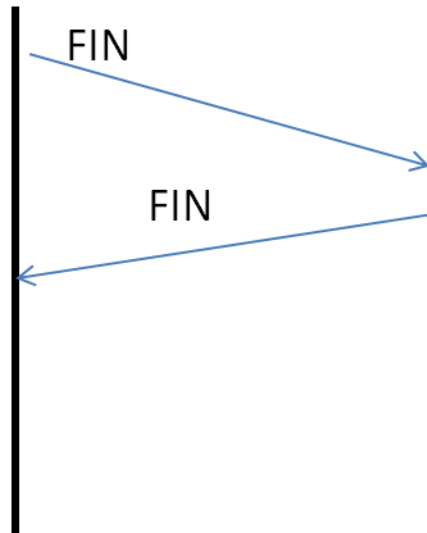
Batch Mode Client



- After user presses, Ctrl-D (EOF), close write half of the socket.
- Also set *stdineof* variable to 1.
 - This will help in inferring the FIN received from server as normal or abnormal termination.
 - In case of normal termination, we received all the replies.

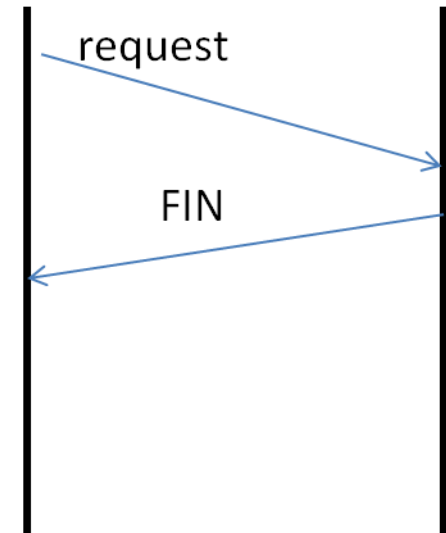
EOF on *stdin*
Set *stdineof*=1
shutdown(*fd*, *SHUT_WR*)

if *stdineof*=1
Receiving *FIN*
from server is
Normal
termination



Data on *stdin*

if *stdineof*=0
Receiving *FIN*
from server is
abnormal
termination





```
1 str_cli(FILE *fp, int sockfd)
2 {
3     int      maxfdp1, stdineof;
4     fd_set   rset;
5     stdineof = 0;
6     FD_ZERO(&rset);
7     for ( ; ; ) {
8         if (stdineof == 0)
9             FD_SET(fileno(fp), &rset);
10        FD_SET(sockfd, &rset);
11        maxfdp1 = max(fileno(fp), sockfd) + 1;
12        select(maxfdp1, &rset, NULL, NULL, NULL);
13        if (FD_ISSET(sockfd, &rset)) { /* socket is readable */
14            if ( (n = Read(sockfd, buf, MAXLINE)) == 0) {
15                if (stdineof == 1)
16                    return; /* normal termination */
17                else
18                    err_quit("str_cli: server terminated prematurely");
19            }
20            Write(fileno(stdout), buf, n);
21        }
22        if (FD_ISSET(fileno(fp), &rset)) { /* input is readable */
23            if ( (n = Read(fileno(fp), buf, MAXLINE)) == 0) {
24                stdineof = 1;
25                shutdown(sockfd, SHUT_WR); /* send FIN */
26                FD_CLR(fileno(fp), &rset);
27                continue;
28            }
29            Writen(sockfd, buf, n);
30        }
31    } }
```

TCP Server Using select()



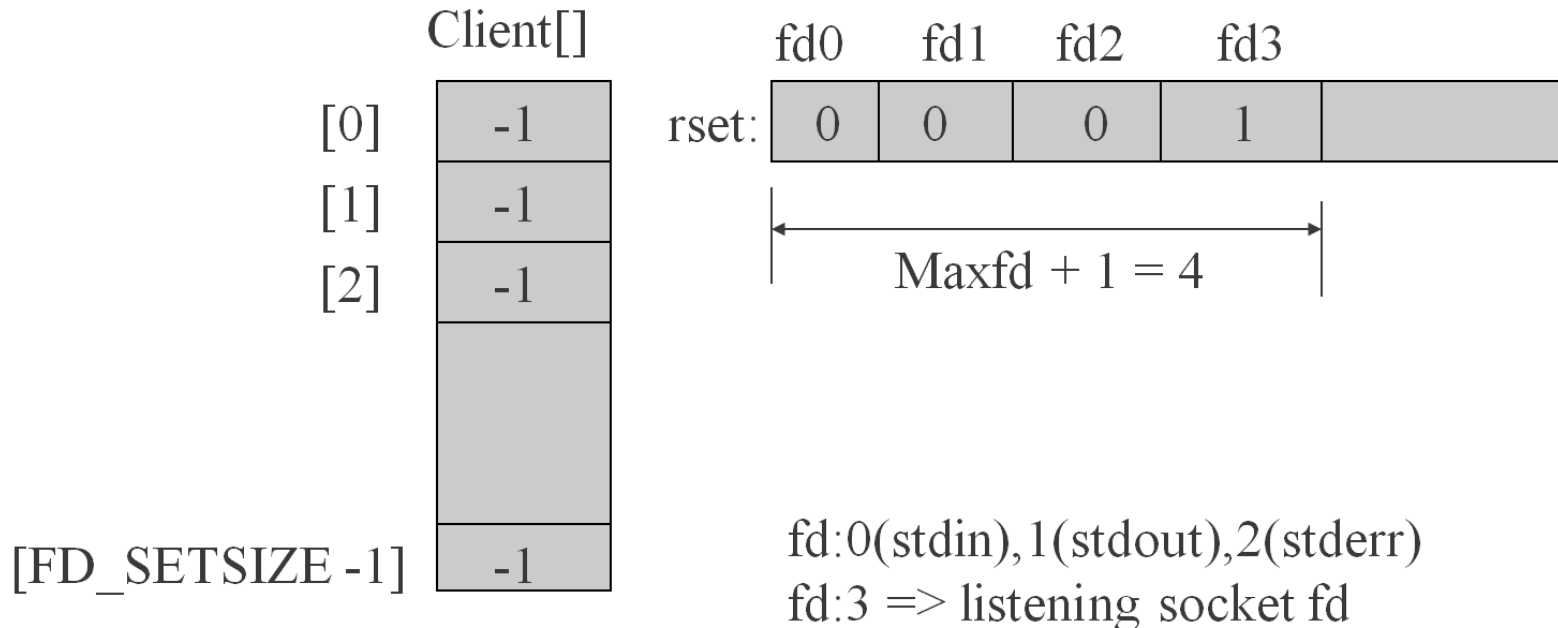
- Single process server that uses select to handle any number of clients, instead of forking one child per client.
- Protocol: echo
- Two data structures:
 - Client array
 - Keeps list of client sockets connected currently
 - fd_set *allset*
 - Keeps list of fds for checking against readability.

TCP Server Using select()



- There are three fds: 0,1,2
- One more fd after creating listening socket.

Before first client has established a connection

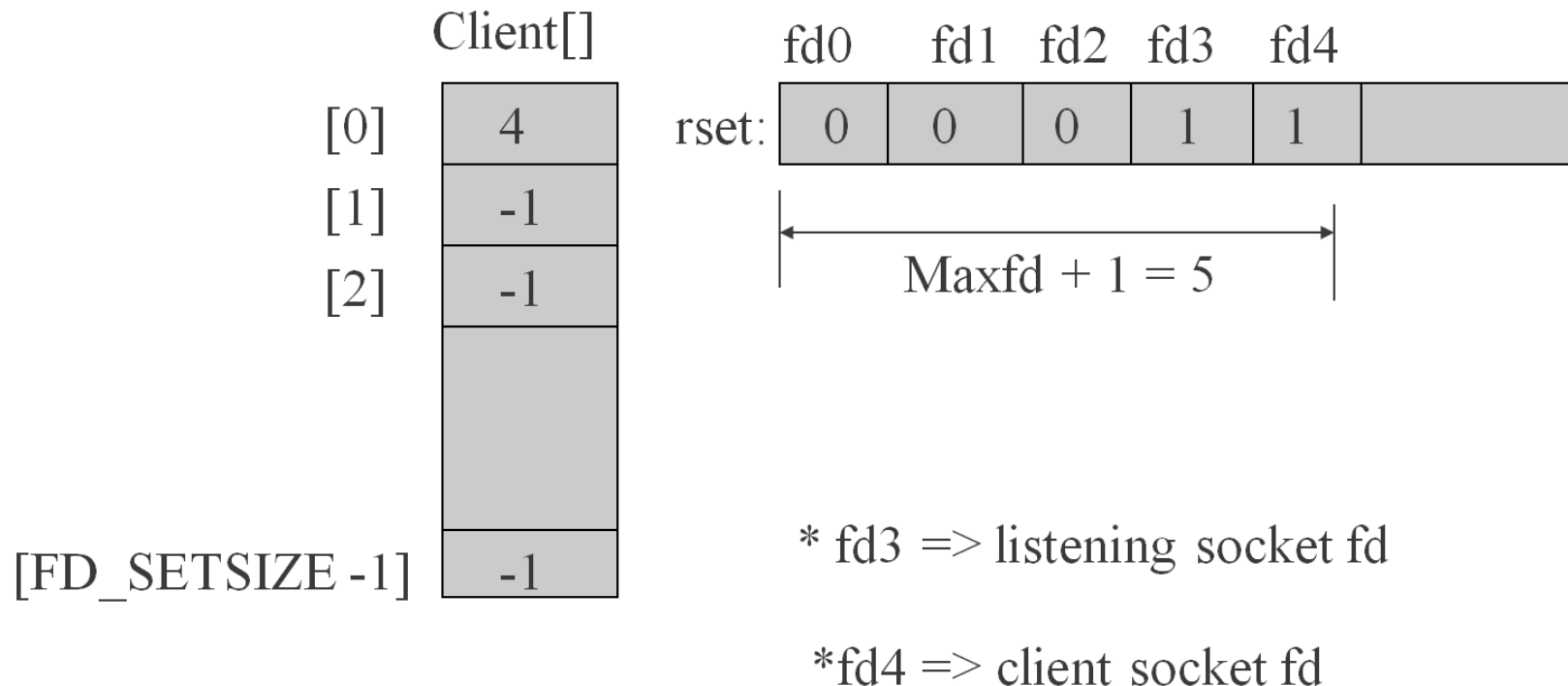


TCP Server Using select()



- When a new client is accepted through accept()
 - A connected socket is added

After first client connection is established

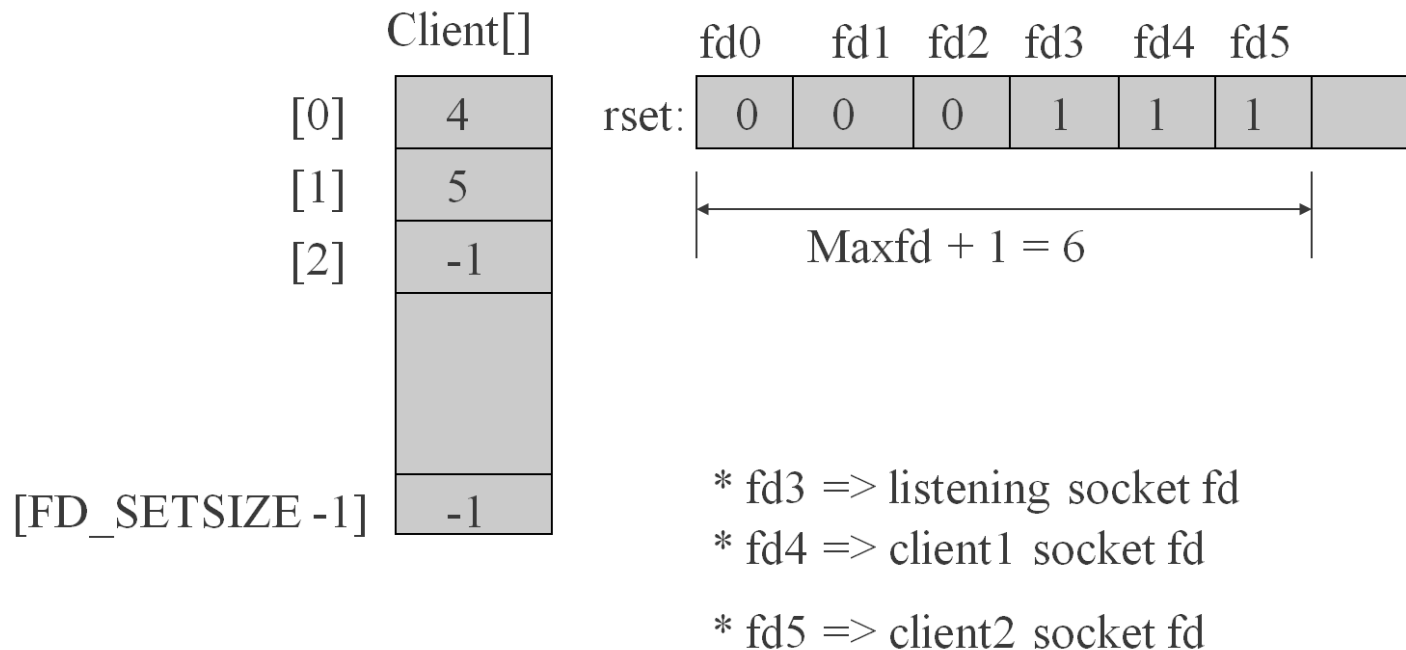


TCP Server Using select()



- When second client is accepted

After second client connection is established

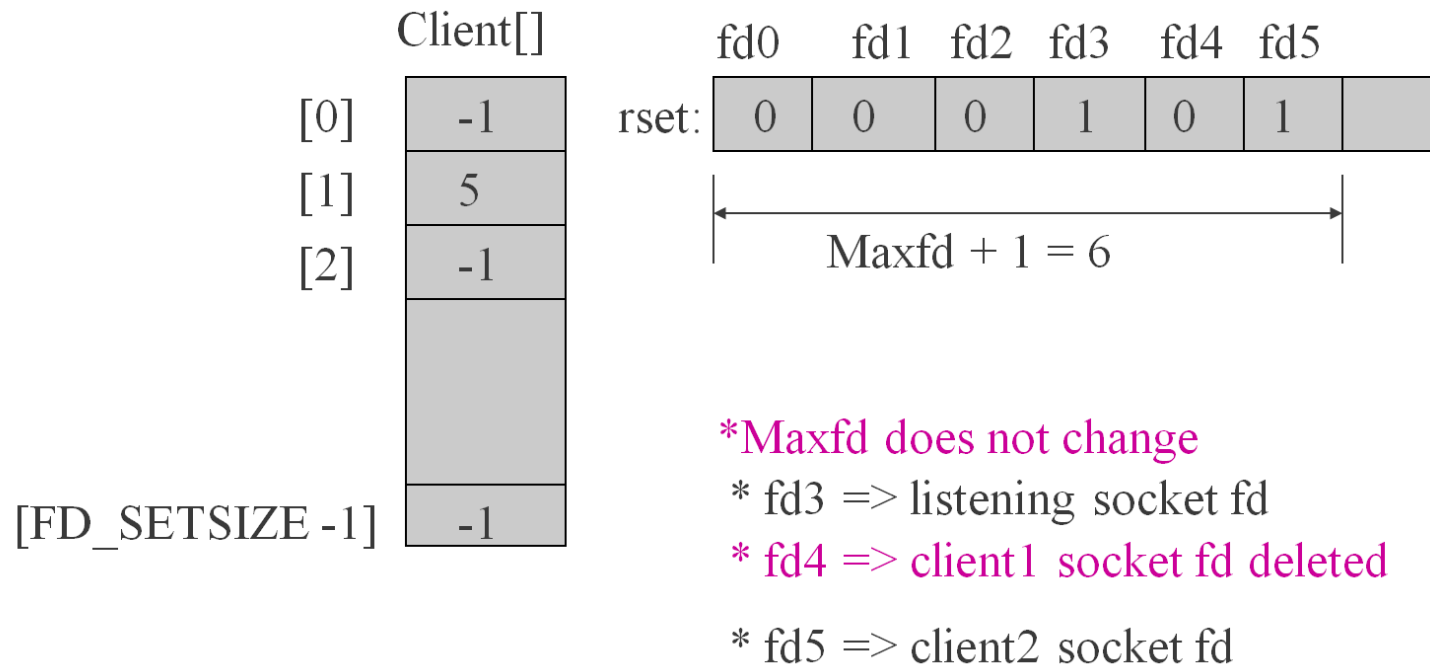


TCP Server Using select()



- When the first client terminates connection
 - This is known when read() returns zero.

After first client terminates its connection



TCP Server Using select()



- Create a passive socket.

```
1  int main(int argc, char **argv)
2  {
3      int          i, maxi, maxfd, listenfd, connfd, sockfd;
4      int          nready, client[FD_SETSIZE];
5      fd_set       rset, allset;
6      struct sockaddr_in cliaddr, servaddr;
7      listenfd = socket(AF_INET, SOCK_STREAM, 0);
8      bzero(&servaddr, sizeof(servaddr));
9      servaddr.sin_family = AF_INET;
10     servaddr.sin_addr.s_addr = htonl(INADDR_ANY);
11     servaddr.sin_port = htons(SERV_PORT);
12     bind(listenfd, (SA *) &servaddr, sizeof(servaddr));
13     listen(listenfd, LISTENQ);
```

TCP Server Using select()



- Handling when listening socket is readable

```
1  maxfd = listenfd;           /* initialize */
2  maxi = -1;                  /* index into client[] array */
3  for (i = 0; i < FD_SETSIZE; i++)
4  client[i] = -1;             /* -1 indicates available entry */
5  FD_ZERO(&allset);
6  FD_SET(listenfd, &allset);
7  for ( ; ; ) {
8      rset = allset;          /* structure assignment */
9      nready = select(maxfd+1, &rset, NULL, NULL, NULL);

10     if (FD_ISSET(listenfd, &rset)) { /* new client connection */
11         cliilen = sizeof(cliaddr);
12         connfd = accept(listenfd, (SA *) &cliaddr, &cliilen);
13         for (i = 0; i < FD_SETSIZE; i++)
14             if (client[i] < 0) {
15                 client[i] = connfd;    /* save descriptor */
16                 break;
17             }
18         if (i == FD_SETSIZE) err_quit("too many clients");
19         FD_SET(connfd, &allset);      /* add new descriptor to set */
20         if (connfd > maxfd)
21             maxfd = connfd;           /* maxfd for select */
22         if (i > maxi)
23             maxi = i;                 /* max index in client[] array */
24         if (--nready <= 0)
25             continue;                /* no more readable descriptors */
26     }
```

TCP Server Using select()



- When a connected socket is readable

```
1  for (i = 0; i <= maxi; i++) { /* check all clients for data */
2      if ((sockfd = client[i]) < 0)
3          continue;
4      if (FD_ISSET(sockfd, &rset)) {
5          if ( (n = Readline(sockfd, line, MAXLINE)) == 0) {
6              /*connection closed by client */
7              close(sockfd);
8              FD_CLR(sockfd, &allset);
9              client[i] = -1;
10         }
11         else
12             Writen(sockfd, line, n);
13         if (--nready <= 0)
14             break; /* no more readable descriptors */
15     }
16 }
```

- This code looks complicated when compared to fork-per-client model. But this design avoids overhead of fork().

Denial-of-Service Attacks



- If malicious client connect to the server, send 1 byte of data (other than a newline), and then goes to sleep.
 - in readline(), server is blocked.
- Solution
 - use nonblocking I/O
 - have each client serviced by a separate thread of control (spawn a process or a thread to service each client)
 - place a timeout on the I/O operation

pselect()

```
1 struct timespec{
2     time_t  tv_sec;    /*seconds*/
3     long    tv_nsec;  /* nanoseconds */
4 };
```



- pselect contains two changes from the normal select function:
 - pselect uses the timespec structure instead of the timeval structure.
 - Accepts signal mask.

```
1 #define _XOPEN_SOURCE 600
2 #include <sys/select.h>
3 int pselect(int  nfds , fd_set * readfds , fd_set * writefds ,
4     fd_set * exceptfds, struct timespec * timeout , const sigset_t * sigmask );
5 //Returns number of ready file descriptors, 0 on timeout, or -1 on error
```

```
2 ready = pselect(nfds, &readfds, &writefds, &exceptfds, timeout, &sigmask);
```

- This call is equivalent to

```
1 sigset_t origmask;
2 sigprocmask(SIG_SETMASK, &sigmask, &origmask);
3 ready = select(nfds, &readfds, &writefds, &exceptfds, timeout);
4 sigprocmask(SIG_SETMASK, &origmask, NULL);          /* Restore signal mask */
```


Problems with select() and poll()



- The select() and poll() system calls are the portable, long-standing, and widely used methods of monitoring multiple file descriptors for readiness.
- Suffer from some problems
 - Kernel must check all the fds to check if they are ready.
 - Each time select() passes data structures which kernel modifies and returns.
 - Once select() returns, the program must inspect the data structure to see which fds are ready.
- Select() scales poorly with the increase of *fds*.
- Signal driven I/O or *epoll* (event poll) provide a scalable solution.



BITS Pilani
Pilani Campus



Non-blocking I/O on Sockets

Socket Operations



- Input operations: read, recv, readv, recvfrom, recvmsg
 - Blocking operations
 - TCP: until a byte arrives.
 - UDP: until a datagram arrives.
 - With non-blocking socket, if no data, return with EWOULDBLOCK error.
- Output Operations: write, send, writev, sendto, sendmsg
 - Blocks if there is no room in socket send buffer.
 - TCP: until all the data is written.
 - UDP: no send buffer present.
 - With non-blocking socket, TCP write will write whatever it can and returns no. of bytes written. If no room at all, it returns with error EWOULDBLOCK.

Socket Operations

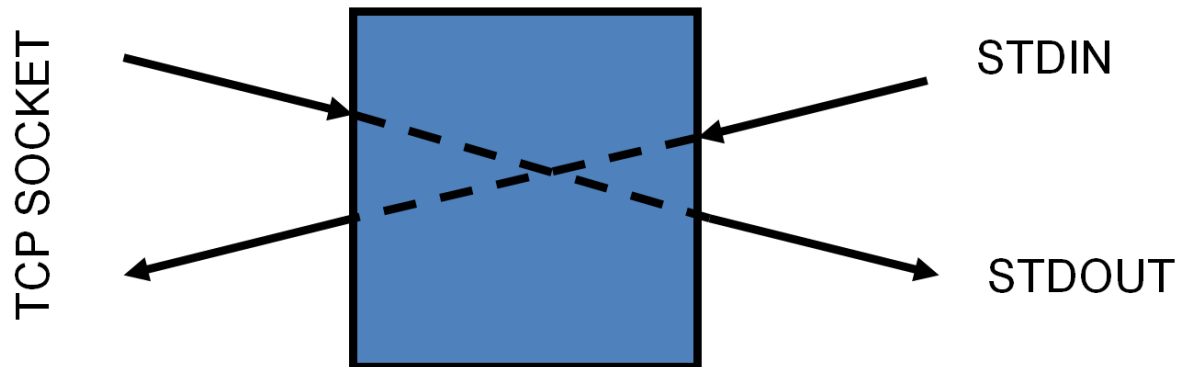


- Accepting incoming connections: `accept`
 - Blocks if no incoming connection.
 - With non-blocking socket, it would return with an error.
- Initiating Connections: `connect`
 - Blocks until client TCP receives ACK.
 - With non-blocking socket, it returns *errno* EINPROGRESS, and continues to establish connection.

Client Handling a socket, stdin, stdout



- A client usually deals with
 - Stdin
 - Stdout
 - Socket



Client Handling a socket, stdin, stdout



- We looked at
 - Stop and wait client
 - Batch mode client - select with blocking I/O
 - Once select() returns, and if socket is readable, read() is called on socket.
 - readline() call gets blocked on socket till it gets required data.
 - During this time, other clients have to wait.
- Now we look at select with non-blocking I/O
 - In this, read() will be a non-blocking operation. It will read whatever data available on socket. It returns.
 - When this fd is readable next time, further data is read.
 - This requires that we track the number of bytes read and the pointers in the buffer.



```
1 str_cli(FILE *fp, int sockfd)
2 {
3     int      maxfdp1, stdineof;
4     fd_set   rset;
5     stdineof = 0;
6     FD_ZERO(&rset);
7     for ( ; ; ) {
8         if (stdineof == 0)
9             FD_SET(fileno(fp), &rset);
10        FD_SET(sockfd, &rset);
11        maxfdp1 = max(fileno(fp), sockfd) + 1;
12        select(maxfdp1, &rset, NULL, NULL, NULL);
13        if (FD_ISSET(sockfd, &rset)) { /* socket is readable */
14            if ( (n = Read(sockfd, buf, MAXLINE)) == 0) {
15                if (stdineof == 1)
16                    return; /* normal termination */
17                else
18                    err_quit("str_cli: server terminated prematurely");
19            }
20            Write(fileno(stdout), buf, n);
21        }
22        if (FD_ISSET(fileno(fp), &rset)) { /* input is readable */
23            if ( (n = Read(fileno(fp), buf, MAXLINE)) == 0) {
24                stdineof = 1;
25                shutdown(sockfd, SHUT_WR); /* send FIN */
26                FD_CLR(fileno(fp), &rset);
27                continue;
28            }
29            Writen(sockfd, buf, n);
30        }
31    }
```

Select with Non-Blocking IO



- Non-blocking IO complicates buffer management.
 - We have to keep track of how much is read and how much is written.
- Two buffers:
 - *to*: reading from standard input and write to socket.
 - *from*: read from socket and write to stdout.

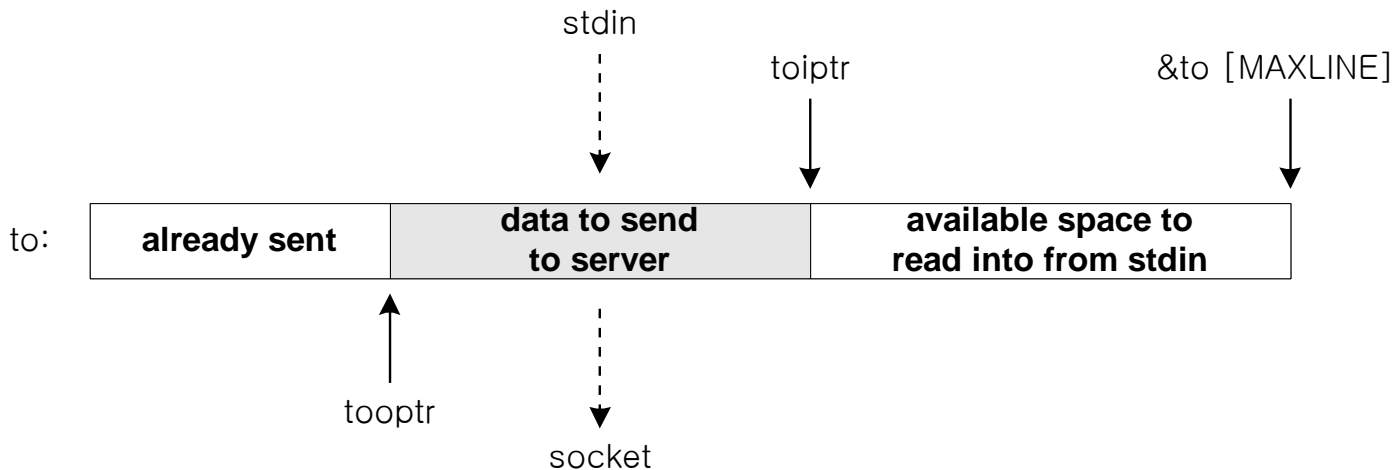


Figure 15.1 Buffer containing data from standard input going to the socket.

Select with Non-Blocking IO

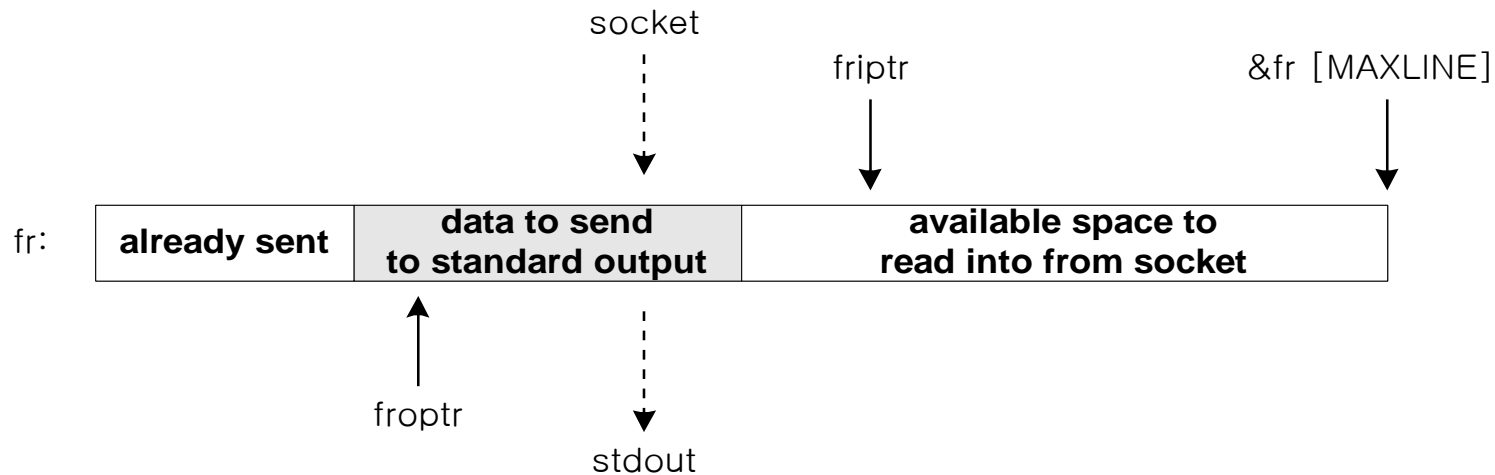


Figure 15.2 Buffer containing data from the socket going to standard output.

- *froptr*: points to the next byte to be sent to stdout.
- *froiptr*: points to the next byte into which next byte can be read..

Select with Non-Blocking IO



```
1 void str_cli(FILE *fp, int sockfd)
2 {
3     int      maxfdp1, val, stdineof;
4     ssize_t n, nwritten;
5     fd_set   rset, wset;
6     char     to[MAXLINE], fr[MAXLINE];
7     char     *toiptr, *tooptr, *friptr, *froptr;
8     val = fcntl(sockfd, F_GETFL, 0);
9     fcntl(sockfd, F_SETFL, val | O_NONBLOCK);
10    val = Fcntl(STDIN_FILENO, F_GETFL, 0);
11    fcntl(STDIN_FILENO, F_SETFL, val | O_NONBLOCK);
12    val = Fcntl(STDOUT_FILENO, F_GETFL, 0);
13    fcntl(STDOUT_FILENO, F_SETFL, val | O_NONBLOCK);
14    toiptr = tooptr = to;      /* initialize buffer pointers */
15    friptr = froptr = fr;
16    stdineof = 0;
17
18    maxfdp1 = max(max(STDIN_FILENO, STDOUT_FILENO), sockfd) + 1;
19    for ( ; ; ) {
20        FD_ZERO(&rset);
21        FD_ZERO(&wset);
22        if (stdineof == 0 && toiptr < &to[MAXLINE])
23            FD_SET(STDIN_FILENO, &rset);      /* read from stdin */
24        if (friptr < &fr[MAXLINE])
25            FD_SET(sockfd, &rset); /* read from socket */
26        if (tooptr != toiptr)
27            FD_SET(sockfd, &wset); /* data to write to socket */
28        if (froptr != friptr)
29            FD_SET(STDOUT_FILENO, &wset); /* data to write to stdout */
30        select(maxfdp1, &rset, &wset, NULL, NULL);
```

reads from standard input

innovate

achieve

lead

```
30 select(maxfdp1, &rset, &wset, NULL, NULL);
31 if (FD_ISSET(STDIN_FILENO, &rset)) {
32     if((n = read(STDIN_FILENO, toiptr, &to[MAXLINE] - toiptr)) < 0) {
33         if (errno != EWOULDBLOCK)
34             err_sys("read error on stdin");
35     } else if (n == 0) {
36         fprintf(stderr, "%s: EOF on stdin\n", gf_time());
37         stdineof = 1; /* all done with stdin */
38         if (tooptr == toiptr)
39             shutdown(sockfd, SHUT_WR); /* send FIN */
40     } else {
41         fprintf(stderr, "%s: read %d bytes from stdin\n", gf_time(),
42                 n);
43         toiptr += n; /* # just read */
44         FD_SET(sockfd, &wset); /* try and write to socket below */
45     }
46 }
```

Amt of space available in to buffer

If user has pressed Ctrl-D,
set stdineof=1
If no outstanding data on
buffer, close the write end.

Increment to pointer
set socket in wset for writability

reads from socket

innovate

achieve

lead

Amt of space available in from buffer

```
47 if (FD_ISSET(sockfd, &rset)) {
48     if ( (n = read(sockfd, friptr, &fr[MAXLINE] - friptr)) < 0) {
49         if (errno != EWOULDBLOCK)
50             err_sys("read error on socket");
51     } else if (n == 0) {
52         fprintf(stderr, "%s: EOF on socket\n", gf_time());
53         if (stdineof)
54             return;      /* normal termination */
55         else
56             err_quit("str_cli: server terminated prematurely");
57     } else {
58         fprintf(stderr, "%s: read %d bytes from socket\n",
59             gf_time(), n);
60         friptr += n;      /* # just read */
61         FD_SET(STDOUT_FILENO, &wset);    /* try and write below */
62     }
63 }
```

Increment friptr.

Add stdout to wset to test for writability.

writes to standard output

innovate

achieve

lead

No of bytes to write >0

```
65 ▾ if (FD_ISSET(STDOUT_FILENO, &wset) && ((n = friptr - froptr) > 0)) {
66 ▾     if ( (nwritten = write(STDOUT_FILENO, froptr, n)) < 0) {
67         if (errno != EWOULDBLOCK)
68             err_sys("write error to stdout");
69 ▾     } else {
70         fprintf(stderr, "%s: wrote %d bytes to stdout\n",
71             gf_time(), nwritten);
72         froptr += nwritten; /* # just written */
73         if (froptr == friptr)
74             froptr = friptr = fr; /* back to beginning of buffer */
75     }
76 }
```

If the write is successful, froptr is incremented by the number of bytes written

Writes to socket

innovate

achieve

lead

No of bytes to write >0

```
78  if (FD_ISSET(sockfd, &wset) && ((n = toiptr - tooptr) > 0)) {
79      if ( (nwritten = write(sockfd, tooptr, n)) < 0) {
80          if (errno != EWOULDBLOCK)
81              err_sys("write error to socket");
82      } else {
83          fprintf(stderr, "%s: wrote %d bytes to socket\n",
84                  gf_time(), nwritten);
85          tooptr += nwritten; /* # just written */
86          if (tooptr == toiptr) {
87              toiptr = tooptr = to; /* back to beginning of buffer */
88              if (stdineof)
89                  shutdown(sockfd, SHUT_WR); /* send FIN */
90          }
91      }
92  }
93  }
94  }
```

If the write is successful, tooptr is incremented by the number of bytes written

if we encountered an EOF on standard input, the FIN can be sent to the server

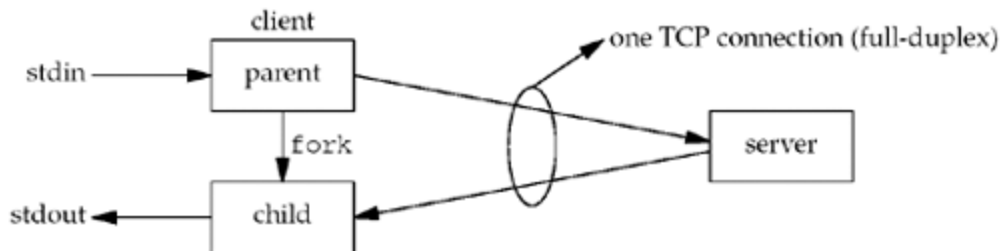


Client using Multiple Processes

Client using Multiple Processes



- Whenever we find the need to use nonblocking I/O, it will usually be simpler to split the application into either processes (using fork) or threads.
 - Parent reads from stdin and writes to socket
 - Child reads from socket and writes to stdout.



- Normal termination occurs when the EOF on standard input is encountered. The parent reads this EOF and calls shutdown to send a FIN.
- If abnormal occurs, the child will read an EOF on the socket. If this happens, the child must tell the parent to stop copying from the standard input to the socket
 - the child sends a signal (e.g. SIGTERM) to the parent.

Client using Multiple Processes



```
1 void str_cli(FILE *fp, int sockfd)
2 {
3     pid_t    pid;
4     char      sendline[MAXLINE], recvline[MAXLINE];
5     if ( (pid = fork()) == 0) { /* child: server -> stdout */
6         while (Readline(sockfd, recvline, MAXLINE) > 0)
7             fputs(recvline, stdout);
8         kill(getppid(), SIGTERM); /* in case parent still running */
9         exit(0);
10    }
11    /* parent: stdin -> server */
12    while (fgets(sendline, MAXLINE, fp) != NULL)
13        Writen(sockfd, sendline, strlen(sendline));
14    shutdown(sockfd, SHUT_WR); /* EOF on stdin, send FIN */
15    pause();
16    return;
17 }
```

Comparing Cleint Designs



- when copying 2,000 lines from a client to a server with an RTT of 175 ms:

Client	Time taken for sending and receiving
stop-and-wait	354.0 sec
select and blocking I/O	12.3 sec
nonblocking I/O	6.9 sec
fork	8.7 sec
threaded version	8.5 sec

- nonblocking I/O version is almost twice as fast as version using blocking I/O with select.
- Version using fork is slower than nonblocking I/O version.
- Nevertheless, given the complexity of the nonblocking I/O code versus the fork code, fork version is simple approach.



BITS Pilani
Pilani Campus



Non-blocking Connect

Nonblocking connect()



- TCP socket nonblocking connect
 - return: an error of EINPROGRESS
 - TCP three-way handshake continues
 - check the connection establishment using select
- There are three uses for a nonblocking connect.
 - We can overlap other processing with the three-way handshake.
 - We can establish multiple connections at the same time using this technique.
 - popular with Web browsers
 - Since we wait for the connection establishment to complete using select, we can specify time limit for select, allowing us to shorten the timeout for the connect.

Nonblocking connect()



- Set the socket to non-blocking.
- Call `connect()`. It will return immediately with error `EINPROGRESS`.
- We use `select()` to check what has happened to `connect()`.
- If the descriptor is readable or writable, we call `getsockopt()` to fetch the socket's pending error (`SO_ERROR`). If the connection completed successfully, this value will be 0.

Nonblocking connect()



```
1 int connect_nonb(int sockfd, const SA *saptr, socklen_t salen, int nsec)
2 {
3     int flags, n, error;
4     socklen_t len;
5     fd_set rset, wset;
6     struct timeval tval;
7     flags = fcntl(sockfd, F_GETFL, 0);
8     fcntl(sockfd, F_SETFL, flags | O_NONBLOCK);
9     error = 0;
10    if ( (n = connect(sockfd, saptr, salen)) < 0)
11        if (errno != EINPROGRESS)
12            return (-1);
13    /* Do whatever we want while the connect is taking place. */
14    if (n == 0)
15        goto done; /* connect completed immediately */
16    FD_ZERO(&rset);
17    FD_SET(sockfd, &rset);
18    wset = rset;
19    tval.tv_sec = nsec;
20    tval.tv_usec = 0;
```

Nonblocking connect()

innovate

achieve

lead

```
22     if ( (n = select(sockfd + 1, &rset, &wset, NULL,
23         nsec ? &tval : NULL)) == 0) {
24         close(sockfd);          /* timeout */
25         errno = ETIMEDOUT;
26         return (-1);
27     }
28     if (FD_ISSET(sockfd, &rset) || FD_ISSET(sockfd, &wset)) {
29         len = sizeof(error);
30         if (getsockopt(sockfd, SOL_SOCKET, SO_ERROR, &error, &len) < 0)
31             return (-1);      /* Solaris pending error */
32     } else
33         err_quit("select error: sockfd not set");
34 done:
35     Fcntl(sockfd, F_SETFL, flags); /* restore file status flags */
36     if (error) {
37         close(sockfd);          /* just in case */
38         errno = error;
39         return (-1);
40     }
41     return (0);
42 }
```



BITS Pilani
Pilani Campus



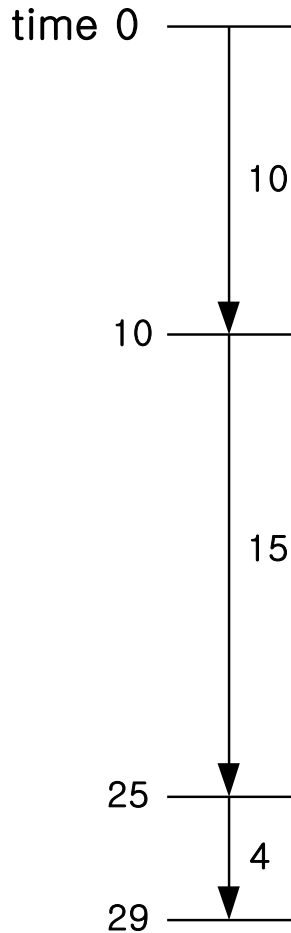
Web client: Non-blocking Connect

nonblocking connect: web client

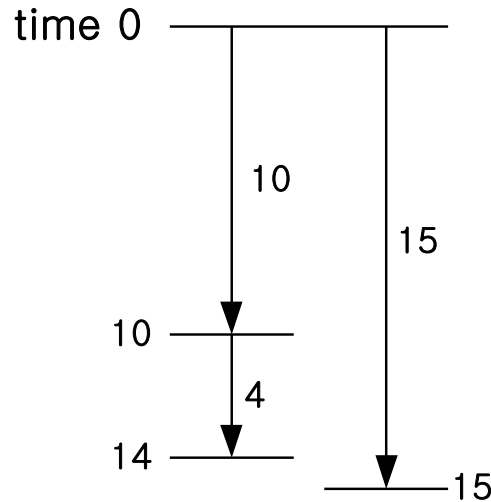


- A real-world example of nonblocking connects started with Netscape Web Client
- The client establishes an HTTP connection with a Web server and fetches a home page.
- On that page are often numerous references to other Web pages.
- Instead of fetching these other pages serially, one at a time, the client can fetch more than one at the same time, using nonblocking connects.

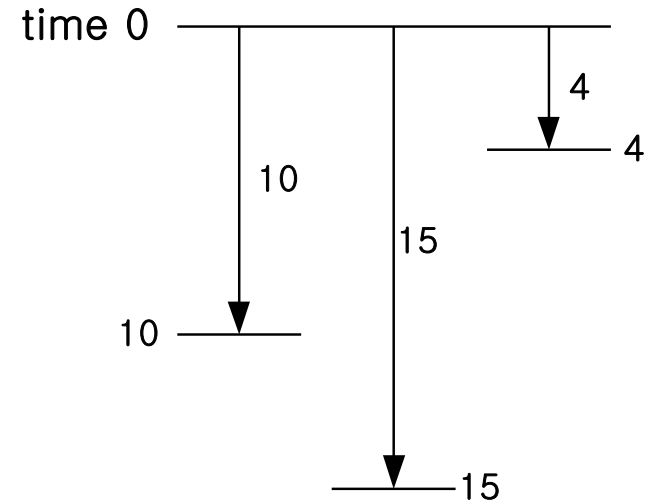
nonblocking connect: web client (cont.)



three connections
done serially



three connections
done in parallel;
maximum of two
connections at a time



three connections
done in parallel;
maximum of three
connections at a time

Non-blocking Connect



- This program will read up to 20 files from a Web server.
- We specify as command-line arguments
 - the maximum number of parallel connections,
 - the server's hostname, and
 - each of the filenames to fetch from the server.

```
2  bash$ web 3 www.foobar.com image1.gif image2.gif image3.gif image4.gif
3  image5.gif image6.gif image7.gif
```

- It means
 - three simultaneous connection
 - server's hostname
 - filename for the home page
 - the files to be read

```
1  #define MAXFILES      20
2  #define SERV          "80"          /* port number or service name */
3  struct file {
4      char    *f_name;                /* filename */
5      char    *f_host;                /* hostname or IPv4/IPv6 address */
6      int     f_fd;                   /* descriptor */
7      int     f_flags;                /* F_xxx below */
8  } file[MAXFILES];
9  #define F_CONNECTING    1           /* connect() in progress */
10 #define F_READING       2           /* connect() complete; now reading */
11 #define F_DONE          4           /* all done */
12 #define GET_CMD         "GET %s HTTP/1.0\r\n\r\n"
13 /* globals */
14 int     nconn, nfiles, nlefttoconn, nlefttoread, maxfd;
15 fd_set  rset, wset;
16 /* function prototypes */
17 void    home_page(const char *, const char *);
18 void    start_connect(struct file *);
19 void    write_get_cmd(struct file *);
```

- Each file has a state and fd.

nonblock/web.c



```
1 main(int argc, char **argv)
2 {
3     int i, fd, n, maxnconn, flags, error;
4     char buf[MAXLINE];
5     fd_set rs, ws;
6     if (argc < 5)
7         err_quit("usage: web <#conns> <hostname> <homepage> <file1> ...");
8     maxnconn = atoi(argv[1]);
9     nfiles = min(argc - 4, MAXFILES);
10    for (i = 0; i < nfiles; i++) {
11        file[i].f_name = argv[i + 4];
12        file[i].f_host = argv[2];
13        file[i].f_flags = 0;
14    }
15    printf("nfiles = %d\n", nfiles);
16    home_page(argv[2], argv[3]);
17    FD_ZERO(&rset);
18    FD_ZERO(&wset);
19    maxfd = -1;
20    nlefttoread = nlefttoconn = nfiles;
21    nconn = 0;
```

- Process command-line arguments
- Read home page
- Initialize globals
 - Fd sets, nconn is current number of connections.

nonblock/start_connect.c



```
1 void start_connect(struct file *fptr)
2 {
3     int    fd, flags, n;
4     ai = Host_serv(fptr->f_host, SERV, 0, SOCK_STREAM);
5     fd = Socket(ai->ai_family, ai->ai_socktype, ai->ai_protocol);
6     fptr->f_fd = fd;
7     printf("start_connect for %s, fd %d\n", fptr->f_name, fd);
8     /* Set socket nonblocking */
9     flags = Fcntl(fd, F_GETFL, 0);
10    Fcntl(fd, F_SETFL, flags | O_NONBLOCK);
11    /* Initiate nonblocking connect to the server. */
12    if ( (n = connect(fd, ai->ai_addr, ai->ai_addrlen)) < 0) {
13        if (errno != EINPROGRESS)
14            err_sys("nonblocking connect error");
15        fptr->f_flags = F_CONNECTING;
16        FD_SET(fd, &rset);    /* select for reading and writing */
17        FD_SET(fd, &wset);
18        if (fd > maxfd)
19            maxfd = fd;
20    } else if (n >= 0)        /* connect is already done */
21        write_get_cmd(fptr); /* write() the GET command */
22 }
```

- Initiate nonblocking connect
- Handle connection complete
- If connect returns successfully, the connection is already complete and the function write_get_cmd ends a command to the server.

nonblock/write_get_cmd.c



```
1  #include    "web.h"
2  void write_get_cmd(struct file *fptr)
3  {
4      int      n;
5      char     line[MAXLINE];
6      n = snprintf(line, sizeof(line), GET_CMD, fptr->f_name);
7      Writen(fptr->f_fd, line, n);
8      printf("wrote %d bytes for %s\n", n, fptr->f_name);
9      fptr->f_flags = F_READING; /* clears F_CONNECTING */
10     FD_SET(fptr->f_fd, &rset); /* will read server's reply */
11     if (fptr->f_fd > maxfd)
12         maxfd = fptr->f_fd;
13 }
```

- Build command and send it
- Set flags

Main function: web.c



```
1 while (nlefttoread > 0) {
2     while (nconn < maxnconn && nlefttoconn > 0) {
3         /* find a file to read */
4         for (i = 0; i < nfiles; i++)
5             if (file[i].f_flags == 0)
6                 break;
7         if (i == nfiles)
8             err_quit("nlefttoconn = %d but nothing found", nlefttoconn);
9         start_connect(&file[i]);
10        nconn++;
11        nlefttoconn--;
12    }
```

- Initiate another connection, if possible

Main function: web.c



```
13  rs = rset;
14  ws = wset;
15  n = select(maxfd + 1, &rs, &ws, NULL, NULL);
16  for (i = 0; i < nfiles; i++) {
17      flags = file[i].f_flags;
18      if (flags == 0 || flags & F_DONE)
19          continue;
20      fd = file[i].f_fd;
21      if (flags & F_CONNECTING &&
22          (FD_ISSET(fd, &rs) || FD_ISSET(fd, &ws))) {
23          n = sizeof(error);
24          if (getsockopt(fd, SOL_SOCKET, SO_ERROR, &error, &n) < 0 ||
25              error != 0) {
26              err_ret("nonblocking connect failed for %s",
27                      file[i].f_name);
28          }
29          /* connection established */
30          printf("connection established for %s\n", file[i].f_name);
31          FD_CLR(fd, &ws); /* no more writability test */
32          write_get_cmd(&file[i]); /* write() the GET command */

```

- select waits for either readability or writability.
 - Descriptors that have a nonblocking connect in progress will be enabled in both sets, while descriptors with a completed connection that are waiting for data from the server will be enabled in just the read set.

Main function: web.c

innovate

achieve

lead

```
29  /* connection established */
30  printf("connection established for %s\n", file[i].f_name);
31  FD_CLR(fd, &wset); /* no more writeability test */
32  write_get_cmd(&file[i]); /* write() the GET command */
33  } else if (flags & F_READING && FD_ISSET(fd, &rs)) {
34  if ( (n = Read(fd, buf, sizeof(buf))) == 0) {
35      printf("end-of-file on %s\n", file[i].f_name);
36      Close(fd);
37      file[i].f_flags = F_DONE; /* clears F_READING */
38      FD_CLR(fd, &rset);
39      nconn--;
40      nlefttoread--;
41  } else {
42      printf("read %d bytes from %s\n", n, file[i].f_name);
43  }
44  }
45  }
46  }
47  exit(0);
48  }
```

- If the F_READING flag is set and the descriptor is ready for reading, we call read.

Performance of Nonblocking Connect



- Table shows the clock time required to fetch a Web server's home page, followed by nine image files from that server.
 - The RTT to the server is about 150 ms.
 - The home page size was 4,017 bytes and the average size of the 9 image files was 1,621 bytes.
 - TCP's segment size was 512 bytes.
- Most of the improvement is obtained with three simultaneous connections.

# simultaneou s connec tions	Clock time (seconds), non blocking	Clock time(sec s) Threads
1	6.0	6.3
2	4.1	4.2
3	3.0	3.1
4	2.8	3.0
5	2.5	2.7
6	2.4	2.5
7	2.3	2.3
8	2.2	2.3
9	2.0	2.2

Q&A





BITS Pilani
Pilani Campus



Thank You