



BITS Pilani
Pilani Campus

Network Programming

K Hari Babu
Department of Computer Science & Information Systems



BITS Pilani
Pilani Campus



Outline

- Advanced I/O Functions
 - How much data is queued?
 - Sockets and standard I/O
- Unix Domain Sockets
 - Socket pair
 - Stream sockets
 - Datagram sockets
 - Passing descriptors
 - Passing Credentials
- Unix I/O Models
 - Blocking
 - Non-blocking
 - I/O multiplexing
 - Signal driven
 - Asynchronous
- I/O Multiplexing
 - select()
 - Client using select()
 - Concurrent server using select()



Advanced I/O

T1: 14.7 - 14.8

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.



BITS Pilani
Pilani Campus

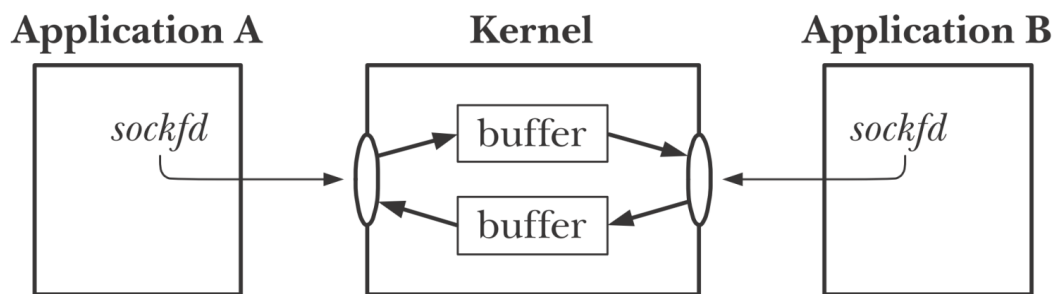


Unix Domain Sockets

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

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()

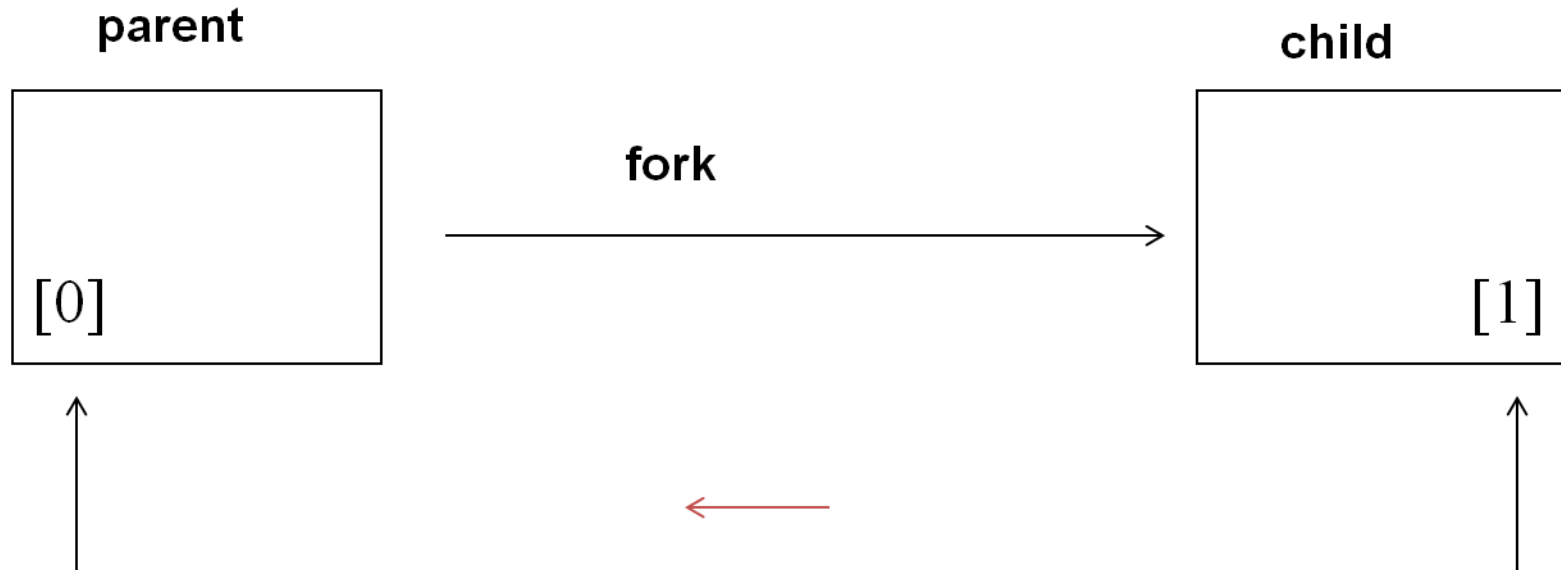


- creates a pair of sockets and connect them together.

```
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

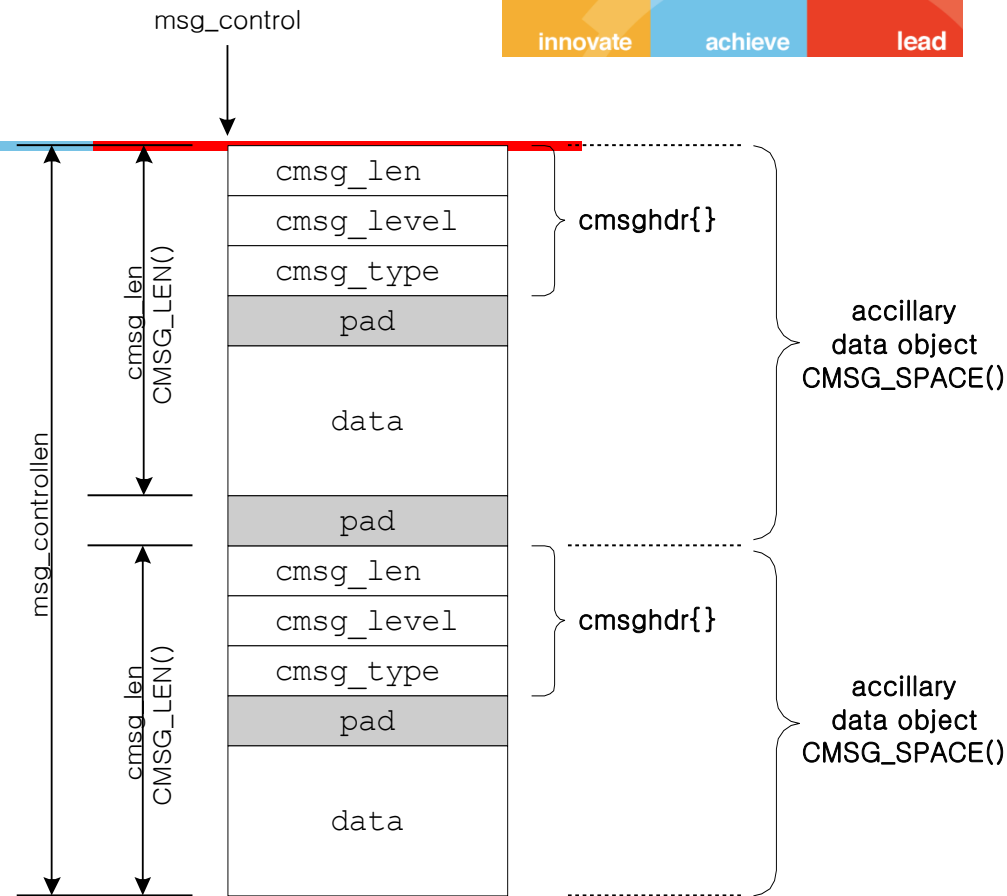


- Steps involved
 - Create a unix domain socket(stream or datagram)
 - one process opens a descriptor by calling any of the unix function that returns a descriptor
 - the sending process build a **msghdr** structure containing the descriptor to be passed
 - Sending process sends ancillary data using **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  cmsgh_len; /* length in bytes, including this structure */  
3     int        cmsgh_level; /* originating protocol */  
4     int        cmsgh_type; /* protocol-specific type */  
5     /* followed by unsigned char cmsgh_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 *mhdrp);
4  //Returns: pointer to first cmsghdr structure or NULL if no ancillary data
5  struct cmsghdr *CMSG_NXTHDR(struct msghdr *mhdrp, struct cmsghdr *cmsgp);
6  //Returns: pointer to next cmsghdr structure or NULL if no more ancillary data
7  unsigned char *CMSG_DATA(struct cmsghdr *cmsgp);
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



```
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_CREDS 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.

Sender

innovate

achieve

lead

```
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

innovate

achieve

lead

```
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");
```



BITS Pilani
Pilani Campus



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

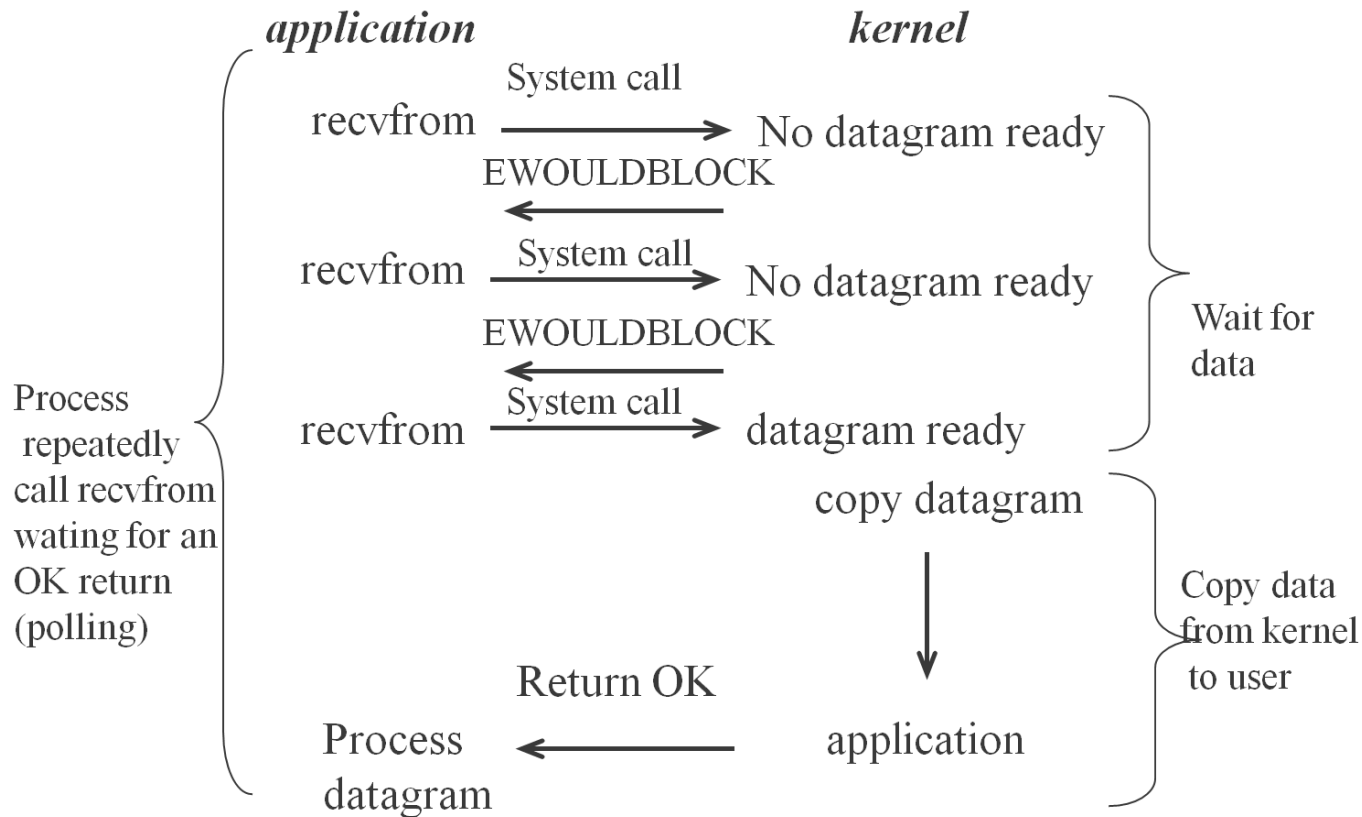


- Most prevalent model

Nonblocking I/O Model



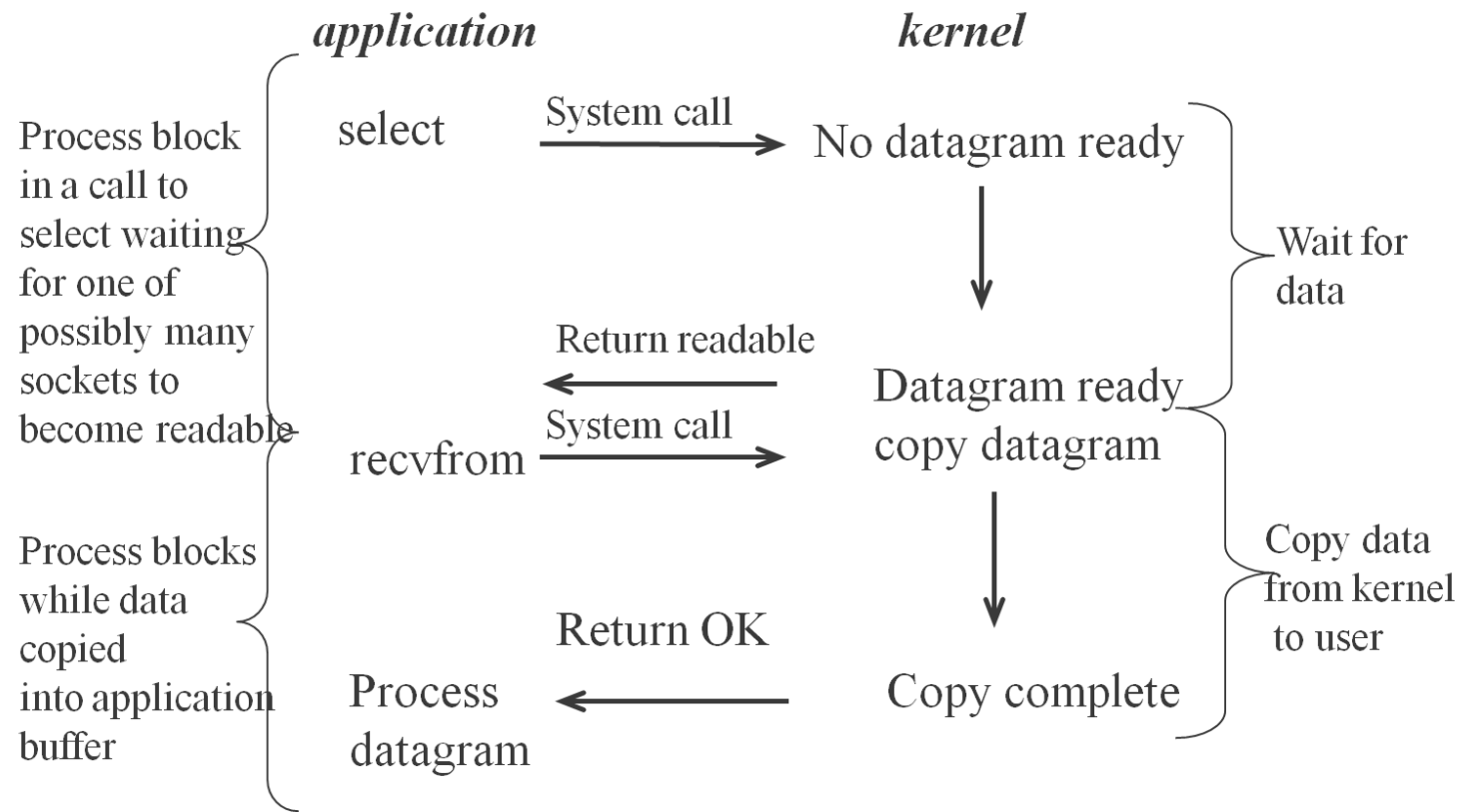
- When the socket is set to be non-blocking,
 - We tell the kernel that do not put the process to sleep if IO can't be completed.



I/O Multiplexing Model



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



I/O Multiplexing Model

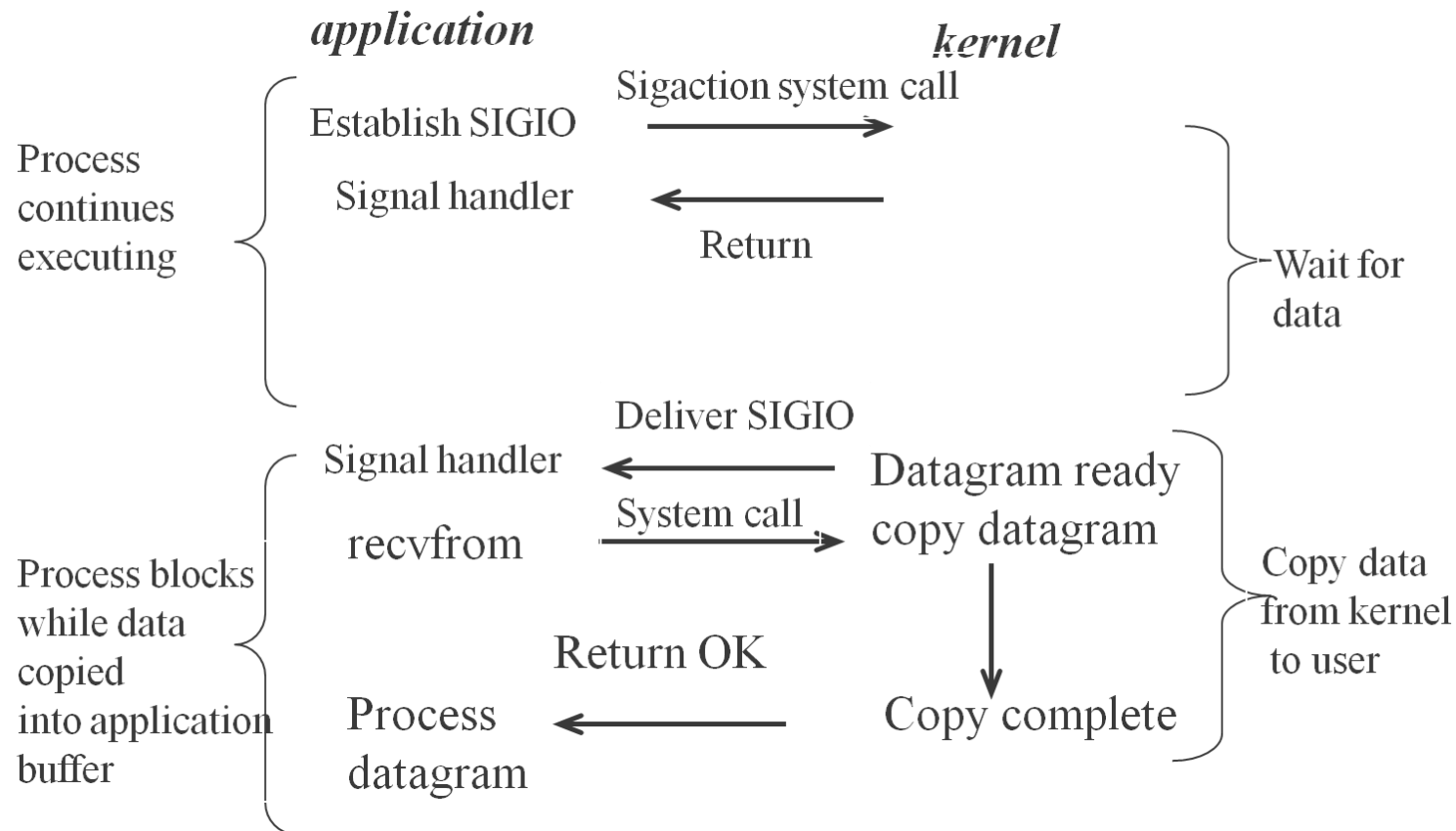


- Blocking and I/O multiplexing seem to be non-different and actually calling two sys calls in I/O multiplexing.
- Advantage with I/O multiplexing is that it can wait for I/O on multiple fds.

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

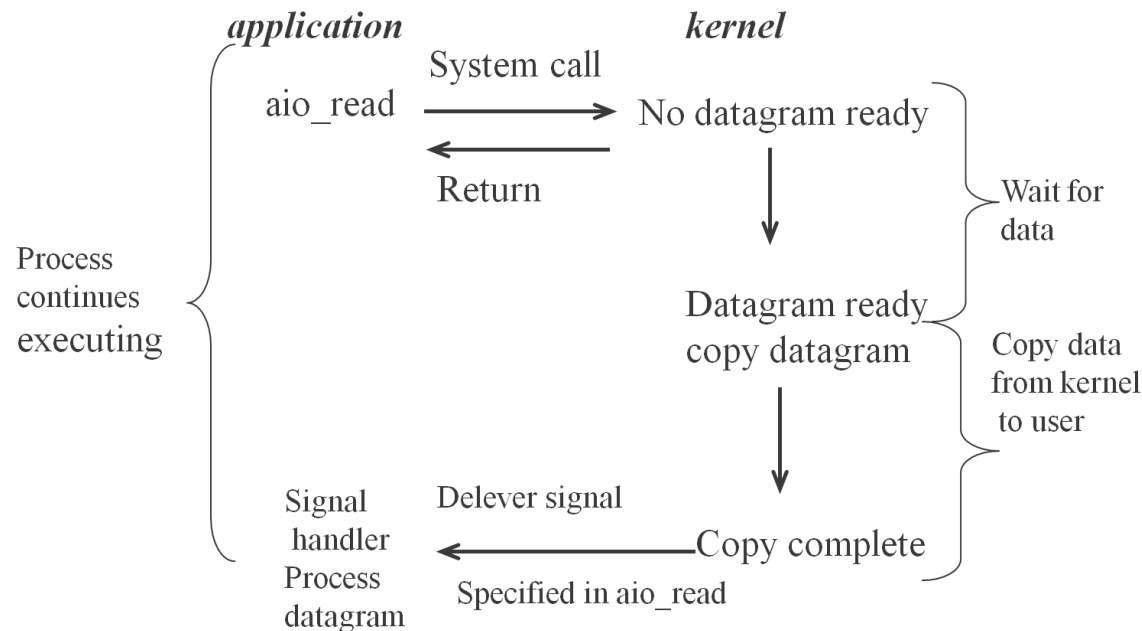


- To use signal-driven I/O with a socket (SIGIO) requires the process to perform the following three steps:
 - A signal handler must be established for the SIGIO signal.
 - The socket owner must be set, normally with the F_SETOWN command of fcntl.
 - Signal-driven I/O must be enabled for the socket,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

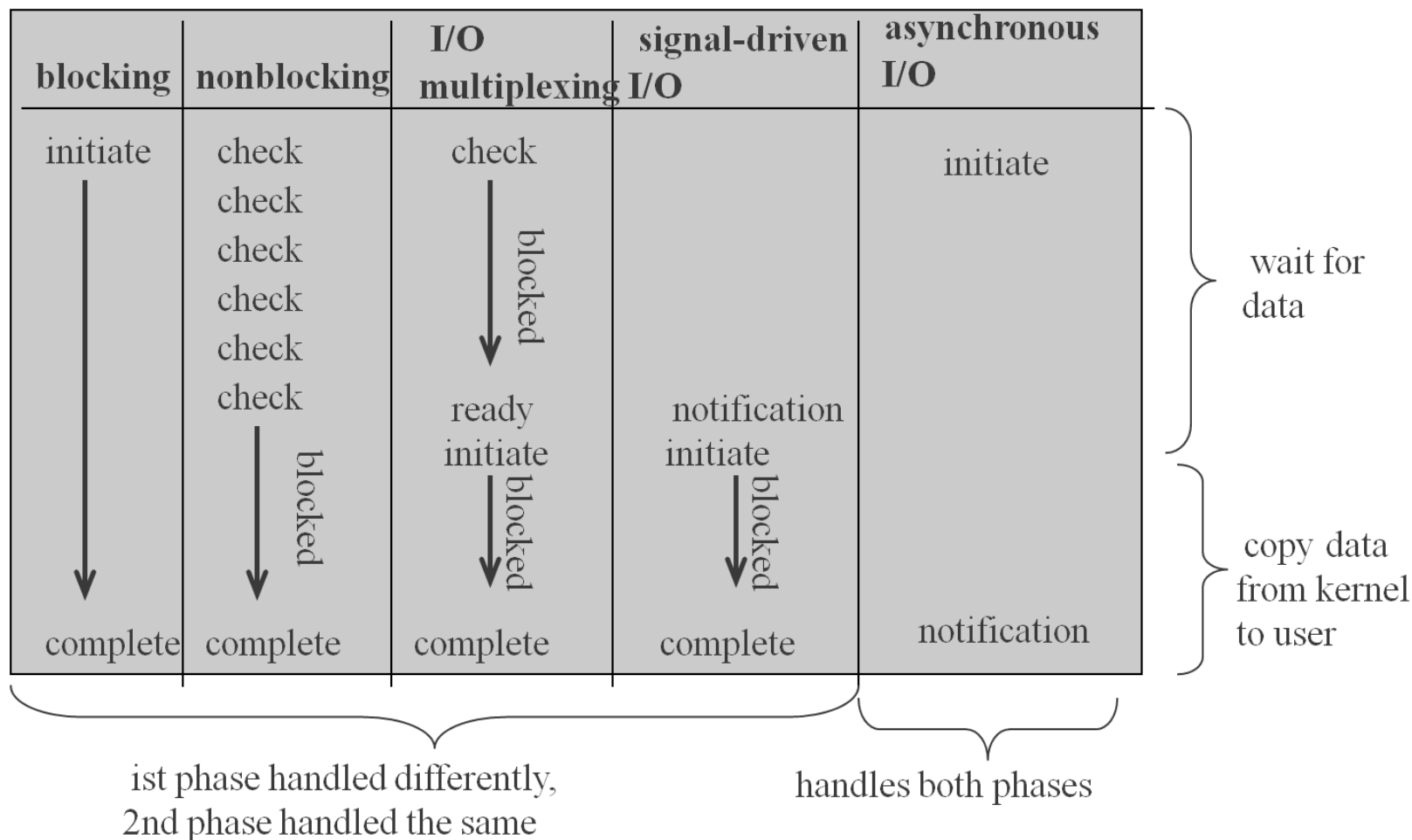
```
1 struct aioctx{
2     int          aio_fdses //    file descriptor
3     off_t        aio_offset //    file offset
4     volatile void* aio_buf //    location of buffer
5     size_t       aio_nbytes //    length of transfer
6     int          aio_reqprio //    request priority offset
7     struct 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 *aioctxp);
```

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

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.

- Readset
 - descriptor for checking readable
- Writerset
 - descriptor for checking writable
- exceptset
 - descriptor 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?



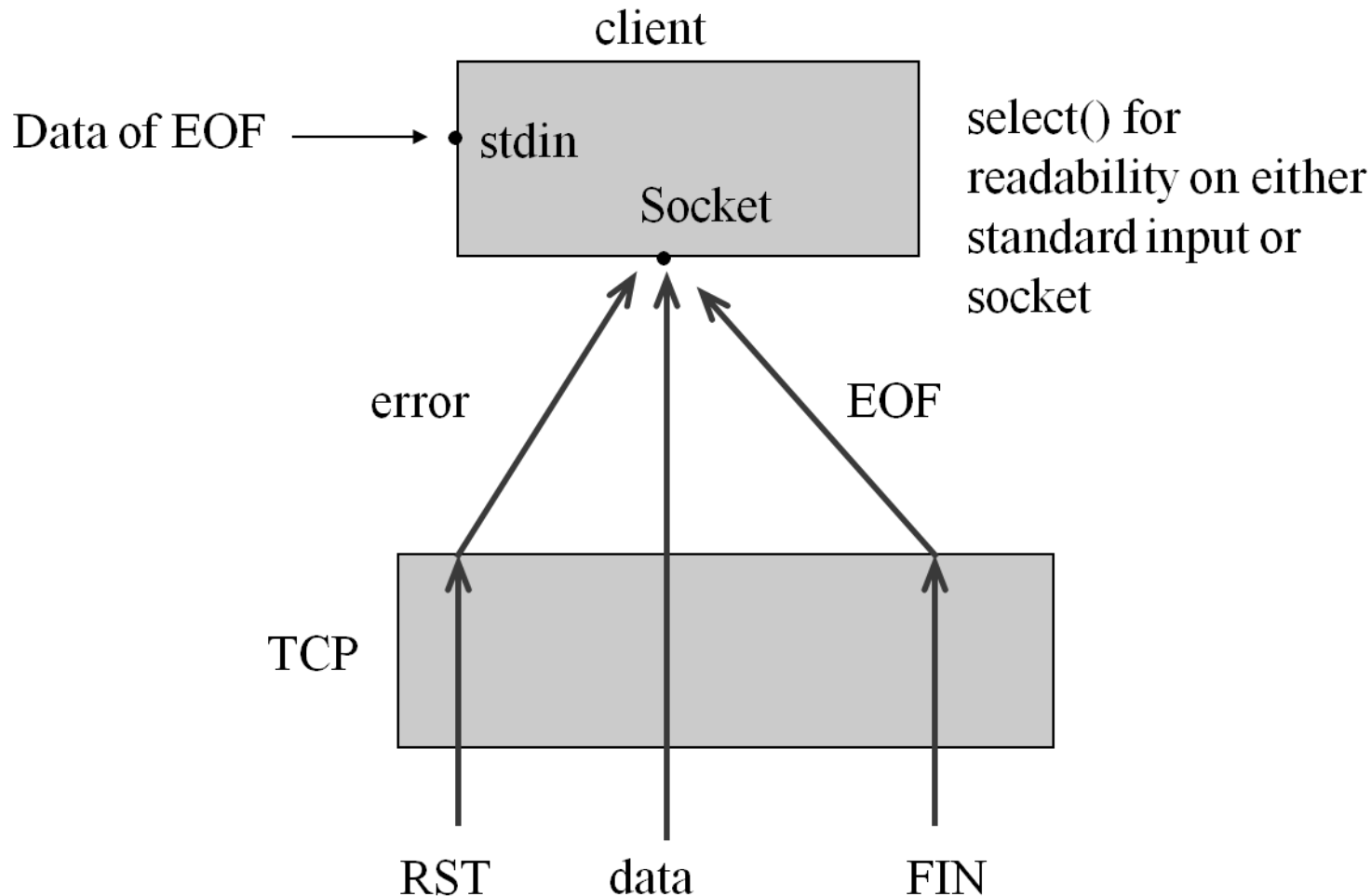
- The number of bytes of data in the socket receive buffer is greater than or equal to the current size of the low-water mark for the socket receive buffer. `SO_RCVLOWAT` socket option. It defaults to 1 for TCP and UDP sockets
- The read half of the connection is closed (i.e., a TCP connection that has 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.

Condition handled by select in a client



Three conditions are handled with the socket



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

```
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
```

Acknowledgements



Q&A





BITS Pilani
Pilani Campus



Thank You