# Bachelor of Computer Application

## BCA-E7
## Network Programming

**Uttar Pradesh Rajarshi Tandon Open University**

**Uttar Pradesh Rajarshi Tandon
Open University**

# Bachelor of Computer Application

## BCA-E7
### Network Programming

**Block**

# 1

# Course Design Committee

| | |
|---|---|
| **Dr. Ashutosh Gupta** | Chairman |
| Director (In-charge) | |
| School of Computer and Information Science, UPRTOU Prayagraj | |
| **Prof. R. S. Yadav** | Member |
| Department of Computer Science and Engineering | |
| MNNIT-Allahabad, Prayagraj | |
| **Ms Marisha** | Member |
| Assistant Professor (Computer Science), | |
| School of Science UPRTOU Prayagraj | |
| **Mr. Manoj Kumar Balwant** | Member |
| Assistant Professor, (Computer Science) | |
| School of Sciences UPRTOU Prayagraj | |

# Course Preparation Committee

| | |
|---|---|
| **Dr. Prabhat Kumar** | Author (Block 1,2) |
| Assistant Professor, Department of IT | |
| NIT Patna | |
| **Dr. Prabhat Ranjan** | Author (Block 3,4) |
| Assistant Professor, Department of Computer Science | |
| Central University of South Bihar | |
| **Dr. Rajiv Mishra** | Editor |
| Associate Professor, Department of CSE | |
| IIT Patna | |
| **Dr. Ashutosh Gupta** (Director in Charge) | |
| School of Computer & Information Sciences, | |
| UPRTOU Prayagraj | |
| **Mr. Manoj Kumar Balwant** | Coordinator |
| Assistant Professor, (Computer Science) | |
| School of sciences UPRTOU Prayagraj | |

[TESLA-010]
BCA-E7/4

# BLOCK INTRODUCTION

**Unit 1:** This unit deals with introduction to network programming. It contains introduction to OSI model, UNIX standards. This unit explains TCP and UDP. This unit tells how to establish the connection and termination of connection in TCP. In this unit, you will learn about buffer sizes and its limitations and standard internet services. The protocol usages by common internet application is described in this unit.

**Unit 2:** This unit deals with elementary sockets. In this unit, you will learn about address structure, value-result arguments, byte ordering and manipulation functions and related functions.

**Unit 3:** This unit deals with elementary TCP sockets. This unit discusses about socket, connect, bind, listen, accept, fork and close functions. In this unit, you will learn about concurrent servers.

**Unit 4:** This unit deals with TCP client/server. This unit tells about TCP Echo server function, Normal start-up. In this unit, you will learn about signal handling server process termination, crashing and rebooting of server host and shutdown of server host.

# UNIT-1 : INTRODUCTION TO NETWORK PROGRAMMING

## Structure

## 1.0 INTRODUCTION

In t his uni t, t he f ocus i s t o pr ovide a ba sic und erstanding of t he technical d esign an d ar chitecture o f t he Internet u sing t wo d ifferent models OSI and TCP/IP. The background of Unix standards, IEEE POSIX and T he O pen G roup's T echnical S tandard de signation t hat w ere l ater converged into The Single Unix Specification Version is discussed.

Most c lient/server a pplications u se either T CP o r U DP as th eir tr ansport layer f or which T CP c onnection e stablishment a nd t ermination a re discussed i n de tail a long w ith de scription of Ipv4, Ipv6 a nd t heir buffer size limitations.

We co ver va rious t opics i n t his uni t that f all in to th is c ategory: TCP's three-way ha ndshake, T CP's c onnection t ermination s equence, pl us T CP, and UDP buffering by the socket layer, and so on.

## 1.1 OBJECTIVES

After the end of this unit, you should be able to:

- Understand t he OSI Model a nd i ts va rious c ommunication processes

- Gain insights regarding the various UNIX standardization schemes

- Differentiate between the TCP and UDP protocols

- Know a bout va rious bu ffer s izes,  standard i nternet s ervices an d popular protocols.

## 1.2  OSI MODEL

OSI Model is an abstract model used to understand a wide range of network a rchitecture.  It was pr oposed a s a  general a pproach t o ne twork models t  o  standardize t  he  communication f  unctions of  a telecommunication or computing system.

The O SI m odel ha s s even l ayers;  Starting a t  the b ottom ( nearest t he physical c onnections), t he l ayers a re: (1 ) P hysical, (2 ) D ata  Link,  (3) Network, (4) Transport, (5) Session, (6) Presentation, and (7) Application.

In the OSI model, control is passed from one layer to the next, starting at the a pplication l ayer i n one  s tation, a nd pr oceeding t o t he bot tom l ayer, over the channel to the next station and back up the hierarchy.

We w ill lo ok a t e ach la yer in  th e O SI mo del  in tu rn, s tarting w ith th e Physical layer. Figure 1.1 shows the layer architecture of OSI model.

| 7 | APPLICATION LAYER |
|---|---|
| 6 | PRESENTATION LAYER |
| 5 | SESSION LAYER |
| 4 | TRANSPORT LAYER |
| 3 | NETWORK LAYER |
| 2 | DATA LINK LAYER |
| 1 | PHYSICAL LAYER |

Figure 1.1 OSI model

- **PHYSICAL LAYER**

    The O SI P hysical la yer d eals w ith th e p hysical a ttributes o f th e actual wired, wireless, fiber optic, or other connection that is used to transport data across a single link. It deals with transmission and reception o f t he uns tructured r aw bi t s tream- electrical i mpulse, light or radio signal over a physical medium.

    It provides the hardware means of sending and receiving data on a carrier, i ncluding d efining  cables, ca rds an d p hysical as pects.  It also doe s B it encoding for faster data transmission. Fast E thernet,

RS232, and ATM are protocols with physical layer components.

- ## DATA LINK LAYER

  The data link layer provides error-free transfer of data frames from one node to another over the physical layer, allowing layers above it to have error-free transmission over the link. Following are some of the functions of Data link layer:

  1. Logical Link establishment between nodes.

  2. Controls Frame traffic by telling transmitting node to "back-off" when no frame buffers are available.

  3. Sequential transfer of frames by defining special sequences to indicate the beginning and end of each packet.

  4. Media access management by defining addresses to the stations communicating.

  5. Frame error checking using Checksum.

  6. Frame acknowledgment and retransmission of non-acknowledged frames and handling duplicate frame receipt.

- ## NETWORK LAYER

  The network layer governs how routers forward packets across multiple hops to get from their source to their destination. It deals with assigning global "routable" addresses to the various systems connected to the network.

  It is also responsible for subnet traffic control instructing a sending station to "throttle back" its frame transmission when the router's buffer fills up. Other functions include:

  1. Frame fragmentation and reassembly at destination station if router's Maximum transmission unit (MTU) is less than frame size.

  2. Logical-physical address mapping that involves translating logical addresses, or names, into physical addresses.

  3. Keeps track of frames forwarded by subnet intermediate systems, to produce billing information.

- ## TRANSPORT LAYER

  The transport layer ensures that messages are delivered error-free, in sequence, and with no losses or duplications. It manages packet loss and retransmission as well as flow control and window size. Services of Transport layer depend upon the services offered by the Network layer. Some of the functions of transport layer are:

  1. Session multiplexing: Multiplex several message streams, or

sessions onto one logical link.

2. End to end message delivery with acknowledgements.

3. Message s egmentation: acc ept m essage f rom t he ( session) layer a bove it, s plits th e me ssage in to s maller units (if n ot already small enough), and passes the smaller units down to the network layer.

- **SESSION LAYER**

  The O SI S ession l ayer ha ndles e stablishing c onnections be tween processes r unning on d ifferent s tations. It a llows tw o a pplication processes o n di fferent m achines t o establish, us e a nd t erminate a connection, c alled a s ession. It a lso pr ovides s upport t o t he sessions established.

- **PRESENTATION LAYER**

  The p resentation l ayer formats t he d ata t o be pr esented t o t he application layer. Like a translator, it translates data from a format used by the application layer into a common format at the sending station, then translate the common format to a format known to the application layer at the receiving station. The key functions of the presentation layer are:

  Data c ompression, D ata e ncryption/decryption, C haracter code translation: for example, ASCII to EBCDIC etc.

- **APPLICATION LAYER**

  The a pplication l ayer a llows us ers a nd a pplication pr ocesses t o access network services. The layer is responsible for functions like Electronic m essaging ( such as m ail), R emote p rinter acces s, Remote f ile acc ess, N etwork m anagement, Inter-process communication, Directory services etc.

  The applications can be Client applications that initiate connection or s erver a pplications t hat r espond t o i ncoming c onnection a nd serve them.

## 1.3   UNIX STANDARDS

The most interesting Unix standardization activity was being done by The Austin Common Standards Revision Group (CSRG) that produced roughly 4,000 pa ges of s pecifications t hat c arry both t he IEEE P OSIX designation as well as The Open Group's Technical Standard designation, thus leading t o m ultiple names t o s ame s tandards, for e xample, ISO/IEC 9945:2002, IEEE S td 1003.1 -2001, a nd t he S ingle U nix S pecification Version 3 are various names of same standard, The POSIX Specification.

**Background on POSIX**

POSIX is an acronym for Portable Operating System Interface. POSIX is

not a single standard, but a set of standards being developed by the Institute for Electrical and Electronics Engineers, Inc., normally called the IEEE. The POSIX standards have also been adopted as international standards by ISO and the International Electrotechnical Commission (IEC), called ISO/IEC.

The interesting history of POSIX standards has been covered only briefly here:

- **IEEE Std 1003.1–1988** (317 pages) was the first POSIX standard. It specified the C language interface into a Unix-like kernel and covered the following areas: process primitives (*fork, exec*, signals, and timers), the environment of a process (user IDs and process groups), files and directories (all the I/O functions), terminal I/O, system databases (password file and group file), and the *tar* and *cpio* archive formats.

The first POSIX standard was a trial-use version in 1986 known as "IEEE-IX." The name "POSIX" was suggested by Richard Stallman.

- **IEEE Std 1003.1–1990** (356 pages) was next, and it was also known as ISO/IEC 9945–1: 1990. Minimal changes were made from the 1988 to the 1990 version. Appended to the title was "Part 1: System Application Program Interface (API) [C Language]," indicating that this standard was the C language API.

- **IEEE Std 1003.2–1992** came next in two volumes (about 1,300 pages). Its title contained "Part 2: Shell and Utilities." This part defined the shell (based on the System V Bourne shell) and about 100 utilities (programs normally executed from a shell, from *awk* and basename to *vi* and *yacc*). Throughout this text, we will refer to this standard as POSIX.2.

- **IEEE Std 1003.1b–1993** (590 pages) was originally known as IEEE P1003.4. This was an update to the 1003.1–1990 standard to include the real-time extensions developed by the P1003.4 working group. The 1003.1b–1993 standard added the following items to the 1990 standard: file synchronization, asynchronous I/O, semaphores, memory management (mmap and shared memory), execution scheduling, clocks and timers, and message queues.

- **IEEE Std 1003.1, 1996** Edition [IEEE 1996] (743 pages) came next and included 1003.1–1990 (the base API), 1003.1b–1993 (real-time extensions), 1003.1c–1995 (pthreads), and 1003.1i–1995 (technical corrections to 1003.1b). This standard was also called ISO/IEC 9945–1: 1996. Three units on threads were added, along with additional sections on thread synchronization (mutexes and condition variables), thread scheduling, and synchronization scheduling. Throughout this text, we will refer to this standard as POSIX.1.

  This standard also contains a Foreword stating that ISO/IEC 9945

consists of the following parts:

Part 1: System API (C language)

Part 2: Shell and utilities

Part 3: System administration (under development) Parts 1 a nd 2 are what we call POSIX.1 and POSIX.2

- **IEEE Std 1003.1g:** Protocol-independent interfaces (PII) became an approved standard in 2000. Until the introduction of The Single Unix Specification Version 3, t his P OSIX w ork w as t he m ost relevant to the topics covered in this book. T his is the networking API s tandard a nd it d efines tw o APIs, w hich it c alls D etailed Network Interfaces (DNIs): 1. D NI/Socket, ba sed on t he 4.4 BSD sockets API 2. DNI/XTI, based on the X/Open XPG4 specification Work on t his s tandard s tarted i n t he l ate 1980s as t he P 1003.12 working group (later renamed P1003.1g). Throughout this text, we will refer to this standard as POSIX.1g.

## Background on The Open Group

The Open Group was formed in 1996 by the consolidation of the X/Open Company ( founded i n 1 984) a nd t he O pen S oftware Foundation ( OSF, founded i n 1988). It i s a n i nternational c onsortium of ve ndors a nd e nd-user customers from industry, government, and academia. Here is a brief background on the standards they produced:

- X/Open published the X/Open Portability Guide, Issue 3 ( XPG3) in 1989.

- Issue 4 w as published i n 1992, f ollowed b y Issue 4, V ersion 2 i n 1994. This latest version was also known as "Spec 1170," with the magic num ber 1,170 b eing t he s um of t he num ber of s ystem interfaces ( 926), t he nu mber of he aders (70), a nd t he num ber of commands ( 174). T he l atest na me for this s et of s pecifications i s the "X/Open Single Unix Specification," although it i s also called "Unix 95."

- In M arch 1997, V ersion 2 of t he Single U nix S pecification w as announced. P roducts c onforming t o t his s pecification w ere c alled "Unix 9 8." We w ill r efer to th is s pecification a s ju st "UNIX 9 8" throughout this text. The number of interfaces required by Unix 98 increases f rom 1,170 t o 1,434, a lthough f or a workstation t his jumps t o 3,030, be cause i t i ncludes t he C ommon D esktop Environment (CDE), which in turn requires the X Window System and the M otif user interface. Details are available in [Josey 1997] and a t ht tp://www.UNIX.org/version2. T he n etworking s ervices that are part of U nix 98 a re de fined for both the sockets and X TI APIs. This specification is nearly identical to POSIX.1g.

## Unification of Standards

Now, Most Unix systems today conform to some version of POSIX.1 and POSIX.2; m any comply w ith T he S ingle U nix S pecification V ersion 3. The focus of this book is on The Single Unix Specification Version 3, with our main focus on the sockets API.

## Internet Engineering Task Force (IETF)

The Internet Engineering Task Force (IETF) is a large, open, international community of ne twork de signers, ope rators, ve ndors, a nd r esearchers concerned w ith t he evolution of t he Internet a rchitecture a nd t he s mooth operation of t he Internet. It i s ope n t o any i nterested i ndividual. T he Internet s tandards pr ocess i s doc umented i n R FC 2026 [Bradner 1996]. Internet s tandards nor mally de al w ith pr otocol i ssues a nd not w ith programming APIs.

Nevertheless, two RFCs (RFC 3493 [Gilligan et al. 2003] and RFC 3542 [Stevens e t a l. 2003] ) s pecify t he s ockets A PI for IPv6. T hese are informational R FCs, no t s tandards, and w ere produced t o s peed t he deployment of portable applications by the numerous vendors working on early releases of IPv6. Although standards bodies tend to take a long time, many APIs were standardized in The Single Unix Specification Version 3.

---

### Check your progress

1. What are the responsibilities of network layer and transport layer?

2. Explain the connection establishment phase of the TCP protocol.

---

## 1.4 TCP AND UDP & TCP CONNECTION ESTABLISHMENT AND FORMAT

This s ection f ocuses on t he following transport l ayer protocols: TCP and UDP.

Most c lient/server applications us e e ither T CP or U DP. A nother pr otocol SCTP i s a ne wer pr otocol, or iginally de signed f or t ransport of t elephony signalling across the Internet. These t ransport pr otocols us e the ne twork-layer pr otocol IP, e ither IPv4 or Ipv6. It i s pos sible for an application to bypass the transport layer and us e IPv4 or IPv6 directly. This i s c alled a raw socket.

UDP is a s imple, u nreliable d atagram p rotocol, w hile T CP is a sophisticated, r eliable b yte s tream pr otocol. Let us l ook i nto bot h t he protocols in detail.

## User Datagram Protocol (UDP)

UDP i s a c onnectionless pr otocol, a nd U DP s ockets a re an e xample of

datagram sockets. There is no guarantee that UDP datagrams ever reach their intended destination. The application sends message to a UDP socket, encapsulated in a UDP datagram, which is then further encapsulated as an IP datagram, which is then sent to its destination. For a UDP datagram reaching its final destination, that order will be preserved across the network, or that datagrams arrive only once is not guaranteed.

Lack of reliability is the drawback we have with network programming with UDP. If a UDP Datagram does not reach its destination or is dropped midway, there is no scope of automatic retransmission.

Each UDP datagram has a length. The length of a datagram is passed to the receiving application along with the data. Being a connectionless service, there need not be any long-term relationship between a UDP client and server. For example, a UDP client can create a socket and send a datagram to a given server and then immediately send another datagram on the same socket to a different server. Similarly, a UDP server can receive several datagrams on a single UDP socket, each from a different client.

## Transmission Control Protocol (TCP)

TCP is described in RFC 793 [Postel 1981c], and updated by RFC 1323 [Jacobson, Braden, and Borman 1992], RFC 2581 [Allman, Paxson, and Stevens 1999], RFC 2988 [Paxson and Allman 2000], and RFC 3390 [Allman, Floyd, and Partridge 2002].

TCP is a connection oriented protocol and provides connections between clients and servers. A TCP client establishes a connection with a given server, exchanges data with that server across the connection, and then terminates the connection.

It provides **reliability** by using acknowledgement in return and if not received retransmitting the data and waiting for a longer duration of time. It does not provide the guarantee to deliver data at the destination. Just delivering data if it can be delivered to a notification to end user if data cannot be sent.

The waiting time for acknowledgement or Roundtrip time(RTT) between Client and server is estimated by the algorithms in TCP.

TCP also sequences the data by associating a sequence number with every byte that it sends. For example, assume an application writes 2,048 bytes to a TCP socket, causing TCP to send two segments, the first containing the data with sequence numbers 1–1,024 and the second containing the data with sequence numbers 1,025–2,048. (A segment is the unit of data that TCP passes to IP.) If the segments arrive out of order, the receiving TCP will reorder the two segments based on their sequence numbers before passing the data to the receiving application. Thus TCP can detect a duplicate data from the sequencing and can discard it.

TCP provides **flow control.** TCP has the advertised window which tells peer how many bytes of data it can accept. It guarantees that the sender

cannot overflow the receiving buffer. The window changes dynamically over time: As data is received from the sender, the window size decreases, but as the receiving application reads data from the buffer, the window size increases. It is posible for the window to reach 0: when TCP's receive buffer for a socket is full and it must wait for the application to read data from the buffer before it can take any more data from the peer.

Finally, a TCP connection is **full-duplex**. This means that an application can send and receive data in both directions on a given connection at any time. This means that TCP must keep track of state information such as sequence numbers and window sizes for each direction of data flow: sending and receiving.

## TCP Connection Establishment and Termination

Let us understand how TCP connections are established and terminated, and TCP's state transition diagram.

### Three-Way Handshake

Figure 1.2 shows the connection establishment of TCP by three-way handshaking.

1. Host A sends a connection request to host B by setting the SYN (a *synchronize* message, used to initiate and establish a connection) bit. Host A also registers its initial sequence number to use (Seq_no fl x).

2. Host B acknowledges the request by setting the ACK (an *acknowledgment*) bit and indicating the next data byte to receive (Ack_no fl x + 1). The "plus one" is needed because the SYN bit consumes one sequence number. At the same time, host B also sends a request by setting the SYN bit and registering its initial sequence number to use (Seq_no fl y).

3. Host A acknowledges the request from B by setting the ACK bit and confirming the next data byte to receive (Ack_no fl y + 1). Note that the sequence number is set to x + 1. On receipt at B the connection is established.

   If during a connection establishment phase, one of the hosts decides to refuse a connection request, it will send a reset segment by setting the RST bit. Each SYN message can specify options such as maximum segment size, window scal- ing, and timestamps. Because TCP segments can be delayed, lost, and duplicated, the initial sequence number should be different each time a host requests a connection.
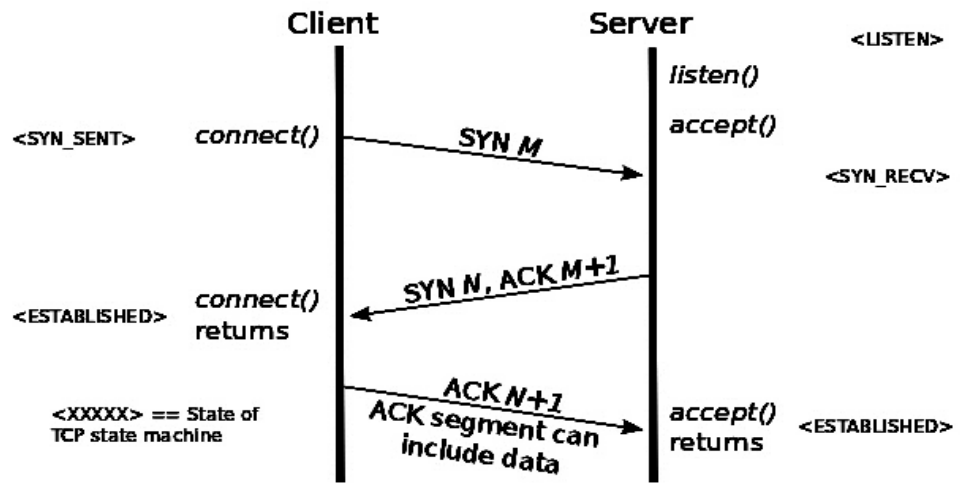
Figure 1.2: TCP Three Way Handshake

## TCP Connection Termination

Figure 1.3 shows the TCP connection termination.

1. One application calls close first, and we say that this end performs the a ctive c lose. T his e nd's T CP s ends a FIN s egment, w hich means it is finished sending data.

2. The ot her e nd t hat r eceives t he F IN p erforms t he p assive cl ose. The received FIN is acknowledged by TCP. The receipt of the FIN is al so p assed t o t he ap plication as an en d of-file ( after an y d ata that may have already been queued for the application to receive), since the receipt of the FIN means the application will not receive any additional data on the connection.

3. Sometime l ater, t he ap plication t hat r eceived t he en d-of-file w ill close its socket. This causes its TCP to send a FIN.

4. The T CP o n t he s ystem t hat r eceives this final FIN (the e nd th at did the active close) acknowledges the FIN.

Figure 1.3: TCP connection termination

# TCP State Transition Diagram

The ope ration of T CP w ith r egard t o c onnection e stablishment a nd connection termination can be specified with a state transition diagram as shown in Figure 1.4.



**Figure 1.4:** State transition diagram

A c onnection p rogresses t hrough a s eries of s tates dur ing i ts l ifetime a nd transition f rom s tate to state is b ased o n th at c urrent s tate a nd segment received in that state.

The s tates are: LISTEN, S YN-SENT, S YNRECEIVED, ES TABLISHED, FIN-WAIT-1, F IN-WAIT-2, C LOSE-WAIT, C LOSING, LAST-ACK, TIME-WAIT, and the fictional state CLOSED.

LISTEN represents waiting for a connection request from any remote TCP and port.

SYN-SENT r epresents waiting f or a m atching connection r equest a fter having sent a connection request. SYN-

RECEIVED represents w aiting f or a confirming c onnection r equest acknowledgment after having both received and sent a connection request.

ESTABLISHED represents a n ope n c onnection, da ta r eceived c an be delivered t o t he us er. T he nor mal s tate f or t he da ta t ransfer pha se of t he connection.

FIN-WAIT-1 represents waiting for a connection termination request from the r emote T CP, or a n acknowledgment of t he c onnection t ermination request previously sent.

FIN-WAIT-2 represents waiting for a connection termination request from the remote TCP.

CLOSE-WAIT r epresents w aiting f or a c onnection t ermination r equest from the local user.

CLOSING r epresents waiting f or a connection t ermination r equest acknowledgment from the remote TCP.

LAST-ACK r epresents w aiting for an a cknowledgment of t he connection termination request previously sent to the remote TCP (which includes an acknowledgment of its connection termination request).

TIME-WAIT r epresents waiting f or e nough t ime t o pa ss t o be s ure t he remote T CP r eceived t he ack nowledgment o f i ts co nnection t ermination request.

CLOSED represents no connection state at all.

## 1.5   BUFFER SIZES AND LIMITATION

The buf fer sizes o f IP Datagrams h ave certain limitations th at a ffect the data an application can transmit. The limitations are as follows:

- IPv4 datagram has a maximum size of 65,535 b ytes, including the IPv4 header. Its 16 bit total length field includes header size.

- IPv6 datagram has a maximum size of 65,575 b ytes, including the 40-byte IPv6 he ader. Its 16 bi t t otal l ength field doe s not i nclude header s ize. O n d atalinks w ith a ma ximum tr ansmission unit (MTU) t hat e xceeds 65,535, IPv6 c an ha ve e xtended pa yload length field of 32 bits.

- MTU i s de pendent on Hardware, f or e xample E thernet M TU i s 1500 bytes whereas Point to Point Protocol has configurable MTU. Minimum link MTU for Ipv4 is 68 b ytes which means, Maximum sized he ader ( 20 b ytes of f ixed he ader, 40 b ytes of options) + minimum s ized fragment ( 8 b ytes) c an be p assed. H owever, Minimum link MTU for Ipv6 is 1,280 bytes.

- The smallest MTU in the path between two hosts is called the **Path MTU**. F or e xample, the E thernet M TU of 1,500 b ytes i s the path

MTU. MTU between two hosts is different in both directions.

- IP Datagrams with size exceeding Link MTU are fragmented at the outgoing interface and reassembled at destination by both IPv4 and Ipv6. IPv4 hosts perform fragmentation on datagrams that they generate and IPv4 routers perform fragmentation on datagrams that they forward. For Ipv6, fragmentation of datagrams is performed only at Ipv6 hosts and not at Ipv6 routers with an exception of routers that generate their own datagrams instead of forwarding.

- Fragmentation fields are included in Ipv4 headers but not in Ipv6 headers. If "don't fragment" (DF) bit is set, it specifies that this datagram must not be fragmented, either by the sending host or by any router. A router that receives an IPv4 datagram with the DF bit set whose size exceeds the outgoing link's MTU generates an ICMPv4 "destination unreachable, fragmentation needed but DF bit set" error message. DF bit is implied with Ipv6 datagrams, so if a Ipv6 router receives a datagram whose size exceeds the outgoing link's MTU, it generates an ICMPv6 "packet too big" error message. This DF bit can also be used to discover the path MTU.

- TCP has a maximum segment size (MSS) that announces to the peer TCP the maximum amount of TCP data that the peer can send per segment. It tells the peer the actual value of the reassembly buffer size tries to avoid fragmentation. The MSS is often set to the interface MTU minus the fixed sizes of the IP and TCP headers.

- On an Ethernet using IPv4, MSS would be 1,460, and on an Ethernet using IPv6, this would be 1,440. (The TCP header is 20 bytes for both, but the IPv4 header is 20 bytes and the IPv6 header is 40 bytes.) The MSS value in the TCP MSS option is a 16-bit field, limiting the value to 65,535. This is fine for IPv4, since the maximum amount of TCP data in an IPv4 datagram is 65,495 (65,535 minus the 20-byte IPv4 header and minus the 20-byte TCP header).

## TCP Output

When an application calls write, the kernel copies all the data from the application buffer into the TCP socket send buffer and returns only when the final byte in the application buffer has been copied into the socket send buffer. Insufficient room in the socket send buffer due to the larger application buffer size or socket send buffer already full, blocks the socket and process is put to sleep.

TCP transmits the data from buffer to peer TCP according to rules of data transmission and discards the data from the buffer only after receiving ACKs from the peer. Data is sent in MSS-sized chunks (announced by peer TCP) from TCP to IP with its header prepended to it. IP then prepends its header to the datagram, searches the appropriate routing table for destination IP and sends the datagram to proper datalink. IP performs fragmentation in case path MTU discovery (in newer implementations) not used or MSS option not used. The output queue associated with each

datalink discards the packet and reports an error to TCP via IP in case it is full.

**UDP Output**

UDP s ocket h as a s end b uffer s ize ( which w e can change w ith t he SO_SNDBUF s ocket opt ion), but t his i s s imply a n uppe r l imit on the maximum-sized U DP d atagram t hat can b e w ritten t o t he s ocket. I f an application w rites a d atagram l arger t han t he s ocket s end b uffer s ize, EMSGSIZE is returned. Since UDP is unreliable, it does not need to keep a copy of the application's data and does not need an actual send buffer.

UDP s imply pr epends i ts 8 -byte h eader an d p asses t he d atagram t o IP. IPv4 or IPv6 pr epends i ts he ader, d etermines t he out going i nterface b y performing the routing function, and then either adds the datagram to the datalink o utput q ueue ( if it f its w ithin th e MTU) o r f ragments th e datagram and adds each fragment to the datalink output queue. If there is no room on the queue for the datagram or one of its fragments, ENOBUFS is often returned to the application.

## 1.6   STANDARD INTERNET SERVICES

Some of t he s tandard s ervices t hat a re pr ovided b y m ost implementations of TCP/IP are following:

| Name | TCP Port | UDP Port | RFC | Description |
|------|----------|----------|-----|-------------|
| Echo | 7 | 7 | 862 | Server returns whatever the client sends. |
| Discard | 9 | 9 | 863 | Server discards whatever the client sends. |
| Daytime | 13 | 13 | 867 | Server returns the time and date in human-readable format. |
| Chargen | 19 | 19 | 864 | TCP server sends a continual stream of characters, until the connection is terminated by the client. UDP server sends a datagram containing a random character (between 0 and 512) each time the client sends a datagram. |
| Time | 37 | 37 | 868 | Server returns the time as a 32-bit binary number. This number represents the number of seconds since midnight January 1, 1990, UTC. |

Figure 1.5 Standard TCP/IP services provided by most implementations.

If w e examine t he p ort n umbers f or t hese s tandard s ervices an d o ther standard TCP/IP services (Telnet, FTP, SMTP, etc.), most are odd numbers. This i s historical a s t hese por t num bers a re de rived f rom t he N CP por t numbers. ( NCP, t he Network C ontrol P rotocol, pr eceded T CP as a transport l ayer pr otocol for the A RPANET.) N CP w as s implex, not full-duplex, s o e ach application r equired t wo c onnections, a nd a n even-odd pair of por t num bers w as r eserved f or each a pplication. W hen T CP a nd UDP became the standard transport layers, only a single port number was needed per application, so the odd port numbers from NCP were used.

Often, *inetd* daemon provides these services on Unix hosts.

---

**Check your progress**

1.    Explain some limitations of buffer sizes of IP datagrams.

2.    Enlist Standard Internet Services provided by TCP/IP.

---

# 1.7   PORT NUMBERS

TCP a nd U DP i dentify a pplications us ing 16 -bit por t num bers called P ort num bers t hat range b etween 1 and 1023. S ervers are represented by their port numbers. For example, a TCP/IP implementation that provides FTP server provides that service TCP port 21, Telnet service is pr ovided on P ort 23, TFTP ( the Trivial F ile Transfer P rotocol) i s on UDP port 69. The well-known ports are managed by the Internet Assigned Numbers Authority (IANA).

A c lient us ually doe sn't c are w hat por t num ber i t us es on i ts e nd. A ll it needs to be certain of is that whatever port number it uses be unique on its host. C lient por t num bers a re c alled e phemeral por ts ( i.e., s hort l ived). This is because a client typically exists only as long as the user running the client needs its service, while servers typically run as long as the host is up.

The w ell-known por t nu mbers a re contained i n the f ile / etc/services o n most Unix systems. To find the port numbers for the Telnet server and the Domain Name System, we can execute

sun % grep telnet /etc/services

telnet 23/tcp                          says it uses TCP port 23

sun % grep domain /etc/services

domain 53/udp                          says it uses UDP port 53

domain 53/tcp                  and TCP port 53.

Port numbers in the range of 1 to 1023 are reserved, and are used by some applications as part of the authentication between the client and server.

# 1.8 PROTOCOL USAGE BY COMMON INTERNET APPLICATIONS

| Application | IP | ICMP | UDP | TCP | SCTP |
|---|---|---|---|---|---|
| ping | | • | | | |
| traceroute | | • | • | | |
| OSPF (routing protocol) | • | | | | |
| RIP (routing protocol) | | | • | | |
| BGP (routing protocol) | | | | • | |
| BOOTP (bootstrap protocol) | | | • | | |
| DHCP (bootstrap protocol) | | | • | | |
| NTP (time protocol) | | | • | | |
| TFTP | | | • | | |
| SNMP (network management) | | | • | | |
| SMTP (electronic mail) | | | | • | |
| Telnet (remote login) | | | | • | |
| SSH (secure remote login) | | | | • | |
| FTP | | | | • | |
| HTTP (the Web) | | | | • | |
| NNTP (network news) | | | | • | |
| LPR (remote printing) | | | | • | |
| DNS | | | • | • | |
| NFS (network filesystem) | | | • | • | |
| Sun RPC | | | • | • | |
| DCE RPC | | | • | • | |
| IUA (ISDN over IP) | | | | | • |
| M2UA,M3UA (SS7 telephony signaling) | | | | | • |
| H.248 (media gateway control) | | | • | • | • |
| H.323 (IP telephony) | | | • | • | • |
| SIP (IP telephony) | | | • | • | • |

Figure 1.6: summarizes the protocol usage of various common Internet applications.

The first two applications, ping and traceroute, are diagnostic applications that use ICMP. traceroute builds its own UDP packets to send and reads ICMP replies. The three popular routing protocols demonstrate the variety of transport protocols used by routing protocols. OSPF uses IP directly, employing a raw socket, while RIP uses UDP and BGP uses TCP. The next five are UDP-based applications, followed by seven TCP applications and four that use both UDP and TCP. The final five are IP telephony applications that use SCTP exclusively or optionally UDP, TCP, or SCTP.

## 1.9 SUMMARY

The unit introduces many of the terms and concepts that shall be expanded on throughout the rest of the book. It also gives an overview of developing protocol-dependent programs.

TCP uses a three-way handshake for establishing connection while a

connection i s t erminated us ing a f our-packet ex change. W hen a T CP connection is established, the connection state is changed from CLOSED to E STABLISHED, a nd upon t ermination, t he s tate i s c hanged t o CLOSED. There are total 11 states in which a TCP connection may reside. A s tate tr ansition d iagram s pecifies th e r ules f or s witching b etween th e states. K nowledge about t he s tate t ransition d iagram i s n ecessary for understanding w hat h appens w hen an a pplication c alls functions s uch as connect, accept, and close.

Unlike T CP, U DP doe sn't e stablish a c onnection be fore s ending da ta, i t just sends. Because of this, UDP is called "Connectionless". UDP packets are o ften called "Datagrams". A n ex ample o f U DP i n act ion i s t he D NS service. DNS servers send and receive DNS requests using UDP.

UDP is a simple, c onnectionless, a nd unr eliable pr otocol, w hile T CP is a complex, c onnection- oriented, a nd r eliable. A lthough m ost of t he applications on t he Internet us e T CP ( the W eb, Telnet, F TP, a nd e mail), there i s a ne ed f or U DP a s w ell. In f urther units, we s hall d iscuss t he reasons to choose UDP instead of TCP.

# 1.10 TERMINAL QUESTIONS

1.  Discuss in detail the layers of OSI model.

2.  Explain TCP/IP layering in detail with neat sketch?

3.  Explain th e T CP s tate tr ansition d iagram with th e h elp o f a diagram

4.  TCP assumes an MSS of 536 if it does not receive an MSS option from the peer. Why is this value used?

5.  UDP i s a s imple, c onnectionless, a nd unr eliable pr otocol, w hile TCP is a c omplex, c onnection- oriented, a nd r eliable. E xplain the statement.

# UNIT 2 : ELEMENTARY SOCKETS

**Structure**

## 2.0    INTRODUCTION

This unit focuses on the description of the sockets API. The socket address structures described here can be passed in two directions: from the process to the kernel, and from the kernel to the process. The latter case is an ex ample o f a v alue-result a rgument, a nd we w ill e ncounter other examples of these arguments throughout the text.

The socket address structure is created by converting a text representation of an address into the binary value using the address conversion functions. Most e xisting IPv4 c ode us es i net_addr a nd inet_ntoa, but t wo n ew functions, inet_pton and inet_ntop, handle both IPv4 and Ipv6.

## 2.1    OBJECTIVES

After the end of this unit, you should be able to:

- Analyse the problem and develop an algorithm for its solution;

- Represent an algorithm in an abstract language (eg. pseudo-code, Structure Diagrams);

- Represent an algorithm with the help of flowchart.

- Understand the fundamental principle of program design.

## 2.2    SOCKET ADDRESS STRUCTURES

Various s tructures a re us ed i n U nix S ocket P rogramming t o hol d information a bout t he address a nd por t, a nd ot her i nformation. M ost socket f unctions r equire a poi nter t o a s ocket address s tructure a s a n

argument. Structures defined in this unit are related to Internet Protocol Family. The names of these structures begin with sockaddr_ and end with a unique suffix for each protocol suite.

The first structure is *sockaddr* that holds the socket information –

```
struct sockaddr {
unsigned short sa_family;
        char sa_data [14];
        };
```

| Attribute | Values | Description |
|---|---|---|
| sa_family | AF_INET<br>AF_UNIX<br>AF_NS<br>AF_IMPLINK | It represents an address family. In most of the Internet-based applications, we use AF_INET. |
| sa_data | Protocol-Specific Address | The content of the 14 bytes of protocol specific address are interpreted according to the type of address. For the Internet family, we will use port number IP address, which is represented by*sockaddr_in* structure defined below. |

This is a generic socket address structure, which will be passed as reference in most of the socket function calls. But any socket function that takes them as pointers must be support socket address structures from any supported protocol families. This generic socket address structure is defined in <sys/socket.h> header.

The below function is an example of function taking pointer to the generic socket address structure.

```
#include <sys/socket.h>
int bind (int sockfd, const struct sockaddr *myaddr, socklen_t addrlen);
Returns: 0 if OK, –1 on error
```

Below example shows how the function is called:

```
struct sockaddr_in serv;              /*IPv4 socket address structure */
/* fill in serv{} */
bind (sockfd, (struct sockaddr *) &serv, sizeof(serv));
```

## IPv4 Socket Address Structure

An IPv4 socket address structure, commonly called an "Internet socket address structure," is named *sockaddr_in* and is defined by including the

<netinet/in.h> header. Both the IPv4 address and the TCP or UDP port number are always stored in the structure in network byte order. The internet (IPv4) socket address structure: *sockaddr_in* has been shown below:

```
struct in_addr {
in_addr_t s_addr; /* 32-bit IPv4 address */
/* network byte ordered */
};
struct sockaddr_in {
uint8_t sin_len; /* length of structure (16) */
sa_family_t sin_family; /* AF_INET */
in_port_t sin_port; /* 16-bit TCP or UDP port number */
/* network byte ordered */
struct in_addr sin_addr; /* 32-bit IPv4 address */
/* network byte ordered */
char sin_zero [8]; /* unused */
};
```

We need not set length field even if it is present unless routing sockets come into picture. Only kernels that deal with socket address structures from various protocol families (e.g., the routing table code) use it. POSIX datatypes are shown for the s_addr, sin_family, and sin_port members in socket address.

The socket functions bind, connect, sendto, and sendmsg, that pass socket address structure all got through an additional sockargs funtion which copies the structure from the process and sets its sin_member to the size of the structure being passed as the argument.

The functions accept, recvfrom, recvmsg, getpeername, and getsockname that pass socket address structure from the kernel to the process too set the sin_len member before returning to the process. The five socket functions that pass a socket address structure from the kernel to the process, accept, recvfrom, recvmsg, getpeername, and getsockname, all set the sin_len member before returning to the process.

The POSIX specification requires only three members in the structure: sin_family, sin_addr, and sin_port. A POSIX-compliant implementation can also define additional structure members, for an Internet socket address structure. Almost all implementations add the sin_zero member so that all socket address structures are at least 16 bytes in size.

| Datatype | Description | Header |
|---|---|---|
| int8_t | Signed 8-bit integer | <sys/types.h> |
| uint8_t | Unsigned 8-bit integer | <sys/types.h> |
| int16_t | Signed 16-bit integer | <sys/types.h> |
| uint16_t | Unsigned 16-bit integer | <sys/types.h> |
| int32_t | Signed 32-bit integer | <sys/types.h> |
| uint32_t | Unsigned 32-bit integer | <sys/types.h> |
| sa_family_t | Address family of socket address structure | <sys/socket.h> |
| socklen_t | Length of socket address structure, normally uint32_t | <sys/socket.h> |
| in_addr_t | IPv4 address, normally uint32_t | <netinet/in.h> |
| in_port_t | TCP or UDP port, normally uint16_t | <netinet/in.h> |

Figure 2.1 lists these three POSIX-defined data types

While a ccessing 32 bi t Ipv4 address, i f s erv is d efined as an Internet socket a ddress s tructure, 32 -bit IPv4 a ddress i n_addr s tructure i s referenced as s erv.sin_addr, w hile s erv.sin_addr.s_addr r eferences t he same 3 2-bit IPv4 a ddress a s a n i n_addr_t ( typically a n uns igned 32 -bit integer).

## IPv6 Socket Address Structure

The IPv6 socket address is defined by including the <netinet/in.h> header file. The IPv6 family is AF_INET6, whereas the IPv4 family is AF_INET. IPv6 socket address structure sockaddr_in6 is shown below.

struct sockaddr_in6

{

    uint8_t         sin6_len;        // sizeof this struct - 28 bytes

    sa_family_t     sin6_family;     // AF_INET6

    in_port_t       sin6_port;

    uint32_t       sin6_flowinfo;

    struct in6_addr    sin6_addr;     // 128 bit IPv6 address

    uint32_t       sin6_scope_id;

};

struct in6_addr

{

    uint8_t    s6_addr [16];    // IPv6 addresss (16 bytes - 128 bits)

};

The m embers i n t his s tructure are o rdered s o t hat i f t he s ockaddr_in6 structure is 64-bit a ligned, s o is the 128-bit s in6_addr m ember. O n s ome 64-bit processors, data accesses of 64-bit values are optimized i f stored on a 64-bit boundary.

## 2.3 VALUE-RESULT ARGUMENTS

When a socket address structure is passed to any socket function, it is always passed by reference (a pointer to the structure is passed). The length of the structure is also passed as an argument.

The way in which the length is passed depends on which direction the structure is being passed:

1. From the **process to the kernel**

2. From the **kernel to the process**

### From process to kernel

*Bind*, *connect*, and *sendto* are the functions that pass a socket address structure from the process to the kernel. Two of the Arguments to these functions are:

- The pointer to the socket address structure

- The integer size of the structure

Because of these two arguments, kernel knows how much data to copy from process to kernel.

```
struct sockaddr_in serv;

/* fill in serv{} */

connect (sockfd, (SA *) &serv, sizeof(serv));
```

The datatype for the size of a socket address structure is actually socklen_t and not int, but the POSIX specification recommends that socklen_t be defined as uint32_t.

### From kernel to process

*Accept*, *recvfrom*, *getsockname*, and *getpeername* are the functions that pass a socket address structure from the kernel to the process.

Two of the Arguments to these functions are:

- The pointer to the socket address structure

- The pointer to an integer containing the size of the structure.

```
struct sockaddr_un  cli;   /* Unix domain */

socklen_t  len;

len = sizeof(cli);        /* len is a value */

getpeername(unixfd, (SA *) &cli, &len);

/* len may have changed */
```

**Value-result argument** (Figure 3.2): the size changes from an integer to be a pointer to an integer because the size is both a value when the function is called and a result when the function returns.

- As a **value**: it tells the kernel the size of the structure so that the kernel does not write past the end of the structure when filling it in

- As a **result**: it tells the process how much information the kernel actually stored in the structure

For two other functions that pass socket address structures, recvmsg and sendmsg, the length field is not a function argument but a structure member.

If the socket address structure is fixed-length, the value returned by the kernel will always be that fixed size: 16 for an Ipv4sockaddr_in and 28 for an Ipv6 sockaddr_in6. But with a variable-length socket address structure (e.g., a Unix domainsockaddr_un), the value returned can be less than the maximum size of the structure.



Figure 2.2: Value-result argument
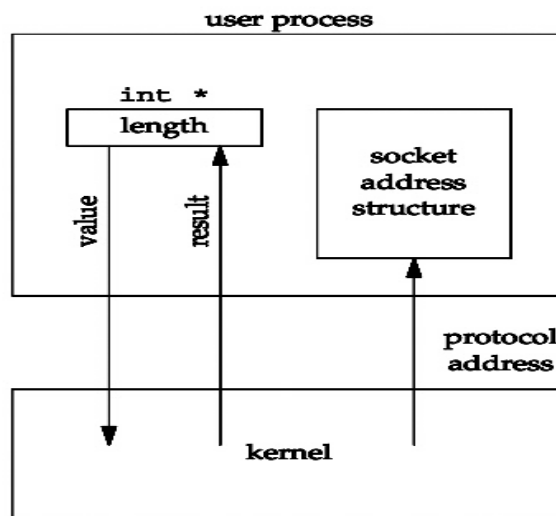
---

**CHECK YOUR PROGRESS**

1. What is *sockaddr*?

2. How is the length of socket address structure sent from a process to kernel?

---

## 2.4  BYTE ORDERING FUNCTIONS

A 16-bit integer made up of 2 bytes can be stored in memory in two ways:

- **Little-endian** order: low-order byte is at the starting address.

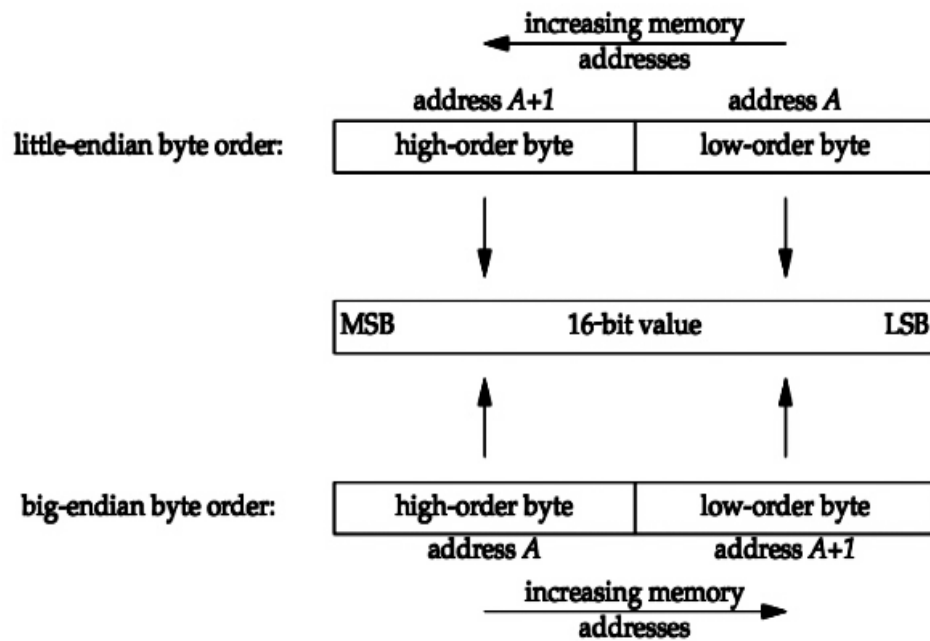- **Big-endian** order: high-order byte is at the starting address.

Figure 2.3: Byte Order

Figure 2.3 shows increasing memory addresses from right to left in the top and left to right in the bottom. Byte ordering used by a given system is called **host byte order.** The same applies to 32-bit integer. There are a variety of systems that can between little and big endian byte order at system reset or run time.

Since, a ll n etworking p rotocols s pecify the ne twork b yte or der w hile transferring d ata, it is imperative for a p rogrammer t o unde rstand t he ordering di fferences. F or e xample, T CP s egments transferred between nodes c ontain a 1 6-bit port num ber a nd a 32 bi t IPV4 ne twork a ddress. The r eceiving and s ending ne twork p rotocol s tacks m ust a gree on t he order in w hich th ese multibyte f ields a re transmitted. The I nternet protocols use big-endian byte ordering for these multibyte integers.

But, both history and the POSIX specification say that certain fields in the socket a ddress s tructures m ust be m aintained i n network b yte or der. We use the following four functions to convert between these two byte orders:

```
#include <netinet/in.h>

uint16_t htons(uint16_t host16bitvalue);

uint32_t htonl(uint32_t host32bitvalue); /* Both return: value in network
byte order */

uint16_t ntohs(uint16_t net16bitvalue);

int32_t ntohl(uint32_t net32bitvalue); /* Both return: value in host byte
order */
```

- h stands for *host*

- n stands for *network*

- s stands for *short* (16-bit value, e.g. TCP or UDP port number)

- l stands for *long* (32-bit value, e.g. IPv4 address)

When using these functions, we do not care about the actual values (big-endian or little-endian) for the host byte order and the network byte order. What we must do is call the appropriate function to convert a given value between the host and network byte order. On those systems that have the same byte ordering as the Internet protocols (big-endian), these four functions are usually defined as null macros. An 8-bit entity is a "Byte" but most Internet standards use the term octet instead of Byte.

## Byte Manipulation Functions

Socket address structures are manipulated using two groups of functions that operate on Multibyte fields. These functions do not interpret data and do not assume data as null-terminated C string. These functions are necessary to manipulate socket assress structures which have IP addresses that have bytes of 0 and not null-terminated C strings. The two groups of function are as follows:

- ❖ Ones whose names start with b (for byte) are from 4.2BSD and are still provided by almost any system that supports the socket functions. Examples for this type are bzero, bcopy and bcmp.

```
void bzero(void *dest, size_tnbytes);

void bcopy(const void *src, void *dest, size_tnbytes);

int bcmp(const void *ptr1, const void *ptr2, size_tnbytes);

                                    /*Returns: 0 if equal, nonzero if unequal*/
```

bzero is used to initialize the socket addresses as it sets the specified number of bytes to 0 at destination, bcopy moves the specified number of bytes from the source to the destination, and bcmp compares two arbitrary byte strings.

- ❖ Ones whose names start with mem(for memory), are from ANSI C standard and are provided with any system that supports an ANSI C library. Examples for this type are memset, memcpy and memcmp.

```
void *memset(void *dest, intc, size_tlen);

void *memcpy(void *dest, const void *src, size_tnbytes);

int memcmp(const void *ptr1, const void *ptr2, size_tnbytes);
            /*Returns: 0 if equal, <0 or >0 if unequal*/
```

memset sets the specified number of bytes to the value c in the destination, memcpy is similar to bcopy, but the order of the two pointer arguments is swapped. bcopy correctly handles overlapping fields, while the behavior of memcpy is undefined if the source and destination overlap.

The two pointers for memcpy are written in the same left-to-right order as an a ssignment s tatement in C . A ll memxxx f unctions r equire a length argument which is th e f inal a rgument me mcmp c ompares tw o a rbitrary byte strings and returns 0 if they are identical. If not identical,      the return va lue is e ither greater than 0 or less than 0, depending on w hether the f irst une qual b yte p ointed t o b y pt r1 i s gr eater t han or l ess t han t he corresponding b yte pointed to b y ptr2. The comparison is done assuming the two unequal bytes are unsigned chars.

## 2.5   RELATED FUNCTIONS

**inet_aton, inet_addr, and inet_ntoa Functions**

We w ill de scribe t wo groups o f a ddress c onversion f unctions i n t his section a nd t he ne xt. They convert Internet ad dresses b etween A SCII strings ( what hum ans pr efer t o us e) a nd network b yte or dered bi nary values (values that are stored in socket address structures).

1. inet_aton, inet_ntoa, and inet_addr convert an IPv4 address from a dotted-decimal string (e.g., "206.168.112.96") to its 32-bit network byte ordered binary value. You will probably encounter these functions in lots of existing code.

2. The newer functions, inet_pton and inet_ntop, handle both IPv4 and IPv6 addresses. We describe these two functions in the next section and use them throughout the text.

---

#include <arpa/inet.h>

int inet_aton(const char *strptr, struct in_addr *addrptr);

Returns: 1 if string was valid, 0 on error

in_addr_t inet_addr(const char *strptr);

Returns:   32-bit binary   network   byte   ordered   IPv4

address;   INADDR_NONE if error

char *inet_ntoa(struct in_addr inaddr);

Returns: pointer to dotted-decimal string

---

The first of these, inet_aton, converts the C character string pointed to by strptr into its 32-bit binary network byte ordered value, which is    stored through the pointer addrptr. If successful, 1 is returned; otherwise, 0 is returned.

An undocumented feature of inet_aton is that if addrptr is a null pointer, the function still performs its validation of the input string but does not store any result.

inet_addr does the same conversion, returning the 32-bit binary network byte ordered value as the return value. The problem with this function is that all $2^{32}$ possible binary values are valid IP addresses (0.0.0.0 through

255.255.255.255), but the function returns the constant INADDR_NONE (typically 32 on e-bits) o n a n e rror. T his m eans t he dotted-decimal s tring 255.255.255.255 (the IPv4 limited broadcast address)

cannot be ha ndled b y t his f unction s ince i ts b inary va lue a ppears t o indicate failure of the function.

A p otential p roblem w ith in et_addr is th at s ome ma n p ages s tate th at it returns 1 on a n error, instead of INADDR_NONE. T his c an l ead t o problems, depending on the C compiler, when comparing the return value of the function (an unsigned value) to a negative constant.

Today, i net_addr i s de precated and a ny ne w c ode s hould us e i net_aton instead. B etter still is to u se th e n ewer f unctions d escribed in the n ext section, which handle both IPv4 and IPv6.

The inet_ntoa function converts a 32-bit binary network byte ordered IPv4 address into its corresponding dotted-decimal string. The string pointed to by t he return va lue of t he function r esides in s tatic m emory. T his m eans the f unction is not r eentrant. F inally, not ice t hat t his f unction t akes a structure as its argument, not a pointer to a structure.

Functions th at ta ke a ctual s tructures a s arguments a re r are. It is mo re common to pass a pointer to the structure.

**inet_pton and inet_ntop Functions**

These two functions are new with IPv6 and work with both IPv4 and IPv6 addresses. We use these two functions throughout the text. The letters "p" and "n" stand for presentation and numeric. The presentation format for an address is often an ASCII string and the numeric format is the binary value that goes into a socket address structure.

```
#include <arpa/inet.h>

int inet_pton(int family, const char *strptr, void *addrptr);

Returns: 1 if OK, 0 if input not a valid presentation format, -1 on error

const char *inet_ntop(int family, const void *addrptr, char *strptr, size_t len);

Returns: pointer to result if OK, NULL on error
```

The family argument for both functions is either AF_INET or AF_INET6. If family is not supported, both functions return an error with errno set to **EAFNOSUPPORT.**

The first function tries to convert the string pointed to by strptr, storing the binary result through the pointer addrptr. If successful, the return value is 1.

If the input string is not a valid presentation format for the specified family, 0 is returned.

inet_ntop doe s t he r everse conversion, f rom num eric ( addrptr) t o presentation (strptr).

The l en a rgument i s t he s ize of t he de stination, t o pr event t he f unction from ove rflowing the cal ler's bu ffer. T o he lp s pecify t his s ize, t he following t wo de finitions a re de fined b y including t he < netinet/in.h> header:

#define INET_ADDRSTRLEN 16 /* for IPv4 dotted-decimal */

#define INET6_ADDRSTRLEN 46 /* for IPv6 hex string */

If len is too small to hold the resulting presentation format, including the terminating null, a null pointer is returned and errno is set to ENOSPC.

The strptr argument to inet_ntop cannot be a null pointer. The caller must allocate memory for the destination and specify its size. On success, this pointer is the return value of the function.

---

**Check your progress**

1.    What is host byte order?

2.    What do inet_pton and inet_ntop Functions do?

---

## 2.6    SUMMARY

Sockets ar e an i ntegral p art o f ev ery n etwork p rogram. T he ad dress structures of the sockets are filled and passed as pointers to the various socket functions. Wh en a p ointer t o o ne o f t hese s tructures i s p assed t o a so cket function, it fills in the contents. These structures are always passed by reference and t he s ize o f t he s tructure is p assed as another argument. Wh en a socket function fills the structure, the length is also passed as reference, so that the value of t he length c an b e upd ated by t he function. T hese ar e termed as value-result arguments.

The address structures are self-defining because they contain a field ("domain") that specifies the address family contained i n the structure. Newer implementations supporting variable-length address structures also contain a length field at the beginning, indicating the length of the entire structure.

The t wo functions t hat c onvert IP a ddresses b etween pr esentation format (what we write, such as ASCII characters) and numeric format (what goes into a s ocket address s tructure) a re i net_pton a nd i net_ntop. T hese t wo functions a re, how ever, pr otocol- dependent. A be tter t echnique i s t o manipulate t he socket a ddress s tructures as opa que obj ects, know ing j ust the pointer to the structure and its size.

## 2.7   TERMINAL QUESTIONS

1.    Define socket and list out its types.

2.    Compare t he  IPV4,  IPV6, U nix dom ain a nd  data l ink s ocket address structures. State your assumptions.

3.    Describe IPV4 and IPV6 socket address structure.

4.    What is byte ordering function?

5.    Explain in detail about address conversion functions.

6.    Explain value-result arguments.

# UNIT-3 :   ELEMENTARY TCP SOCKETS

## Structure

## 3.0   INTRODUCTION

This uni t  describes t  he s  ocket ne  twork  Inter P  rocess Communication ( IPC) i nterface, w  hich c  an be  us ed b  y p  rocesses t o communicate with other processes, regardless of where they are running, i.e. t he s ame i nterfaces c an be  us ed f or bot h i nter a nd i ntra m achine communication. T he s ocket i nterface c an be  us ed t o c ommunicate us ing many di fferent ne twork pr otocols. H owever, our  di scussion s hall be restricted to the TCP/IP protocol suite, since it is the de facto standard for communicating over the Internet.

## 3.1   OBJECTIVES

After the end of this unit, you should be able to:

- Understand t  he di  fferent e  lementary f  unctions r  equired f  or establish a TCP connection between a server and client

- Understand the difference between concurrent and iterative servers

- Implement a concurrent server

- Gain an insight of other functions related to socket communication

## 3.2   ELEMENTARY   TCP   SOCKETS   AND SOCKET FUNCTION

The timeline of a typical scenario that takes place between a TCP client and server has been shown in Figure 3.1. First, the server is started, and then sometime later, a client is started that connects to the server. We assume that the client sends a r equest to the s erver, the s erver p rocesses the request, and the server sends a reply back to the client. This continues until the client closes its end of the connection, which sends an end-of-file notification to the server. The server then closes its end of the connection and either terminates or waits for a new client connection.
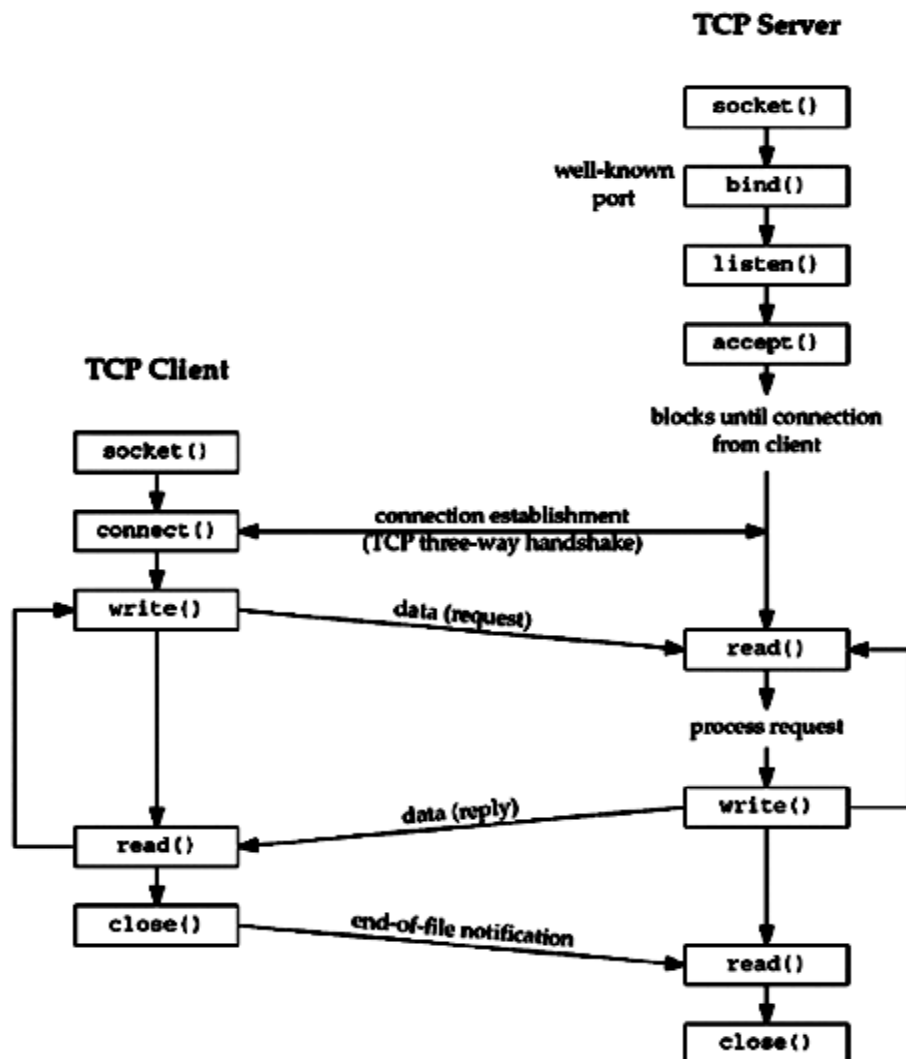


Fig. 3.1: TCP Client and Server

A socket is an abstraction of a communication endpoint. Just as they would use file d escriptors  to acc ess  files, ap plications u se so cket d escriptors t o ac cess sockets. S ocket d escriptors ar e  implemented as f ile d escriptors i n t he U NIX

System. Indeed, many of the functions that deal with file descriptors, such as read and write, will work with a socket descriptor.

To pe rform network I/O, the first t hing a process must d o is call t he socket function, s pecifying t he t ype of c ommunication pr otocol de sired ( TCP using IPv4, UDP using IPv6, Unix domain stream protocol, etc.).

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

int socket(int domain, int type, int protocol);

Returns: file (socket) descriptor if OK, −1 on error
```

The *domain* parameter specifies the nature of the communication, as well as t                                                                        he address f ormat. T able 3. 1 l ists t he dom ains s pecified b y P OSIX.1. T he constants start with AF_ (for address family) because each domain has its own representation format for an address.

| Domain | Description |
|---|---|
| AF_INET | IPv4 Internet domain |
| AF_INET6 | IPv6 Internet domain |
| AF_Local | UNIX domain |
| AF_ROUTE | Routing sockets |
| AF_KEY | Key socket |
| AF_UNSPEC | Unspecified |

Table 3.1: Socket Communication Domains

The s econd p arameter *type* specifies the t ype of t he s ocket, which further determines th e communication characteristics. T able 3 .2 s ummarizes t he socket types defined by POSIX.1.

| Type | Description |
|---|---|
| SOCK_DGRAM | fixed-length, connectionless, unreliable messages |
| SOCK_SEQPACKET | fixed-length, s equenced, r eliable, c onnection-oriented messages |
| SOCK_STREAM | sequenced, r eliable, bi directional, c onnection-oriented byte stream |
| SOCK_RAW | datagram interface to IP (optional in POSIX.1) |

Table 3.2: Socket Types

The third argument *protocol* is usually zero, to select the default protocol for t he g iven dom ain a nd s ocket t ype. W hen m ultiple pr otocols a re supported for the same domain and socket type, we c an use the *protocol* argument t o s elect a p articular pr otocol. T he de fault p rotocol f or a SOCK_STREAM socket in the AF_INET communication domain is TCP (Transmission C ontrol P rotocol). T he de fault pr otocol f or a SOCK_DGRAM socket in the AF_INET communication domain is UDP (User D atagram P rotocol). T able 3.3 l ists t he pr otocols de fined f or t he Internet domain sockets.

| Protocol | Description |
|---|---|
| IPPROTO_TCP | TCP transport protocol |
| IPPROTO_UDP | UDP transport protocol |
| IP_PROTO_SCTP | SCTP transport protocol |

Table 3.3: Protocols defined for Internet domain sockets

On success, the *socket* function returns a small non-negative integer value. This is termed as socket descriptor, denoted by *sockfd*. This socket descriptor depends upon the protocol family (IPv4, IPv6, or Unix) and the type of the (stream, datagram, or raw).

These sockets support bi directional communication. The I/O operations on a socket can be disabled by using the *shutdown* function.

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

int shutdown(int sockfd, int how);

Returns: 0 if OK, −1 on error
```

If *how* is set to SHUT_RD, then reading from the socket is disabled. If *how* is SHUT_WR, then the socket can't be used for transmitting data. The value of *how* can be set to SHUT_RDWR to disable both data transmission and reception.

If a socket can be closed then, why a shutdown is required? There can be several reasons.

- First, close will deallocate the network endpoint only when the last active reference is closed. If we duplicate the socket (with dup, for example), the socket won't be deallocated until we close the last file descriptor referring to it. The shutdown function allows us to deactivate a socket independently of the number of active file descriptors referencing it.

- Second, it is sometimes convenient to shut a socket down in one direction only. For example, we can shut a socket down for writing if we want the process we are communicating with to be able to tell when we are done transmitting data, while still allowing us to use the socket to receive data sent to us by the process.

## 3.3   CONNECT FUNCTION

A connection-oriented network service (SOCK_STREAM or SOCK_SEQPACKET) requires that before data is exchanged, a connection must be established between the socket of the process requesting the service (the client) and the process providing the service (the server). The connect function to create a connection.

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

int connect(int sockfd, const struct sockaddr *addr, socklen_t len);

Returns: 0 if OK, −1 on error
```

Here, s ockfd i s a s ocket d escriptor r eturned b y the s ocket function. The second and third arguments are a pointer to a socket address structure and its size. The socket address structure must contain the IP address and port number of the server.

While c onnecting to a s erver, the c onnect r equest mig ht f ail for mu ltiple reasons. F or a co nnect r equest t o s ucceed, t he m achine t o w hich w e ar e trying to connect must be up and running, the server must be bound to the address w e a re t rying t o c ontact, and there m ust be r oom i n t he s erver's pending connect queue.

---

**Check your progress**

1.    What are the reasons of shutting down a socket?

2.    Enlist some reasons of the connect request failure.

---

## 3.4   BIND FUNCTION

The bi nd f unction a ssigns a l ocal pr otocol a ddress t o a s ocket. With t he Internet pr otocols, t he pr otocol a ddress i s t he c ombination of either a 3 2-bit IPv4 address or a 128-bit IPv6 address, along with a 16-bit TCP or UDP port number.

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

int bind (int sockfd, const struct sockaddr *myaddr, socklen_t addrlen);

                                    Returns: 0 if OK, −1 on error
```

The second a rgument i s a poi nter t o a pr otocol-specific ad dress, an d t he third argument is the size of this address structure.

The bind function lets us specify the IP address, the port, both, or neither. Table 3.4 e nlists the va lues to w hich s in_addr a nd s in_port, or s in6_addr and sin6_port, can be set as per the requirement.

| Process Specifies | | Result |
|---|---|---|
| **IP address** | **Port** | |
| Wildcard | 0 | Kernel chooses IP address and port |
| Wildcard | nonzero | Kernel chooses IP address, process specifies port |
| Local IP address | 0 | Process specifies IP address, kernel chooses port |
| Local IP address | nonzero | Process specifies IP address and port |

Table 3.4: Result when specifying IP address and/or port number to bind

If IPv4 is being used, the wildcard address can be denoted by the constant INADDR_ANY, whose value is normally 0. This informs the kernel to choose the IP address.

```
struct sockaddr_in servaddr;

servaddr.sin_addr.s_addr = htonl (INADDR_ANY);          /* wildcard */
```

This technique cannot be used with IPv6, since the length of IPv6 address is 128-bit which can be stored in a structure. (The C language does not allow a constant structure on the right-hand side of an assignment.) The following code snippet shows the method of assigning wildcard address in case of IPv6.

```
struct sockaddr_in6 serv;

serv.sin6_addr = in6addr_any;                           /* wildcard */
```

The extern declaration of the variable in6addr_only is present in the <netinet/in.h> header file. The system shall allocate the memory and initialize the in6addr_any variable to the constant IN6ADDR_ANY_INIT.

The value of INADDR_ANY (0) shall be the same in either host or network. This eliminates the need of htonl. However, since the header <netinet/in.h> defines all the INADDR_ constants in host byte order, the function htonl should be used with any of these constants.

## 3.5   LISTEN FUNCTION

A server announces that it is willing to accept connect requests by calling the listen function. The call to the socket function always created an active socket, i.e. a client socket that can issue a connect. The listen function converts an unconnected socket into a passive socket. This denotes that the kernel should accept incoming connection requests directed to this socket.

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

int listen(int sockfd, int backlog);

Returns: 0 if OK, −1 on error
```

The backlog argument provides a hint to the system regarding the number of outstanding connect requests that it should enqueue on behalf of the process. The actual value is determined by the system, but the upper limit is specified as SOMAXCONN in <sys/socket.h>.

To understand the backlog argument, we must realize that for a given listening socket, the kernel maintains two queues:

An incomplete connection queue, which contains an entry for each SYN that has arrived from a client for which the server is awaiting completion of the TCP three- way handshake. These sockets are in the SYN_RCVD state.

A completed connection queue, which contains an entry for each client with whom the TCP three-way handshake has completed. These sockets are in the ESTABLISHED state.

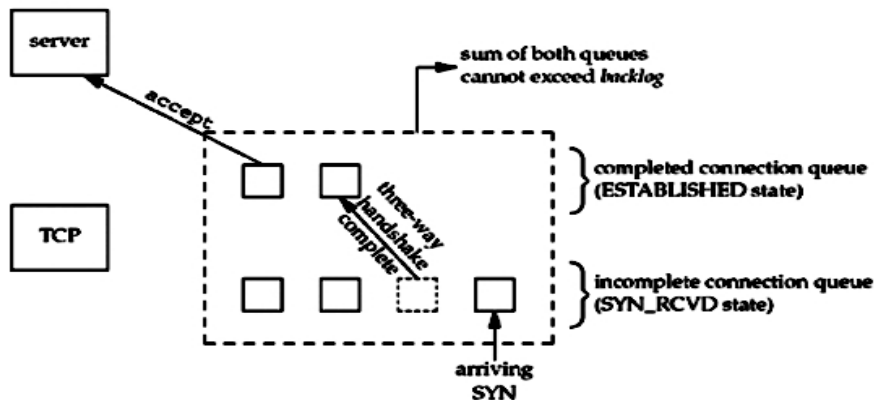Figure 3.2 depicts these two queues for a given listening socket.



Fig 3.2: Queues maintained by TCP for a listening socket

Whenever an entry is created in the incomplete queue, the arguments from the listen socket are copied to the newly created connection. Once the queue is full, the system will reject additional connect requests, so the backlog value must be chosen based on the expected load of the server and the amount of processing it must do to accept a connect request and start the service.

## 3.6   ACCEPT FUNCTION

Once a server has called listen, the socket can receive connect requests. The accept function is used to retrieve a connect request and establish a connection.

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

int accept(int sockfd, struct sockaddr *restrict addr, socklen_t *restrict len);

Returns: file (socket) descriptor if OK, −1 on error
```

If the call to accept is successful, a "new" descriptor is returned by the kernel. This new descriptor identifies the TCP connection that has been established with the client. The sockfd argument specifies the listening socket (created upon successful call to the socket function) while the

return v alue o f a ccept i s t ermed as  t he  connected s ocket.  A s erver normally creates only one listening socket which exists for the lifetime of the server. A connected socket is created for each client connection that is accepted (i.e., for which the TCP three-way handshake completes). When the server is finished serving a client, the connected socket is closed.

This f unction r eturns up  t o t hree v alues: an  i nteger r eturn co de t hat i s either a new socket descriptor or an error indication, the protocol address of t he c lient pr ocess ( through t he c liaddr poi nter), a nd t he s ize of  t his address (through the addrlen pointer). If the protocol address of the client is not required, both cliaddr and addrlen are set to null pointers.

## 3.7   FORK AND EXEC FUNCTION

The fork function is used for creating a new process in Unix.

```
#include <unistd.h>

pid_t fork(void);

                Returns: 0 in child, process ID of child in parent, −1 on error
```

The fork function is "called once but returns twice". The process ID of the newly created child process is returned to the parent (the calling process) while the value 0 is returned to the child process.

The reason why 0 i s returned to the child, instead of the process ID of the parent, i s t hat a  c hild ha s onl y one  pa rent a nd i t c an a lways obt ain t he parent's pr ocess  ID b y  calling  getppid f unction. A  pa rent, on t  he ot her hand, can have any number of children, and there is no way to obtain the process IDs of its children. If a parent wants to keep track of the process IDs of all its children, it must record the return values from fork.

Any descriptor opened by the parent process before calling fork shall be shared with the child process after fork executes successfully. This feature is u sed b y t he n etwork s ervers w here t he p arent f irst cal ls ac cept t o establish a connection with the client and then calls fork. This ensures that the c onnected s ocket i s s hared be tween t he pa rent a nd t he  child pr ocess. The child process can then read and write on the connected socket and the parent can close the same connected socket.

The fork function is generally used for the following purpose:

1.    Making a copy of the process ensures that so that one copy handles one operation while the other copy performs another task. This is, generally, the case in network servers.

2.    When a p rocess w ants t o ex ecute an other p rogram, i t f irst cal ls
       fork t o m ake a c opy o f itself, a nd t hen one of t he c opies ( child)
       calls ex ec t o replace itself w ith the new p rogram. This i s the cas e
       for programs such as shells.

The only way in which an executable program file on disk can be executed
by Unix i s for a n existing p rocess t o c all one o f t he s ix e xec functions.
exec r eplaces t he cu rrent p rocess i mage w ith t he n ew p rogram f ile, an d
this ne w pr ogram nor mally s tarts a t t he m ain f unction. T he pr ocess ID
does n ot ch ange. W e refer t o t he p rocess t hat c alls ex ec as t he cal ling
process and the newly executed program as the new program.

The differences in the six exec functions are:

   **a)**    whether the program file to execute is specified by a filename or a
              pathname;

   **b)**    whether the arguments to the new program are listed one by one or
              referenced through an array of pointers; and

   **c)**    whether t he environment of t he c alling p rocess i s pa ssed t o t he
              new program or whether a new environment is specified.

The six variants of the exec function have been shown below:

```
#include <unistd.h>

int execl (const char *pathname, const char *arg0, ... /* (char *) 0 */ );

int execv (const char *pathname, char *const argv[]);

int execle (const char *pathname, const char *arg0, ... /* (char *) 0,

            char *const envp[] */ );

int execve (const char *pathname, char *const argv[], char *const envp[]);

int execlp (const char *filename, const char *arg0, ... /* (char *) 0 */ );

int execvp (const char *filename, char *const argv[]);

                All six return: -1 on error, no return on success
```

These f unctions r eturn t o t he c aller onl y i f a n error oc curs. O therwise,
control passes to the start of the new program, normally the main function.

Figure 3.3 depicts t he r elationship a mong t hese s ix va riants of t he e xec
function. It s hould be not ed t hat, onl y e xecve i s a s ystem c all w hile t he
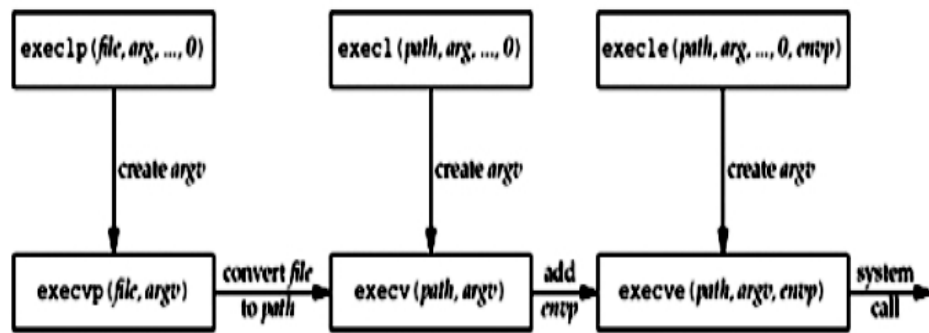remaining five are library functions that internally call execve.

Figure 3.3: Relationship among the six exec functions

There exist following differences among these variants of exec functions:

The functions execlp, execl and execle considers each string parameter as a separate parameter to the exec function, with a null pointer terminating the variable number of parameters. The functions execvp, execv and execve have an argv array, that contains pointers to the string parameters. The argv array must contain a null pointer to specify its end, since a count is not specified.

The functions execlp and execvp require a file parameter specifying the filename. This is converted into a pathname by using the current PATH environment variable. However, if the filename parameter to execlp or execvp contains a slash (/) anywhere in the string, the PATH variable is not used. The remaining functions require a fully qualified pathname argument.

The functions execlp, execl, execvp, and execv do not require an explicit environment pointer. The current value of the external variable environ is used for building an environment list that is passed to the new program. The functions execle and execve require an explicit environment list. The parameter envp is an array of pointers terminated by a null pointer.

---

**Check your progress**

1.    What is the default value of INADDR_ANY?

2.    What is the use of Fork function?

---

# 3.8   CONCURRENT SERVERS

When a client request requires longer time to complete, it is not feasible to dedicate a single server for one client. The server should be able to serve multiple client requests at the same time. Such type of servers is termed as concurrent servers. The simplest method to implement

a co ncurrent s erver i s t o f ork a ch ild p rocess for s erving e ach client request. T he f ollowing c ode s nippet s hows t he i mplementation f or a typical concurrent server.

```
pid_t pid;
int listenfd, connfd;
listenfd = Socket( ... );
/* fill in sockaddr_in{} with server's well-known port */
Bind(listenfd, ... );
Listen(listenfd, LISTENQ);
for ( ; ; ) {
        connfd = Accept (listenfd, ... ); /* probably blocks */
        if( (pid = Fork()) == 0) {
                Close(listenfd); /* child closes listening socket */
                doit(connfd); /* process the request */
                Close(connfd); /* done with this client */
                exit(0); /* child terminates */
        }
        Close(connfd); /* parent closes connected socket */
}
```

When a connection is established, accept returns, the server calls fork, and the child process services the client (on connfd, the connected socket) and the parent process waits for another connection (on listenfd, the listening socket). The parent closes the connected socket since the child handles the new client.

The f unction does whatever i s r equired t o s ervice t he c lient. W hen t his function returns, we explicitly close the connected socket in the child. This is n ot r equired s ince t he n ext s tatement cal ls e xit, an d p art o f p rocess termination is to c lose all o pen d escriptors b y t he ke rnel. W hether t o include th is e xplicit c all to c lose o r n ot is a ma tter o f p ersonal programming taste.

## 3.9   CLOSE FUNCTION

The c lose function i s us ed t o c lose a n ope n s ocket a nd t erminate the TCP connection.

```
#include <unistd.h>

int close (int sockfd);

Returns: 0 if OK, -1 on error
```

The de fault a ction of c lose f unction w ith a T CP s ocket i s t o m ark t he socket as closed and return to the calling process immediately. The socket descriptor s hall no l onger be us able b y t he pr ocess. It c annot be f urther used for e ither da ta t ransmission or r eception. H owever, T CP w ill t ry t o send any data that has already been queued to be sent, after which the TCP connection termination procedure takes place.

## 3.10  RELATED FUNCTIONS

These getsockname f unction r eturns t he l ocal protocol a ddress while t he f unction ge tpeername r eturns t he foreign pr otocol address associated with a socket.

```
#include <unistd.h>

int g etsockname(int s ockfd, s truct s ockaddr *l ocaladdr, s ocklen_t
*addrlen);

int g etpeername(int s ockfd, s truct s ockaddr *pe eraddr, s ocklen_t
*addrlen);

                                      Both return: 0 if OK, -1 on error
```

The f unctions r eturn t he c ombination of a n IP a ddress a nd por t num ber associated with one of the two ends of a network connection.

These two functions are required for the following reasons:

- After connect successfully returns in a TCP client that does not call bind, g etsockname r eturns t he l ocal IP a ddress and l ocal por t number assigned to the connection by the kernel.

- After c alling bi nd w ith a por t num ber of 0 ( telling t he ke rnel t o choose the l ocal port number), getsockname r eturns the l ocal por t number that was assigned.

- getsockname can be called to obtain the address family of a socket.

- In a T CP s erver t hat binds t he w ildcard IP address, onc e a connection i s es tablished w ith a cl ient (accept r eturns successfully), t he s erver can cal l getsockname t o o btain t he l ocal IP a ddress a ssigned t o t he c onnection. T he s ocket de scriptor argument in this call must be that of the connected socket, and not the listening socket.

- When a s erver i s ex eced b y the p rocess that c alls accept, the o nly way th e s erver c an o btain th e id entity o f th e client is to c all getpeername.

Example: Obtaining the Address Family of a Socket

The sockfd_to_family function returns the address family of a socket. The following c ode s nippet de monstrates r eturning t he a ddress f amily o f a socket.

```
1 #include "unp.h"
2 int sockfd_to_family(int sockfd)
4 {
5       struct sockaddr_storage ss;
6       socklen_t len;
7       len = sizeof(ss);
8       if (getsockname(sockfd, (SA *) &ss, &len) < 0)
9               return (-1);
10      return (ss.ss_family);
11 }
```

## Check your progress

Write a p rogram t o ex change o ne h ello m essage b etween server an d client to demonstrate the client/server model.

## 3.11 SUMMARY

A call to the socket function returns a socket descriptor which can be used for inter process communication between two different machines or on the same machine. Clients wishing to establish a connection with the server call t he c onnect function w hile s ervers c all t he bi nd, l isten, a nd accept function to accept connections from the client. Open socket scan be closed by issuing a call to the standard close function, although there also exist a shutdown function for the similar purpose.

TCP servers should be able to serve concurrent requests from the clients. This is achieved by calling fork function for every client connection being handled by the server. However, UDP servers are, generally, iterative in nature.

## 3.12 TERMINAL QUESTIONS

1.  Describe the procedure and sequence of function calls required for establishing a TCP connection between a client and the server.

2.  In Section 3.4, we stated that the INADDR_ constants defined by the <netinet/in.h> header are in host byte order. How can we tell this?

3.  An iterative server waits for the child to execute the command and exit before accepting the next connect request. Write a pseudocode for the server so that the time to service one request doesn't delay the processing of incoming connect requests.

4. Refer to code for concurrent servers in Section 3.7. Assume the child runs first after the call to fork. The child then completes the service of the client before the call to fork returns to the parent. What happens in the two calls to close?

5. Write a program to implement a chat room environment for one client and one server.

6. Write a program to implement a chat room environment for multiple client and one server.

# UNIT-4 :     TCP Client/Server

## Structure

## 4.0    INTRODUCTION

A s erver pr ovides a s ervice on a given por t b y waiting f or connections from future clients. A client can connect to a service once the server i s r eady t o a ccept co nnections ( accept). In or der t o m ake a connection, the client must know the IP number of the server machine and the port number of the service. If the client does not know the IP number, it ne eds t o r equest na me/number r esolution. O nce t he c onnection i s accepted b y t he s erver, each p rogram can co mmunicate v ia i nput-output channels over the sockets created at both ends.

## 4.1    OBJECTIVES

After the end of this unit, you should be able to:

- Understand the various functions of TCP Echo server

- Gain i nsights r egarding nor mal s tartup, t erminate a nd s ignal handling server process termination

- Differentiate between crashing and rebooting of server host

- Know about shutdown of server host

## 4.2    TCP ECHO SERVER FUNCTION

**TCP Echo Server: Main Function**

The concurrent server program has been represented through the following code:

```
#include "unp.h"

Int main(int argx, char **argv)

{

        Int listenfd, connfd;

        pid_t childpid;

        socklen_t clien;

        struct sockaddr_in cliaddr,servaddr;

        listenfd=Socket(AF_INET,SOCK_STREAM,0);

        bzero(&servaddr,sizeof(servaddr));

        servaddr.sin_family=AF_INET;

        servaddr.sin_addr.s_addr=htonl(INADDR_ANY);

        servaddr.sin_port=htons(SERV_PORT);

        Bind(listenfd,LISTENQ);

        for( ; ;) {

        clien = sizeof(cliaddr);

        connfd = Accept(listenfd, (SA*) &cliaddr, &clilen);

        if((childpid = Fork()) == 0) { /* child process */

            Close(listenfd); /*close listening socket*/

             str_echo(connfd); /*process the request */

              exit(0);

        }

        Close(connfd); /*parent closes connected socket */

    }

}
```

The actions processed by the code are as follows:

- **Create socket, bind server's well-known port**

    ❖    A TCP socket is created.

    ❖    An Internet s ocket a ddress s tructure is f illed in w ith th e wildcard ad dress ( INADDR_ANY) a nd t he s erver's well-known port (SERV_PORT, here defined as 9877 in header).

Binding th e w ildcard a ddress te lls th e s ystem th at w e w ill accept a c onnection de stined for a ny l ocal i nterface, i n c ase the s ystem i s multihome. It s hould be greater than 1023 ( we do not ne ed a reserved por t), greater t han 500 0 ( to a void conflict w ith th e e phemeral p orts a llocated b y many Berkeley-derived implementations), less than 49152 (to avoid conflict with t he " correct" r ange of ephemeral por ts), a nd i t should not conflict with any registered port.

❖ The socket is converted into a listening socket by listen.

- ## Wait for client connection to complete

  ❖ The s erver b locks i n the c all to accept, w aiting f or a client connection to complete.

- ## Concurrent server

  ❖ For each client, fork spawns a child, and the child handles the new c lient. T he c hild closes t he l istening s ocket a nd t he parent closes the connected socket.

## TCP Echo Server: str_echo Function

The function str_echo is responsible for conducting the server processing for ea ch cl ient. It r eads d ata f rom t he cl ient an d ech oes i t b ack t o t he client.

```
#include   "unp.h"

void

str_echo(int sockfd)

{

   ssize_t   n;

   char      buf[MAXLINE];

again:
   while ( (n = read(sockfd, buf, MAXLINE)) > 0)
      Writen(sockfd, buf, n);

   if (n < 0 && errno == EINTR)
      goto again;

   else if (n < 0)
      err_sys("str_echo: read error");

}
```

The code provided above processes the following actions:

- **Read a buffer and echo the buffer**
  - ❖ Read reads data from the socket and the line is echoed back to the client by writen. If the client closes the connection (the normal scenario), the receipt of the client's FIN causes the child's read to return 0. This causes the str_echo function to return, which terminates the child.

---

### Check your progress

1. What are the steps involved in startup of TCP client/server?

2. What are choices of disposition?

---

## 4.3   NORMAL START-UP

Normal start-up of TCP client/server includes following steps:

- Start TCP server in background on host system.

  *linux % tcpserv01 &*

  [1] 17870

  When the server starts, it calls *socket, bind, listen*, and *accept*, blocking in the call to *accept*. State of the server's listening socket is verified through *netstat* program. *netstat* command is used along with *-a* flag to see all listening sockets.

- Start the client on the same host, specifying the server's IP address of 127.0.0.1 (the loopback address).

- The client calls socket and connect, the latter causing TCP's three-way handshake to take place.

- When the three-way handshake completes, connect returns in the client and accept returns in the server and connection is established.

- The client calls *str_cli*, which will block in the call to *fgets*, because we have not typed a line of input yet.

- When accept returns in the server, it calls fork and the child calls *str_echo*. This function calls *readline*, which calls read, which blocks while waiting for a line to be sent from the client.

- The server parent, on the other hand, calls *accept* again, and blocks while waiting for the next client connection.

## 4.4 TERMINATE AND SIGNAL HANDLING SERVER PROCESS TERMINATION

### Normal Termination

At this point, the connection is established and whatever we type to the client is echoed back.

```
linux % tcpcli01 127.0.0.1   # this line has been represented

earlier hello, world            # now this is typed

hello, world                    # the line is echoed

good bye

good bye

^D                      # Control-D is the terminal EOF character
```

If netstat is executed immediately, the following is received:

```
linux % netstat -a | grep 9877

tcp     0     0 *:9877              *:*              LISTEN

tcp     0     0 localhost:42758     localhost:9877   TIME_WAIT
```

The output of the netstat is piped into grep. This prints only the lines in possession of the port acquainted with the server:

- The client's side of the connection (since the local port is 42758) enters the TIME_WAIT state

- The listening server still waits for another client connection.

The following steps are involved in the normal termination of client and server:

- When we type our EOF character, fgets returns a null pointer and the function str_cli returns.

- str_cli returns to the client main function, which terminates by calling exit.

- Part of process termination is the closing of all open descriptors, so the client socket is closed by the kernel. This sends a FIN to the server, to which the server TCP responds with an ACK. This is the first half of the TCP connection termination sequence. At this point, the server socket is in the CLOSE_WAIT state and the client socket is in the FIN_WAIT_2 state.

- When the server TCP receives the FIN, the server child is blocked in a call to read, and read then returns 0. This causes the str_echo function to return to the server child main.

- The server child terminates by calling exit.

- All open descriptors in the server child are closed.

- The closing of the connected socket by the child causes the final two segments of the TCP connection termination to take place: a FIN from the server to the client, and an ACK from the client.

- Finally, the SIGCHLD signal is sent to the parent when the server child terminates.

- This occurs in this example, but we do not catch the signal in our code, and the default action of the signal is to be ignored. Thus, the child enters the zombie state. We can verify this with the ps command.

```
linux % ps -t pts/6 -o pid,ppid,tty,stat,args,wchan

 PID  PPID TT     STAT COMMAND        WCHAN

22038 22036 pts/6   S   -bash          read_chan

17870 22038 pts/6   S   ./tcpserv01     wait_for_connect

19315 17870 pts/6   Z   [tcpserv01 <defu do_exit
```

The STAT of the child is now Z (for zombie).

**Zombie Process:** It is a process that has completed execution (via the exit system call) but still has an entry in the process table. This process is in the "Terminated State". This occurs for child processes, where the entry is still needed to allow the parent process to read its child's exit status: once the exit status is read via the wait system call, the zombie's entry is removed from the process table and it is said to be "reaped". A child process always first becomes a zombie before being removed from the resource table. In most cases, under normal system operation zombies are immediately waited on by their parent and then reaped by the system – processes that stay zombies for a long time are generally an error and cause a resource leak.

We need to clean up our zombie processes and doing this requires dealing with Unix signals. The following section sheds light on signal handling.

**POSIX Signal Handling**

A signal is a notification to a process regarding occurrence of an event. Often, signals are regarded as software interrupts. Signals usually occur asynchronously. This means that a process is not provided any information regarding the time of occurrence of a signal before its actual occurrence. There are different types of signal as following:

| Name | Description | Default action |
|---|---|---|
| SIGABRT | Abnormal Termination (abort) | Terminate + core |
| SIGALRM | Timer expired (alarm) | terminate |
| SIGBUS | Hardware Fault | Terminate + core |
| SIGCANCEL | Threads library internal use | Ignore |
| SIGCOUNT | Continue stopped process | Continue/ignore |
| SIGEMT | Hardware fault | Terminate + core |
| SIGFPE | Arithmetic exception | Terminate +core |
| SIGFREEZE | Checkpoint freeze | Ignore |
| SIGHUP | Hang-up | terminate |
| SIGILL | Illegal instruction | Terminate + core |
| SIGINFO | Status request from keyword | ignore |
| SIGINT | Terminal interrupt character | terminate |
| SIGIO | Asynchronous I/O | Terminate / Ignore |
| SIGIOT | Hardware fault | Terminate + core |
| SIGJVM1 | Java virtual machine internal use | ignore |
| SIGKILL | Termination | terminate |
| SIGLOST | Resource lost | terminate |
| SIGLWP | Threads library internal use | Terminate/ignore |
| SIGPIPE | Write to pipe with no readers | terminate |
| SIGPOLL | Pollable event (poll) | terminate |
| SIGVTALRM | Virtual time alarm(settimer) | terminate |
| SIGWAITING | threads library internal use | Ignore |

Table 4.1: UNIX System signal

Signals can be sent:

- By one process to another process (or to itself)

- By the kernel to a process.

  ❖ For example, whenever a process terminates, the kernel send a SIGCHLD signal to the parent of the terminating process.

Every signal has a disposition, which is also called the action associated with the signal. The disposition of a signal is set by calling the sigaction function. Following are the three choices for the disposition:

1. Catching a signal. We can provide a function called a signal handler that is called whenever a specific signal occurs. The two signals SIGKILL and SIGSTOP cannot be caught. Our function is called with a single integer argument that is the signal number and the function returns nothing. Its function prototype is therefore:

2. void handler (int signo);

For most signals, we can call sigaction and specify the signal handler to catch it. A few signals, SIGIO, SIGPOLL, and SIGURG, all require additional actions on the part of the process to catch the signal.

3. Ignoring a signal. We can ignore a signal by setting its disposition to SIG_IGN. The two signals SIGKILL and SIGSTOP cannot be ignored.

4. Setting the default disposition for a signal. This can be done by setting its disposition to SIG_DFL. The default is normally to terminate a process on receipt of a signal, with certain signals also generating a core image of the process in its current working directory. There are a few signals whose default disposition is to be ignored: SIGCHLD and SIGURG (sent on the arrival of out -of-band data) are two that we will encounter in this text.

**Signal Function**

The POSIX way to establish the disposition of a signal is to call the sigaction function, which is complicated in that one argument to the function is a structure (struct sigaction) that we must allocate and fill in.

An easier way to set the disposition of a signal is to call the signal function. The first argument is the signal name and the second argument is either a pointer to a function or one of the constants SIG_IGN or SIG_DFL.

However, signal is an historical function that predates POSIX. Different implementations provide different signal semantics when it is called, providing backward compatibility, whereas POSIX explicitly spells out the semantics when sigaction is called.

The solution is to define our own function named signal that just calls the POSIX sigaction function. This provides a simple interface with the

desired P OSIX s emantics. W e i nclude t his f unction i n our  ow n l ibrary, along with our err_XXX functions and our wrapper functions.

```
#include   "unp.h"
Sigfunc *
signal(int signo, Sigfunc *func)
{
    struct sigaction   act, oact;
    act.sa_handler = func;
    sigemptyset(&act.sa_mask);
    act.sa_flags = 0;
    if (signo == SIGALRM) {
#ifdef SA_INTERRUPT
        act.sa_flags |= SA_INTERRUPT;   /* SunOS 4.x */
#endif
    } else {
#ifdef SA_RESTART
        act.sa_flags |= SA_RESTART;    /* SVR4, 44BSD */
#endif
    }
    if (sigaction(signo, &act, &oact) < 0)
        return(SIG_ERR);
    return(oact.sa_handler);
}
/* end signal */
Sigfunc *
Signal(int signo, Sigfunc *func)   /* for our signal() function */
{
    Sigfunc *sigfunc;
    if ( (sigfunc = signal(signo, func)) == SIG_ERR)
        err_sys("signal error");
    return(sigfunc);
}
```

Simplify function prototype using typedef

The nor mal f unction pr ototype f or signal is c omplicated b y th e le vel  of nested parentheses.

```
void (*signal (int signo, void (*func) (int))) (int);
```

To simplify this, we define the Sigfunc type in our unp.h header as

```
typedef   void   Sigfunc(int);
```

stating that signal handlers are functions with an integer argument and the function returns nothing (void). The function prototype then becomes

```
Sigfunc *signal (int signo, Sigfunc *func);
```

A poi nter t o a s ignal h andling function i s t he second a rgument t o t he function, as well as the return value from the function.

**Set handler**

The sa_handler member of    t    he sigaction structure is    s    et t    o the *func* argument.

Set signal mask for handler

POSIX allows us to specify a set of signals that will be blocked when our signal handler is called. Any signal that is blocked cannot be delivered to a process. W e s et the sa_mask member t o the empty s et, w hich m eans that no a dditional s ignals  will be  bl ocked w hile our  s ignal ha ndler i s running. POSIX guarantees that the signal being caught is always blocked while its handler is executing.

**Set SA_RESTART flag**

SA_RESTART is a n opt ional f lag. W hen t he  flag i s s et, a s ystem cal l interrupted by this signal will be automatically restarted by the kernel.

If t  he s  ignal be  ing c  aught i  s not    SIGALRM, w  e s   pecify the SA_RESTART flag, i  f de fined. T his i s  because t  he pur pose o f generating the SIGALRM signal is normally to  place a timeout on an  I/O operation, i n w hich c ase, w  e w ant t he bl ocked s ystem c all t o be interrupted by the signal.

**Call sigaction**

We call sigaction and then return the old action for the signal as the return value of the signal function.

Throughout t his t ext,  we w ill u se th e signal function f rom t he a bove definition.

**Handling SIGCHLD Signals**

The zombie state is to maintain information about the child for the parent to fetch later, which includes:

- Process ID of the child,

- Termination status,

- Information on the resource utilization of the child.

If a parent process of zombie children terminates, the parent process ID of all the zombie children is set to 1 (the init process), which will inherit the children a nd c lean t hem up ( init will wait for t hem, w hich r emoves t he zombie).

## Handling Zombies

Zombies t ake u p s pace in t he k ernel an d ev entually we can r un o ut o f processes. W henever w e fork children, w e m ust wait for t hem t o p revent them f rom b ecoming z ombies. W e can es tablish a s ignal h andler t o catch SIGCHLD and call wait within the handler. W e establish the signal handler b y a dding t he f ollowing f unction c all a fter t he call to listen (in server's main function; it must be done before forking the first child and needs to be done only once.):

```
Signal (SIGCHLD, sig_chld);
```

The signal handler, the function sig_chld, is defined below:

```
#include   "unp.h"
void
sig_chld(int signo)
{
   pid_t   pid;
   int    stat;
   pid = wait(&stat);
   printf("child %d terminated\n", pid);
   return;

}
```

Note that calling standard I/O functions such as printf in a s ignal h andler is not recommended. We cal l printf here as a di agnostic tool to see when the child terminates.

## Compiling and running the program on Solaris

This pr ogram ( tcpcliserv/tcpserv02.c) i s c ompiled on S olaris 9 a nd us es the signal function from the system library.

```
solaris % tcpserv02 &            # start server in background

[2] 16939

solaris % tcpcli01 127.0.0.1        # then start client in foreground

hi there                 # we type this

hi there                 # and this is echoed

^D                     # we type our EOF character

child 16942 terminated            # output by printf in signal handler

accept error: Interrupted system call # main function aborts
```

The sequence of steps is as follows:

1. We terminate the client by typing our EOF character. The client TCP sends a FIN to the server and the server responds with an ACK.

2. The receipt of the FIN delivers an EOF to the child's pending readline. The child terminates.

3. The parent is blocked in its call to accept when the SIGCHLD signal is delivered. The sig_chld function executes (our signal handler), wait fetches the child's PID and termination status, and printf is called from the signal handler. The signal handler returns.

4. Since the signal was caught by the parent while the parent was blocked in a slow system call (accept), the kernel causes the accept to return an error of EINTR (interrupted system call). The parent does not handle this error (see server's main function), so it aborts.

From this example, we know that when writing network programs that catch signals, we must be cognizant of interrupted system calls, and we must handle them. In this example, the signal function provided in the standard C library does not cause an interrupted system call to be automatically restarted by the kernel. Some other systems automatically restart the interrupted system call. If we run the same example under BSD, using its library version of the signal function, the kernel restarts the interrupted system call and accept does not return an error. To handle this potential problem between different operating systems is one reason we define our own version of the signal function.

As part of the coding conventions used in this text, we always code an explicit return in our signal handlers, even though this is unnecessary for a function returning void. This reads as a reminder that the return may interrupt a system call.

**Handling Interrupted System Calls**

The term "slow system call" is used to describe any system call that can block forever, such as accept. That is, the system call need never return. Most networking functions fall into this category. Examples are:

- Accept: there is no guarantee that a server's call to accept will ever return, if there are no clients that will connect to the server.

- Read: the server's call to read in server's str_echo function will never return if the client never sends a line for the server to echo.

Other ex amples o f s low s ystem cal ls ar e r eads and w rites of pi pes a nd terminal devices. A notable exception is disk I/O, which usually returns to the caller (assuming no catastrophic hardware failure).

When a process is blocked in a slow system call and the process catches a signal a nd t he s ignal handler r eturns, t he s ystem cal l can r eturn an er ror of EINT. S ome k ernels au tomatically r estart s ome i nterrupted s ystem calls. For portability, when we write a program that catches signals (most concurrent s ervers c atch SIGCHLD), w e m ust be pr epared f or slow system calls to return EINTR.

To ha ndle a n i nterrupted accept, w e ch ange t he cal l to accept in server's main function, t he be ginning of t he f or l oop, t o t he following:

```
for ( ; ; ) {
    clilen = sizeof (cliaddr);
    if ( (connfd = accept (listenfd, (SA *) &cliaddr, &clilen)) < 0) {
        if (errno == EINTR)
            continue;        /* back to for () */
        else
            err_sys ("accept error");
    }
```

Note th at th is accept is not our w rapper function Accept, s ince w e m ust handle the failure of the function ourselves.

Restarting the interrupted system call is fine for:

- Accept

- Read

- Write

- Select

- Open

However, t here i s one function t hat w e cannot r estart: connect. If t his function returns EINTR, we cannot call it again, as doing so will return an immediate error. When connect is interrupted by a caught signal and is not automatically r estarted, we m ust cal l select to w ait f or t he c onnection t o complete.

wait and waitpid Functions

We can call wait function to handle the terminated child.

```
#include <sys/wait.h>
pid_t wait (int *statloc);
pid_t waitpid (pid_t pid, int *statloc, int options);


/* Both return: process ID if OK, 0 or–1 on error */
```

wait and waitpid both return two values: the return value of the function is the process ID of the terminated child, and the termination status of the child (an integer) is returned through the statloc pointer.

There are three macros that we can call that examine the termination status (see APUE):

- WIFEXITED: tells if the child terminated normally

- WIFSIGNALED: tells if the child was killed by a signal

- WIFSTOPPED: tells if the child was just stopped by job control

Additional macros let us then fetch the exit status of the child, or the value of the signal that killed the child, or the value of the job-control signal that stopped the child. We will use the WIFEXITED and WEXITSTATUS macros for this purpose.

If there are no terminated children for the process calling wait, but the process has one or more children that are still executing, then wait blocks until the first of the existing children terminates.

waitpid has more control over which process to wait for and whether or not to block:

- The *pid* argument specifies the process ID that we want to wait for. A value of -1 says to wait for the first of our children to terminate.

- The *options* argument specifies additional options. The most common option is WNOHANG, which tells the kernel not to block if there are no terminated children.

## Difference between wait and waitpid

The following example illustrates the difference between the wait and waitpid functions when used to clean up terminated children.

We modify our TCP client as below, which establishes five connections with the server and then uses only the first one ( sockfd[0]) in the call to str_cli. The purpose of establishing multiple connections is to spawn multiple children from the concurrent server.

```
#include    "unp.h"
int main(int argc, char **argv)
{
    int  i, sockfd[5];
    struct sockaddr_in  servaddr;

    if (argc != 2)
        err_quit("usage: tcpcli <IPaddress>");

    for (i = 0; i < 5; i++) {
        sockfd[i] = Socket(AF_INET, SOCK_STREAM, 0);

        bzero(&servaddr, sizeof(servaddr));
        servaddr.sin_family = AF_INET;
        servaddr.sin_port = htons(SERV_PORT);
        Inet_pton(AF_INET, argv[1], &servaddr.sin_addr);
        Connect(sockfd[i], (SA *) &servaddr, sizeof(servaddr));
    }
    str_cli(stdin, sockfd[0]);      /* do it all */
    exit(0);
}
```

When the client terminates, all open descriptors are closed automatically
by the kernel (we do not call close, only exit), and all five connections are
terminated at about the same time. This causes five FINs to be sent, one on
each connection, which in turn causes all five server children to terminate
at about the same time. This causes five SIGCHLD signals to be delivered
to the parent at about the same time. This causes the problem under
discussion.

We first run the server in the background and then our new client:

```
linux % tcpserv03 &
[1] 20419
linux % tcpcli04 127.0.0.1
hello                 # we type this
hello                 # and it is echoed
^D                    # we then type our EOF character
child 20426 terminated     # output by server
```

Only on e printf is out put, w hen w e e xpect a ll f ive c hildren t o ha ve terminated. If we execute ps, we see that the other four children still exist as zombies.

| PID TTY | TIME CMD |
|---------|----------|
| 20419 pts/6 | 00:00:00 tcpserv03 |
| 20421 pts/6 | 00:00:00 tcpserv03 <defunct> |
| 20422 pts/6 | 00:00:00 tcpserv03 <defunct> |
| 20423 pts/6 | 00:00:00 tcpserv03 <defunct> |

Establishing a s ignal h andler a nd c alling w ait f rom th at h andler a re insufficient for preventing zombies. The problem is that all five signals are generated before the signal handler is executed, and the signal handler is executed o nly o ne t ime b ecause U nix s ignals a re nor mally not que ued. This pr oblem i s nonde terministic. D ependent on t he t iming o f t he FINs arriving a t t he s erver h ost, t he s ignal ha ndler i s e xecuted t wo, t hree or even four times.

The c orrect s olution i s to c all waitpid instead of wait. T he code be low shows t he ve rsion of our sig_chld function t hat handles SIGCHLD correctly. T his v ersion w orks be cause we call waitpid within a l oop, f etching t he s tatus of a ny of our c hildren t hat have t erminated, with t he WNOHANG option, w hich t ells waitpid not to block if there are running children that have not yet terminated. We cannot call wait in a loop, because there is no way to prevent wait from blocking if there are running children that have not yet terminated.

```
#include   "unp.h"
void
sig_chld(int signo)
{
    pid_t   pid;
    int     stat;

    while ( (pid = waitpid(-1, &stat, WNOHANG)) > 0)
        printf("child %d terminated\n", pid);
    return;
}
```

The code below shows the final version of our server. It correctly handles a return of EINTR from accept and it establishes a signal handler (code above) that calls waitpid for all terminated children.
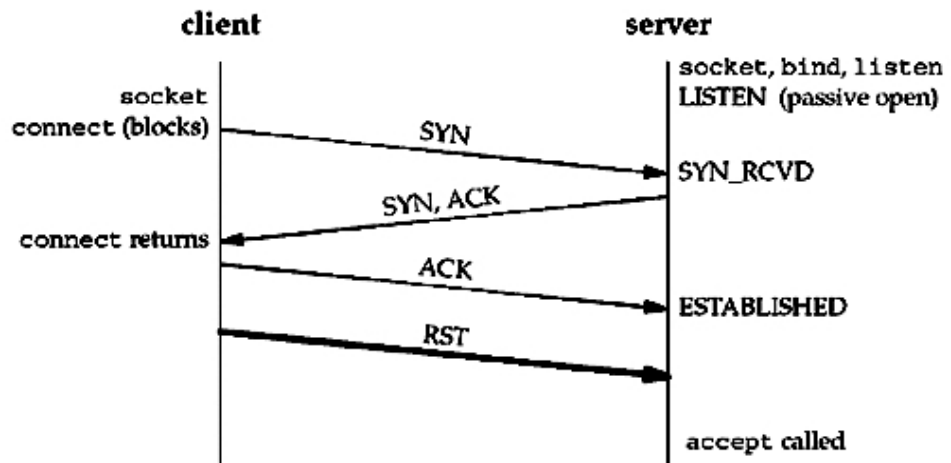
```
#include    "unp.h"
int main(int argc, char **argv)
{
    int    listenfd, connfd;
    pid_t           childpid;
    socklen_t       clilen;
    struct sockaddr_in  cliaddr, servaddr;
    void            sig_chld(int);
    listenfd = Socket(AF_INET, SOCK_STREAM, 0);
    bzero(&servaddr, sizeof(servaddr));
    servaddr.sin_family     = AF_INET;
    servaddr.sin_addr.s_addr = htonl(INADDR_ANY);
    servaddr.sin_port       = htons(SERV_PORT);
    Bind(listenfd, (SA *) &servaddr, sizeof(servaddr));
    Listen(listenfd, LISTENQ);
    Signal(SIGCHLD, sig_chld);  /* must call waitpid() */
    for ( ; ; ) {
        clilen = sizeof(cliaddr);
        if ( (connfd = accept(listenfd, (SA *) &cliaddr, &clilen)) < 0) {
            if (errno == EINTR)
                continue;      /* back to for() */
            else
                err_sys("accept error");
        }
        if ( (childpid = Fork()) == 0) {    /* child process */
            Close(listenfd);   /* close listening socket */
            str_echo(connfd);  /* process the request */
            exit(0);
        }
        Close(connfd);         /* parent closes connected socket */
    }
}
```

The pur pose of this section ha s be en t o de monstrate three s cenarios t hat we can encounter with network programming:

- We must catch the SIGCHLD signal when forking child processes.

- We must handle interrupted system calls when we catch signals.

- A SIGCHLD handler m ust be c oded correctly us ing waitpid to prevent any zombies from being left around.

## Connection Abort before accept Returns

There is a nother c ondition s imilar to the in terrupted s ystem c all th at c an cause accept to re turn a n onfatal error, in w hich c ase w e s hould j ust call accept again. The sequence of packets shown below has been seen on busy s ervers ( typically b usy W eb s ervers), w here t he s erver r eceives an RST for an ESTABLISHED connection before accept is called.



The t hree-way handshake c ompletes, t he c onnection i s e stablished, a nd then t he cl ient T CP s ends an R ST ( reset). On t he s erver s ide, t he connection i s qu eued b y i ts T CP, w aiting f or t he s erver pr ocess t o call accept w hen t he R ST ar rives. S ometime l ater, t he s erver p rocess cal ls accept.

An e asy w ay to s imulate th is s cenario is to s tart th e s erver, h ave it call socket, bind, and listen, and then go to sleep for a short period of time before cal ling accept. W hile t he s erver p rocess i s as leep, s tart t he cl ient and h ave i t cal l socket and connect. As s oon a s connect returns, s et the SO_LINGER socket option to generate the RST and terminate.

## Termination of Server Process

We will now start our c lient/server and then ki ll the s erver child process, which simulates the crashing of the server process. We must be careful to distinguish between the crashing of the server *process* and the crashing of the server *host*.

The following steps take place:

1. We start the server and client and type one line to the client to verify that all is okay. That line is echoed normally by the server child.

2. We find the process ID of the server child and kill it. As part of process termination, all open descriptors in the child are closed. This causes a FIN to be sent to the client, and the client TCP responds with an ACK. This is the first half of the TCP connection termination.

3. The SIGCHLD signal is sent to the server parent and handled correctly.

4. Nothing happens at the client. The client TCP receives the FIN from the server TCP and responds with an ACK, but the problem is that the client process is blocked in the call to fgets waiting for a line from the terminal.

5. Running netstat at this point shows the state of the sockets.

```
linux % netstat -a | grep 9877

tcp     0     0 *:9877            *:*              LISTEN

tcp     0     0 localhost:9877      localhost:43604    FIN_WAIT2

tcp     1     0 localhost:43604     localhost:9877     CLOSE_WAIT
```

6. We can still type a line of input to the client. Here is what happens at the client starting from Step 1:

```
linux % tcpcli01 127.0.0.1  # start client

hello              # the first line that we type

hello                # is echoed correctly we kill the server child on the
server host

another line       # we then type a second line to the client

str_cli : server terminated prematurely
```

When we type "another line," str_cli calls written and the client TCP sends the data to the server. This is allowed by TCP because the receipt of the FIN by the client TCP only indicates that the server process has closed its end of the connection and will not be sending any more data. The receipt of the FIN does not tell the client TCP that the server process has terminated (which in this case, it has).

When the server TCP receives the data from the client, it responds with an RST since the process that had that socket open has terminated. We can verify that the RST was sent by watching the packets with tcpdump.

7. The client process will not see the RST because it calls readline immediately after the call to writen and readline returns 0 (EOF) immediately because of the FIN that was received in Step 2. Our client is not expecting to receive an EOF at this point (str_cli) so it quits with the error message "server terminated prematurely."

8. When the client terminates (by calling err_quit in str_cli), all its open descriptors are closed.

   ❖ If the readline happens before the RST is received (as shown in this example), the result is an unexpected EOF in the client.

   ❖ If the RST arrives first, the result is an ECONNRESET ("Connection reset by peer") error return from readline.

The problem in this example is that the client is blocked in the call to fgets when the FIN arrives on the socket. The client is really working with two descriptors, the socket and the user input. Instead of blocking on input from only one of the two sources, it should block on input from either source.

## SIGPIPE Signal

The rules are:

- When a process writes to a socket that has received an RST, the SIGPIPE signal is sent to the process. The default action of this signal is to terminate the process, so the process must catch the signal to avoid being involuntarily terminated.

- If the process either catches the signal and returns from the signal handler, or ignores the signal, the write operation returns EPIPE.

We can simulate this from the client by performing two writes to the server (which has sent FIN to the client) before reading anything back, with the first write eliciting the RST (causing the server to send an RST to the client). We must use two writes to obtain the signal, because the first write elicits the RST and the second write elicits the signal. It is okay to write to a socket that has received a FIN, but it is an error to write to a socket that has received an RST.

We modify our client as below:

```
#include    "unp.h"
void
str_cli(FILE *fp, int sockfd)
{
  char    sendline[MAXLINE], recvline[MAXLINE];

  while (Fgets(sendline, MAXLINE, fp) != NULL) {

    Writen(sockfd, sendline, 1);
    sleep(1);
    Writen(sockfd, sendline+1, strlen(sendline)-1);

    if (Readline(sockfd, recvline, MAXLINE) == 0)
      err_quit("str_cli: server terminated prematurely");

    Fputs(recvline, stdout);
  }
}
```

The writen is called two times. The intent is for the first writen to elicit the RST and then for the second writen to generate SIGPIPE.

Run the program on the Linux host:

```
linux % tcpclill 127.0.0.1

hi there       # we type this line

hi there       # this is echoed by the server

               # here we kill the server child

bye            # then we type this line

Broken pipe    # this is printed by the shell
```

We s tart t he cl ient, t ype i n o ne l ine, s ee t hat line ech oed co rrectly, an d then t erminate t he s erver c hild on t he s erver hos t. W e t hen t ype a nother line ("bye") and the shell tells us the process died with a SIGPIPE signal.

The r ecommended w ay t o ha ndle SIGPIPE depends on w hat t he application wants to do when this occurs:

- If there is nothing special to do, then setting the signal disposition to SIG_IGN is easy, assuming that subsequent output operations will catch the error of EPIPE and terminate.

- If special actions are needed, when the signal occurs (writing to a log file perhaps), then the signal should be caught and any desired actions can be performed in the signal handler.

- If multiple sockets are in use, the delivery of the signal will not tell us which socket encountered the error. If we need to know which write caused the error, then we must either ignore the signal or return from the signal handler and handle EPIPE from the write.

---

**Check Your Progress**

- What will this program print?

1.  #include<stdio.h>

2.  #include<signal.h>

3.  #include<unistd.h>

4.

5.  void response (int);

6.  void response (int sig_no)

7.  {

8.  printf("%s is working\n",sys_siglist[sig_no]);

9.  }

10.  int main()

11.  {

12.  alarm(5);

13.   sleep(50);

14.  printf("Sanfoundry\n");

15.   signal(SIGALRM,response);

16.  return 0;

17.  }

---

## 4.5 CRASHING AND REBOOTING OF SERVER HOST

This section discusses the case when we will establish a connection between the client and server and then assume the server host crashes and reboots. The easiest way to simulate this is to establish the connection, disconnect the server from the network, shut down the server host and

then reboot it, and then reconnect the server host to the network. We do not want the client to see the server host shut down. The following steps take place:

- We start the server and then the client. We type a line to verify that the connection is established.

- The server host crashes and reboots.

- We type a line of input to the client, which is sent as a TCP data segment to the server host.

- When the server host reboots after crashing, its TCP loses all information about connections that existed before the crash. Therefore, the server TCP responds to the received data segment from the client with an RST.

- Our client is blocked in the call to *readline* when the RST is received, causing *readline* to return the error *ECONNRESET*.

---

**Check your progress**

1.  How are zombies handled?

2.  What happens when a system does not catch SIGTERM signal?

---

## 4.6    SHUTDOWN OF SERVER HOST

When a Unix system is shutdown, the *init* process normally sends the *SIGTERM* signal to all processes (this signal can be caught), waits some fixed amount of time (often between 5 and 20 seconds), and then sends *SIGKILL* signal (which we cannot catch) to any process still running. This gives all running processes a short amount of time to clean up and terminate.

If we do not catch *SIGTERM* and terminate, our server will be terminated by *SIGKILL* signal. When the process terminates, all the open descriptors are closed, and we then follow the same sequence of steps discussed under —termination of server process. We need to select the select or poll function in the client to have the client detect the termination of the server process as soon it occurs.

**Problem**:  To write an algorithm for TCP echo client server

**Solution**: Server-

STEP 1: Start
STEP 2: Declare the variables for the socket
STEP 3: Specify the family, protocol, IP address and port number
STEP 4: Create a socket using socket() function
STEP 5: Bind the IP address and Port number
STEP 6: Listen and accept the client's request for the connection

STEP 7: Read the client's message
STEP 8: Display the client's message
STEP 9: Close the socket
STEP 10: Stop

Client-

STEP 1: Start
STEP 2: Declare the variables for the socket
STEP 3:  Specify the family, protocol, IP address and port number
STEP 4: Create a socket using socket() function
STEP 5: Call the connect() function
STEP 6: Read the input message
STEP 7: Send the input message to the server
STEP 8: Display the server's echo
STEP 9: Close the socket
STEP 10: Stop

## 4.7   SUMMARY

In this unit, we discuss TCP client/server mechanism with concept of T CP e cho s erver. C lient s erver s cenario i s e xplained us ing s ocket programming i n  Linux  host. V arious  conditions t hat m ay o ccur dur ing course of execution of s erver and c lient, s tarting from nor mal start u p to server f ailure.  Unit  also c overs va rious s ignal ha ndling f unction a nd interrupt handling in server.

## 4.8   TERMINAL QUESTIONS

1.    Which IP address for server is used at time of client initialization and why?

2.    Write a code segment for sigaction function to establish disposition of a signal.

3.    Explain the significance of wait and wait_pid along with their limitations.

4.    Discuss the cases of crashing of server and server host.

5.    Distinguish the behavior of server at time of server shut down and server reboot.

6.    What is the output of this program?

    #include<stdio.h>

    #include<signal.h>

    #include<stdlib.h>

    void response (int);

    void response (int sig_no)

```c
{
    printf("%s\n",sys_siglist[sig_no]);
    printf("This is singal handler\n");
}
int main()
{
    pid_t child;
    int status;
    child = fork();
    switch (child){
        case -1 :
            perror("fork");
            exit (1);
        case 0 :
            kill(getppid(),SIGKILL);
            printf("I am an orphan process because my parent has been killed by me\n");
            printf("Handler failed\n");
            break;
        default :
            signal(SIGKILL,response);
            wait(&status);
            printf("The parent process is still alive\n");
            break;
    }
    return 0;
}
```

7.  What is the output of this program?

```c
#include<stdio.h>
#include<signal.h>
void response (int);
void response (int sig_no)
{
    printf("%s\n",sys_siglist[sig_no]);
}
```

```c
int main()
{
    pid_t child;
    int status;
    child = fork();
    switch(child){
        case -1:
            perror("fork");
        case 0:
            break;
        default :
            signal(SIGCHLD,response);
            wait(&status);
            break;
    }
}
```

**Bachelor of Computer Application**

**BCA-E7**
**Network Programming**

Uttar Pradesh Rajarshi Tandon
Open University

Block

# 2

# Course Design Committee

# Course Preparation Committee

# BLOCK INTRODUCTION

**Unit 5:** This unit deals with I/O multiplexing wherein the different I/O models are discussed. It also explains the select function, Batch Input and Buffering, Shutdown Function and Poll Function.

**Unit 6:** The various socket options are discussed in this unit. You will learn about the socket states, Generic socket options, IPV6 socket options, ICMP6 socket options along with TCP socket options.

**Unit 7:** The elementary UDP sockets are briefed in this unit. Echo server functions, lost datagram, lack of flow control in UDP and determining out going interface with UDP are also the sub-topics to be focused.

**Unit 8:** This unit covers the DNS protocol of application layer that is responsible for name and address conversion. The gethost by Name function, Resolver options, Functions and IPV6 support, Uname function and other networking information is also discussed.

# UNIT-5 :  I/O MULTIPLEXING

## Structure

## 5.0     INTRODUCTION

I/O mu ltiplexing is  th e c apability o f h andling mu ltiple  I/O conditions (i.e., input is ready to be read or the descriptor is ready to take more output). There are various situations where I/O multiplexing is being required.

- When a client is handling multiple descriptors.

- When a client is to handle multiple sockets at the same time.

- If a TCP server handles both a listening socket and its connected sockets.

- If a server handles both TCP and UDP.

- If  a s erver ha ndles  multiple s ervices a nd pe rhaps m ultiple protocols.

## 5.1     OBJECTIVE

This unit imparts the basic knowledge of I/O multiplexing.

- The different kinds of I/O model and select function is discussed.

- Buffering and Batch Input are explained.

- Use and working of Shutdown and Poll functions is discussed.

## 5.2 I/O MODELS

There are five basic I/O models that are available in Unix.

- Blocking I/O

- Nonblocking I/O

- I/O multiplexing (select and poll)

- Signal driven I/O (SIGIO)

- Asynchronous I/O (the POSIX aio_ functions)

There are normally two distinct phases for an input operation:

1. Waiting for the data to be ready:

   This involves waiting for data to arrive on the network. When it arrives, it is copied into a buffer within the kernel.

2. Copying the data from the kernel to the process.

   This means copying the (ready) data from the kernel's buffer into our application buffer.

### Blocking I/O Model

The most prevalent model for I/O is the blocking I/O model. By default, all sockets are blocking. The scenario is shown in the figure below:



Figure 5.1: Blocking I/O Model

We u se U DP f or this e xample i nstead of T CP because w ith U DP, t he concept of data being "ready" to read is simple, either an entire datagram has b een r eceived o r i t has n ot. W ith T CP i t g ets m ore co mplicated, as additional variables such as the socket's low-water mark come into play.

We al so r efer t o r ecvfrom as a s ystem c all t o d ifferentiate b etween o ur application and the kernel, regardless of how recvfrom is implemented.

In t he figure above, t he process calls r ecvfrom a nd t he s ystem c all d oes not r eturn unt il t he da tagram a rrives a nd i s c opied i nto our a pplication buffer, o r a n error oc curs. W e s ay t hat t he pr ocess i s bl ocked t he e ntire time from when it calls recvfrom until it returns. When recvfrom returns successfully, our application processes the datagram.

## Nonblocking I/O Model

When a socket is set to be nonblocking, we are telling the kernel "when an I/O op eration t hat I r equest c annot be c ompleted w ithout put ting t he process t o s leep, do no t put t he pr ocess t o s leep, but r eturn an e rror instead". The figure is below:



Figure 5.2: Nonblocking I/O Model

- For t he f irst t hree r ecvfrom, t here i s no da ta t o r eturn a nd t he kernel immediately returns an error of EWOULDBLOCK.

- For t he f ourth t ime w e cal l r ecvfrom, a d atagram is r eady, it is copied i nto our a pplication buf fer, a nd recvfrom r eturns successfully. We then process the data.

When a n a pplication s its i n a l oop c alling r ecvfrom on a nonbl ocking descriptor lik e th is, it is c alled **polling**. T he application i s c ontinually polling the kernel to see if some operation is ready. This is often a waste of C PU t ime, but t his model i s oc casionally encountered, nor mally on systems dedicated to one function.

## I/O Multiplexing Model

With **I/O multiplexing**, w e c all s elect or pol l a nd bl ock i n on e of t hese two s ystem c alls, in stead o f b locking in th e actual I/O s ystem c all. T he figure is a summary of the I/O multiplexing model:



Figure 5.3: I/O Multiplexing Model

We b lock in a c all to s elect, w aiting f or th e d atagram s ocket to b e readable. W hen s elect r eturns t hat t he s ocket i s r eadable, w e t hen c all recvfrom to copy the datagram into our application buffer.

## Multithreading with blocking I/O

Another c losely r elated I/O m odel i s t o use m ultithreading w ith bl ocking I/O. That model very closely resembles the model described above, except that in stead o f u sing s elect to b lock o n mu ltiple f ile d escriptors, th e program uses multiple threads (one per file descriptor), and each thread is then free to call blocking system calls like recvfrom.

## Signal-Driven I/O Model

The **signal-driven I/O model** uses signals, telling the kernel to notify us with the SIGIO signal when the descriptor is ready. The figure is below:



Figure 5.4: Signal-Driven I/O Model

- We first enable the socket for signal-driven I/O and install a signal handler us ing t he s igaction system c all. T he r eturn f rom th is system c all is imme diate a nd o ur p rocess c ontinues; it is n ot blocked.

- When t he d atagram i s r eady t o b e r ead, t he SIGIO s ignal i s generated for our process. We can either:

  ❖ Read t he d atagram f rom t he s ignal h andler by calling recvfrom a nd t hen not ify the m ain l oop t hat t he da ta i s ready to be processed

  ❖ Notify the main loop and let it read the datagram.

The advantage to this model is that we are not blocked while waiting for the datagram to arrive. The main loop can continue executing and just wait to be notified by the signal handler that either the data is ready to process or the datagram is ready to be read.

## Asynchronous I/O Model

Asynchronous I /O is de fined b y t he P OSIX s pecification, a nd va rious differences i n t he *real-time* functions t hat a ppeared i n t he v arious

standards w hich cam e t ogether t o f orm t he current P OSIX s pecification have been reconciled.

These f unctions w ork b y telling t he ke rnel t o s tart t he op eration and t o notify us w hen the entire ope ration ( including the c opy of the da ta from the ke rnel t o our bu ffer) i s c omplete. T he m ain difference be tween t his model and the signal-driven I/O model is that with signal-driven I/O, the kernel t ells us w hen an I/O ope ration c an be i nitiated, but w ith asynchronous I/O, the k ernel t ells us w hen an I/O ope ration is c omplete. See the figure below for example:



Figure 5.5: Asynchronous I/O Model

- We c all a io_read ( the P OSIX asynchronous I/O f unctions be gin with aio_ or lio_) and pass the kernel the following:

    ❖ Descriptor, buffer poi nter, buf fer s ize ( the s ame t hree arguments for read),

    ❖ File offset (similar to lseek),

    ❖ Method to notify us when the entire operation is complete.

    This s ystem c all r eturns imme diately and o ur p rocess is n ot blocked while waiting for the I/O to complete.

- We assume in this example that we ask the kernel to generate some signal when the operation is complete. This signal is not generated

until the data has been copied into our application buffer, which is different from the signal-driven I/O model.

## Comparison of the I/O Models

The main difference between the first four models is the first phase, as the second phase in the first four models is the same: the process is blocked in a call to recvfrom while the data is copied from the kernel to the caller's buffer. Asynchronous I/O, however, handles both phases and is different from the first four. The figure below is a comparison of the five different I/O models.

## Synchronous I/O versus 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.

| blocking | nonblocking | I/O multiplexing | signal-driven I/O | asynchronous I/O | |
|---|---|---|---|---|---|
| initiate | check | check | | initiate | wait for data |
| | check | | | | |
| | check | | | | |
| | check | blocked | | | |
| | check | | | | |
| | check | | | | |
| | check | | | | |
| blocked | check | | | | |
| | check | ready initiate | notification initiate | | |
| | blocked | blocked | blocked | | copy data from kernel to user |
| complete | complete | complete | complete | notification | |

1st phase handled differently, 2nd phase handled the same (blocked in call to recvfrom) | handles both phases

Figure 5.6: I/O Model Comparison

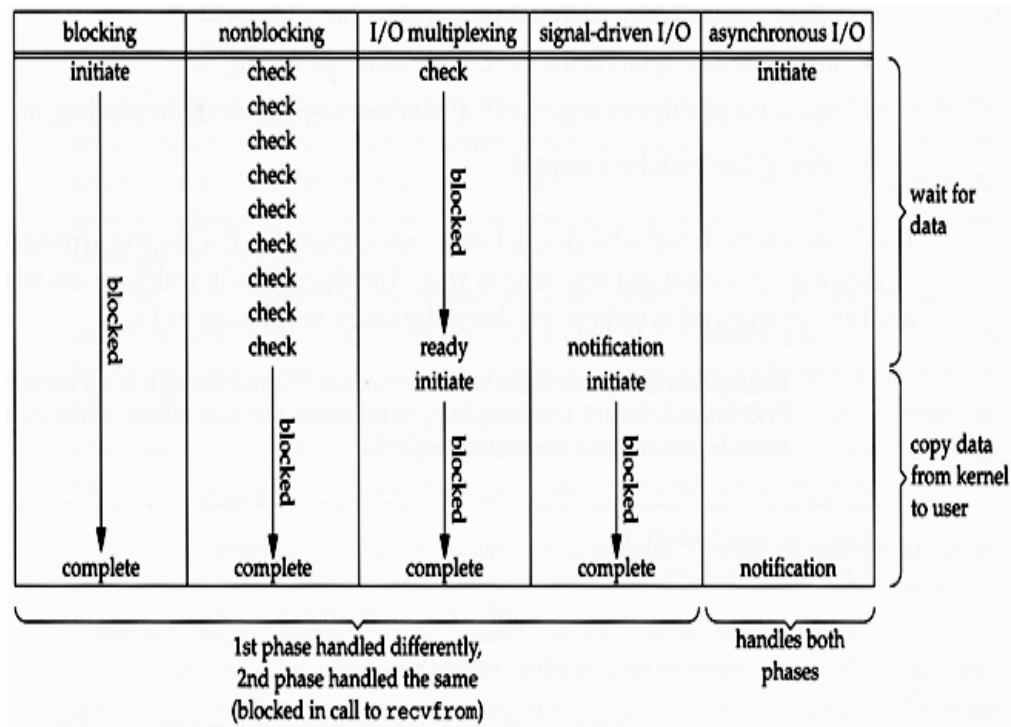Using these definitions, the first four I/O models (blocking, nonblocking, I/O multiplexing, and signal-driven I/O) are all synchronous because the actual I/O operation (recvfrom) blocks the process. Only the asynchronous I/O model matches the asynchronous I/O definition.

## 5.3  SELECT FUNCTION

The select function allows the process to instruct the kernel to either:

- Wait for any one of multiple events to occur and to wake up the process only when one or more of these events occurs, or

- When a specified amount of time has passed.

This means that we tell the kernel what descriptors we are interested in (for reading, writing, or an exception condition) and how long to wait. The descriptors in which we are interested are not restricted to sockets; any descriptor can be tested using select.

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

#include <sys/time.h>

int select(int maxfdp1, fd_set *readset, fd_set *writeset, fd_set *exceptset,
        const struct timeval *timeout);
/* Returns: positive count of ready descriptors, 0 on timeout, −1 on error */
```

### The *timeout* argument

The *timeout* argument tells the kernel how long to wait for one of the specified descriptors to become ready. A timeval structure specifies the number of seconds and microseconds.

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

There are three possibilities for the *timeout*:

1. **Wait forever** (*timeout* is specified as a null pointer). Return only when one of the specified descriptors is ready for I/O.

2. **Wait up to a fixed amount of time** (*timeout* points to a timeval structure). Return when one of the specified descriptors is ready for I/O, but do not wait beyond the number of seconds and microseconds specified in the timeval structure.

3. **Do not wait at all** (*timeout* points to a timeval structure and the timer value is 0, i.e. the number of seconds and microseconds specified by the structure are 0). Return immediately after checking the descriptors. This is called polling.

# Note:

- The w ait in th e f irst tw o s cenarios is n ormally interrupted i f t he process c atches a s ignal a nd r eturns f rom t he s ignal h andler. For portability, w e m ust b e pr epared f or select to r eturn a n e rror of EINTRif w e a re c atching s ignals. Berkeley-derived k ernels never automatically restart select.

- Although the timeval structure has a microsecond field tv_usec, the actual r esolution s upported b y t he ke rnel i s of ten m ore c oarse. Many Unix kernels round the timeout value up to a multiple of 10 ms. There is also a scheduling latency involved, meaning it takes some time after the timer expires before the kernel schedules this process to run.

- On s ome s ystems, t he timeval structure c an r epresent v alues t hat are not s upported by select; it w ill f ail w ith EINVAL if the tv_sec field in the timeout is over 100 million seconds.

- The const qualifier on the *timeout* argument me ans it is n ot modified by select on return.

## The descriptor sets arguments *

The t hree m iddle arguments, *readset*, *writeset*, a nd *exceptset*, s pecify t he descriptors t hat w e w ant t he ke rnel t o t est f or r eading, w riting, a nd exception c onditions. T here are onl y t wo exception c onditions c urrently supported:

- The arrival of out-of-band data for a socket.

- The pr esence of c ontrol s tatus i nformation t o b e read f rom t he master s ide o f a p seudo-terminal t hat ha s be en put i nto pa cket mode. (Not covered in UNP)

select uses descriptor s ets, t ypically an ar ray o f i ntegers, w ith each b it i n each i nteger corresponding t o a de scriptor. For e xample, us ing 32 -bit integers, the first element of the array corresponds to descriptors 0 through 31, the second element of the array corresponds to descriptors 32 t hrough 63, a nd s o on. A ll t he i mplementation de tails a re i rrelevant t o t he application a nd are hi dden i n t he fd_set datatype a nd t he following four macros:

```
void FD_ZERO(fd_set *fdset);        /* clear all bits in fdset */

void FD_SET(int fd, fd_set *fdset);  /* turn on the bit for fd in fdset */

void FD_CLR(int fd, fd_set *fdset);  /* turn off the bit for fd in fdset */

int FD_ISSET(int fd, fd_set *fdset); /* is the bit for fd on in fdset ? */
```

We allocate a descriptor set of the fd_set datatype, we set and test the bits in t he s et us ing t hese m acros, a nd w e c an a lso a ssign i t t o a nother descriptor set across an equals sign (=) in C.

An array of integers using one bit per descriptor, is just one possible way to implement select. Nevertheless, it is common to refer to the individual descriptors within a descriptor set as bits, as in "turn on the bit for the listening descriptor in the read set."

The following example defines a variable of type fd_set and then turns on the bits for descriptors 1, 4, and 5:

```
fd_set rset;

FD_ZERO(&rset);        /* initialize the set: all bits off */

FD_SET(1, &rset);      /* turn on bit for fd 1 */

FD_SET(4, &rset);      /* turn on bit for fd 4 */

FD_SET(5, &rset);      /* turn on bit for fd 5 */
```

It is important to initialize the set, since unpredictable results can occur if the set is allocated as an automatic variable and not initialized.

Any of the middle three arguments to select, *readset*, *writeset*, or *exceptset*, can be specified as a null pointer if we are not interested in that condition. Indeed, if all three pointers are null, then we have a higher precision timer than the normal Unix sleep function. The poll function provides similar functionality.

## The *maxfdp1* argument

The *maxfdp1* argument specifies the number of descriptors to be tested. Its value is the maximum descriptor to be tested plus one. The descriptors 0, 1, 2, up through and including *maxfdp1*−1 are tested.

The constant FD_SETSIZE, defined by including <sys/select.h>, is the number of descriptors in the fd_set datatype. Its value is often 1024, but few programs use that many descriptors.

The reason the *maxfdp1* argument exists, along with the burden of calculating its value, is for efficiency. Although each fd_set has room for many descriptors, typically 1,024, this is much more than the number used by a typical process. The kernel gains efficiency by not copying unneeded portions of the descriptor set between the process and the kernel, and by not testing bits that are always 0.

## *readset*, *writeset*, and *exceptset* as value-result arguments

select modifies the descriptor sets pointed to by the *readset*, *writeset*, and *exceptset* pointers. These three arguments are value-result arguments. When we call the function, we specify the values of the descriptors that we are interested in, and on return, the result indicates which descriptors are ready. We use the FD_ISSET macro on return to test a specific descriptor in an fd_set structure. Any descriptor that is not ready on return will have its corresponding bit cleared in the descriptor set. To handle this, we turn on all the bits in which we are interested in all the descriptor sets each time we call select.

## Return value of select

The return value from this function indicates the total number of bits that are ready across all the descriptor sets. If the timer value expires before any of the descriptors are ready, a value of 0 is returned. A return value of −1 indicates an error (which can happen, for example, if the function is interrupted by a caught signal).

---

**Check your progress**

1.     Explain nonblocking I/O model.

2.     What are the possibilities for timeouts?

---

# 5.4   BATCH INPUT AND BUFFERING

If we consider the network between the client and server as a full-duplex pipe, with requests going from the client to the server and replies in the reverse direction, then the following figure shows our stop-and-wait mode:
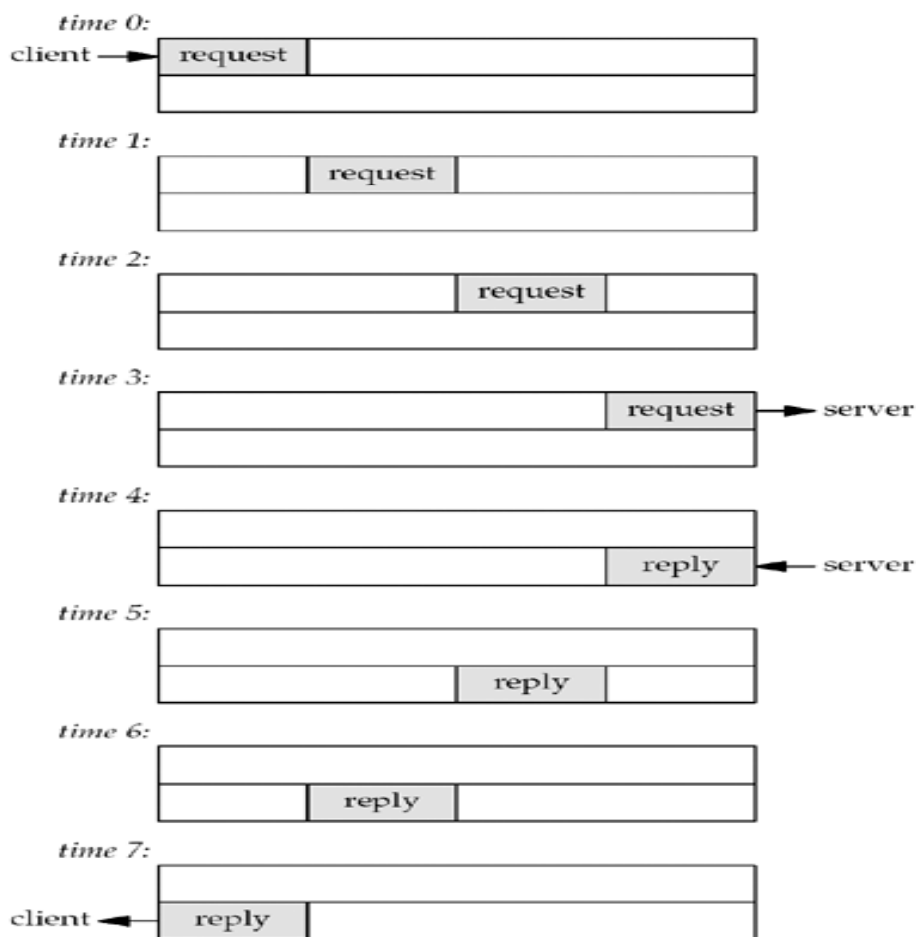


Figure 5.7: Stop and Wait Mode

Note that this figure:

- Assumes that there is no server processing time and that the size of the request is the same as the reply

- Shows s how onl y t he da ta pa ckets, i gnoring t he T CP acknowledgments that are also going across the network

A request is sent by the client at time 0 and we assume an RTT of 8 units of time. The reply sent at time 4 is received at time 7.

This s top-and-wait mo de is f ine f or in teractive in put. T he p roblem is : if we run our client in a batch mode, when we redirect the input and output, however, the resulting output file is always smaller than the input file (and they should be identical for an echo server).

## Batch mode

To s ee w hat's h appening, r ealize t hat i n a b atch m ode, w e can k eep sending r equests a s fast as t he n etwork can accept t hem. T he s erver processes them and sends back the replies at the same rate. This leads to the full pipe at time 7, as shown below:

*Time 7:*

| request 8 | request 7 | request 6 | request 5 |
|-----------|-----------|-----------|-----------|
| reply 1   | reply 2   | reply 3   | reply 4   |

*Time 8:*

| request 9 | request 8 | request 7 | request 6 |
|-----------|-----------|-----------|-----------|
| reply 2   | reply 3   | reply 4   | reply 5   |

Figure 5.8: Batch Mode

We assume:

- After s ending th e f irst r equest, w e imme diately s end a nother, a nd then another

- We can k eep s ending requests as f ast as t he n etwork can accept them, along with processing replies as fast as the network supplies them.

Assume that the input file contains only nine lines. The last line is sent at time 8, as shown in the above figure. But we cannot close the connection after writing this request because there are still other requests and replies in the pipe. The cause of the problem is our handling of an EOF on input: The function returns to the main function, which then terminates. But in a batch mode, an EOF on input does not imply that we have finished reading from the socket; there might still be requests on the way to the server, or replies on the way back from the server.

The solution is to close one-half of the TCP connection by sending a FIN to the server, telling it we have finished sending data, but leave the socket descriptor op en f or r eading. T his i s don e w ith t he shutdown f unction, described in the next section.

## Buffering concerns

When s everal l ines of i nputs a re a vailable f rom t he s tandard i nput select will cause the code to read the input using fgets which will read the available lines into a buffer used by stdio. But, fgets only returns a single line a nd le aves a ny remaining d ata s itting in th e s tdio buf fer. T he following code writes that single line to the server and then select is called again to wait for more work, even if there are additional lines to consume in the stdio buffer. The reason is that select knows nothing of the buffers used b y s tdio; it w ill only s how r eadability from th e v iewpoint o f the read system c all, n ot c alls lik e fgets. T hus, mixing s tdio a nd select is considered very error-prone and should only be done with great care.

The s ame p roblem e xists w ith readline in th is e xample (str_cli function). Instead of da ta be ing h idden f rom select in a s tdio buf fer, i t i s hi dden in readline's b uffer. A function th at g ives v isibility in to readline's buf fer, so one pos sible s olution is to modify our c ode to use that function before calling select to s ee i f da ta ha s a lready b een read but not c onsumed. But again, the complexity grows out of hand quickly when we have to handle the case where the readline buffer contains a partial line (meaning we still need to read more) as well as when it contains one or more complete lines (which we can consume).

## pselect Function

The pselect function w as i nvented b y P OSIX a nd i s now s upported b y many of the Unix variants.

```
#include <sys/select.h>
#include <signal.h>
#include <time.h>
int pselect (int maxfdp1, fd_set *readset, fd_set *writeset, fd_set *exceptset,
        const struct timespec *timeout, const sigset_t *sigmask);
/* Returns: count of ready descriptors, 0 on timeout, −1 on error */
```

pselect contains two changes from the normal select function:

1. pselect uses t he timespec structure ( another P OSIX i nvention) instead of the timeval structure. The tv_nsec member o f the n ewer structure s pecifies n anoseconds, w hereas t he tv_usec member o f the older structure specifies microseconds.

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

2.  pselect adds a sixth argument: a pointer to a signal mask. This allows the program to disable the delivery of certain signals, test some global variables that are set by the handlers for these now-disabled signals, and then call pselect, telling it to reset the signal mask.

With regard to the second point, consider the following example:

Our program's signal handler for SIGINT just sets the global intr_flag and returns. If our process is blocked in a call to select, the return from the signal handler causes the function to return with errno set to EINTR. But when select is called, the code looks like the following:

```
if (intr_flag)
   handle_intr();     /* handle the signal */
/* signals occurring in here are lost */
if ( (nready = select( ... )) < 0) {
   if (errno == EINTR) {
     if (intr_flag)
        handle_intr();
   }
   ...
}
```

The problem is that between the test of intr_flag and the call to select, if the signal occurs, it will be lost if select blocks forever.

With pselect, we can now code this example reliably as:

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

```
        handle_intr ();

   }

   ...

  }
```

Before t esting t he intr_flag variable, w e b lock SIGINT. W hen pselect is called, i t r eplaces t he s ignal m ask o f t he p rocess w ith an em pty s et (i.e., zeromask) a nd t hen c hecks t he de scriptors, possibly going t o s leep. But w hen pselect returns, t he s ignal m ask of t he pr ocess i s r eset t o i ts value before pselect was called (i.e., SIGINT is blocked).

## 5.5   SHUTDOWN FUNCTION

The n ormal w ay to terminate a n etwork connection is to c all the close function. But, there are tw o limita tions w ith close that c an be avoided with shutdown:

1.    close decrements t he d escriptor's r eference co unt an d cl oses t he socket only if the count reaches 0. W ith shutdown, we c an initiate TCP's normal c onnection termination sequence (the four segments beginning with a FIN in), regardless of the reference count.

2.    close terminates b oth d irections o f d ata t ransfer, r eading an d writing. S ince a TCP c onnection i s f ull-duplex, t here a re t imes when we want to tell the other end that we have finished sending, even though that end might have more data to send us. This is the scenario we encountered in the previous section with batch input to our str_cli function. T he f igure be low s hows t he t ypical f unction calls in this scenario.
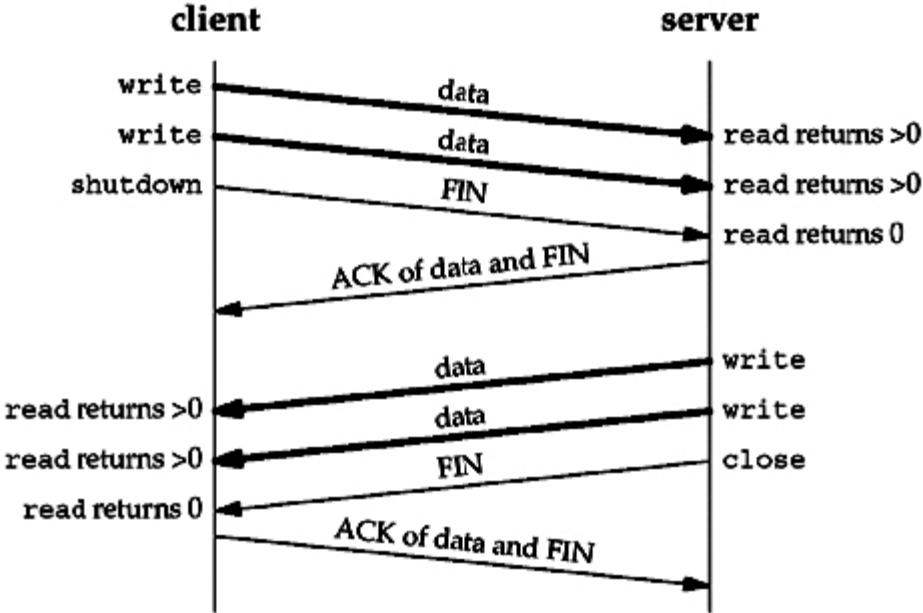


Figure 5.9: Network Termination using shutdown Function

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


int shutdown(int sockfd, int howto);


/* Returns: 0 if OK, −1 on error */
```

The action of the function depends on the value of the *howto* argument:

- SHUT_RD: **The read half of the connection is closed.** No more data can be received on the socket and any data currently in the socket receive buffer is discarded. The process can no longer issue any of the read functions on the socket. Any data received after this call for a TCP socket is acknowledged and then silently discarded.

- SHUT_WR: **The write half of the connection is closed.** In the case of TCP, this is called a **half-close**. Any data currently in the socket send buffer will be sent, followed by TCP's normal connection termination sequence. As we mentioned earlier, this closing of the write half is done regardless of whether or not the socket descriptor's reference count is currently greater than 0. The process can no longer issue any of the write functions on the socket.

- SHUT_RDWR: **The read half and the write half of the connection are both closed.** This is equivalent to calling shutdown twice: first with SHUT_RD and then with SHUT_WR.

The three SHUT_*xxx* names are defined by the POSIX specification. Typical values for the howto argument that you will encounter will be 0 (close the read half), 1 (close the write half), and 2 (close the read half and the write half).

---

## Check your progress

1. What are the possible values of *howto* argument?

2. Explain *pselect* function?

---

## 5.6   POLL FUNCTION

Poll provides functionality that is similar to select, but poll provides additional information when dealing with STREAMS devices.

```
#include <poll.h>
int poll (struct pollfd *fdarray, unsigned long nfds, int timeout);
/* Returns: count of ready descriptors, 0 on timeout, −1 on error */
```

Arguments:

The first argument (*fdarray*) is a pointer to the first element of an array of structures. Each element is apollfd structure that specifies the conditions to be tested for a given descriptor, fd.

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

The conditions to be tested are specified by the events member, and the function returns the status for that descriptor in the corresponding revents member. This data structure (having two variables per descriptor, one a value and one a result) avoids value-result arguments (the middle three arguments forselect are value-result). Each of these two members is composed of one or more bits that specify a certain condition. The following figure shows the constants used to specify the events flag and to test the revents flag.

| Constant | events | revents | Description |
| --- | --- | --- | --- |
| POLLIN | x | x | normal or priority band to read |
| POLLRDNORM | x | x | normal data to read |
| POLLRDBAND | x | x | priority band data to read |
| POLLPRI | x | x | high-priority data to read |
| POLLOUT | x | x | normal data can be written |
| POLLWRNORM | x | x | normal data can be written |
| POLLWRBAND | x | x | priority band data can be written |
| POLLERR | | x | an error has occurred |
| POLLHUP | | x | hangup has occurred |
| POLLNVAL | | x | descriptor is not an open file |

Figure 5.10: Summary of Constants specifying events and revents flags

The first four constants deal with input, the next three deals with output, and the final three deals with errors. The final three cannot be set in events, but are always returned in revents when the corresponding condition exists.

With regard to TCP and UDP sockets, the following conditions cause poll to return the specified revent. Unfortunately, POSIX leaves many holes (optional ways to return the same condition) in its definition ofpoll.

- All regular TCP data and all UDP data is considered normal.

- TCP's out-of-band data is considered priority band.

- When the read half of a TCP connection is closed (e.g., a FIN is received), this is also considered normal data and a subsequent read operation will return 0.

- The presence of an error for a TCP connection can be considered either normal data or an error (POLLERR). In either case, a subsequent read will return −1 with errno set to the appropriate value. This handles conditions such as the receipt of an RST or a timeout.

- The availability of a new connection on a listening socket can be considered either normal data or priority data. Most implementations consider this normal data.

- The completion of a nonblocking connect is considered to make a socket writable.

The number of elements in the array of structures is specified by the *nfds* argument.

The *timeout* argument specifies how long the function is to wait before returning. A positive value specifies the number of milliseconds to wait. The constant INFTIM (wait forever) is defined to be a negative value.

Return values from poll:

- −1 if an error occurred

- 0 if no descriptors are ready before the timer expires

- Otherwise, it is the number of descriptors that have a nonzero revents member.

If we are no longer interested in a particular descriptor, we just set the fd member of the pollfd structure to a negative value. Then the events member is ignored and the revents member is set to 0 on return.

This section is discusses the TCP echo server from using poll instead of select.

In the select version we allocate a client array along with a descriptor set named rset. With poll, we must allocate an array of pollfd structures to maintain the client information instead of allocating another array. We handle the fd member of this array the same way we handled the client array in the selection version: a value of −1 means the entry is not in use; otherwise, it is the descriptor value. Any entry in the array

of pollfd structures passed to poll with a negative value for the fd member is just ignored.

```
/* include fig01 */
#include    "unp.h"
#include    <limits.h>      /* for OPEN_MAX */
Int main(int argc, char **argv)
{
    int             i, maxi, listenfd, connfd, sockfd;
    int             nready;
    ssize_t         n;
    char            buf[MAXLINE];
    socklen_t       clilen;
    struct pollfd    client[OPEN_MAX];
    struct sockaddr_in  cliaddr, servaddr;
    listenfd = Socket(AF_INET, SOCK_STREAM, 0);
    bzero(&servaddr, sizeof(servaddr));
    servaddr.sin_family     = AF_INET;
    servaddr.sin_addr.s_addr = htonl(INADDR_ANY);
    servaddr.sin_port       = htons(SERV_PORT);
    Bind(listenfd, (SA *) &servaddr, sizeof(servaddr));
    Listen(listenfd, LISTENQ);
    client[0].fd = listenfd;
    client[0].events = POLLRDNORM;
    for (i = 1; i < OPEN_MAX; i++)
        client[i].fd = -1;     /* -1 indicates available entry */
    maxi = 0;                 /* max index into client[] array */
/* end fig01 */
/* include fig02 */
    for ( ; ; ) {
        nready = Poll(client, maxi+1, INFTIM);
        if (client[0].revents & POLLRDNORM) {   /* new client connection */
            clilen = sizeof(cliaddr);
            connfd = Accept(listenfd, (SA *) &cliaddr, &clilen);
#ifdef  NOTDEF
            printf("new client: %s\n", Sock_ntop((SA *) &cliaddr, clilen));
```

```c
#endif
        for (i = 1; i < OPEN_MAX; i++)
          if (client[i].fd < 0) {
              client[i].fd = connfd;  /* save descriptor */
              break;
          }
        if (i == OPEN_MAX)
          err_quit("too many clients");
        client[i].events = POLLRDNORM;
        if (i > maxi)
          maxi = i;           /* max index in client[] array */
        if (--nready <= 0)
          continue;           /* no more readable descriptors */
    }
    for (i = 1; i <= maxi; i++) {   /* check all clients for data */
        if ( (sockfd = client[i].fd) < 0)
          continue;
        if (client[i].revents & (POLLRDNORM | POLLERR)) {
          if ( (n = read(sockfd, buf, MAXLINE)) < 0) {
              if (errno == ECONNRESET) {
                  /* connection reset by client */
#ifdef  NOTDEF
                  printf("client[%d] aborted connection\n", i);
#endif
                  Close(sockfd);
                  client[i].fd = -1;
              } else
                  err_sys("read error");
          } else if (n == 0) {
              /* connection closed by client */
#ifdef  NOTDEF
              printf("client[%d] closed connection\n", i);
#endif
              Close(sockfd);
              client[i].fd = -1;
```

```
        } else

            Writen(sockfd, buf, n);

        if (--nready <= 0)

            break;           /* no more readable descriptors */

        }

    }

  }

}
```

This code does the following:

- **Allocate array of pollfd structures.**
  We declare OPEN_MAX elements in our array of pollfd
  structures. Determining the maximum number of descriptors
  that a process can have opened at any one time is difficult. One
  way is to call the POSIX sysconf function with
  an argument of _SC_OPEN_MAX (as described in
  APUE) and then dynamically allocate an array of the
  appropriate size.

- **Initialize.** We u se th e f irst e ntry in  th e client array fo r t he
  listening socket and set the descriptor for the remaining entries to –
  1. We also set the POLLRDNORM event for this descriptor, to be
  notified b ypoll when a  ne w c onnection i s r eady t o b e a ccepted.
  The v ariable maxi contains t he l argest i ndex o f t he client array
  currently in use.

- **Call poll, check for new connection.** We cal l poll to w ait
  for either a new connection or data on existing connection.

  ❖ When a  ne w c onnection  is accep ted, w  e f ind t he f irst
    available entry in the client array by looking for the first one
    with a negative descriptor.

  ❖ We start the search with the index of 1, since client[0] is used
    for the listening socket.

  ❖ When an available entry is found, we save the descriptor and
    set the POLLRDNORM event.

- **Check for data on an existing connection.** The two return
  events that we check for are POLLRDNORM and POLLERR. We
  did not s et POLLERR in t he ev ents m ember b ecause i t i s al ways
  returned w hen t  he c ondition i s  true. T  he r eason w  e ch  eck
  for POLLERR is be cause s ome i mplementations r eturn t his e vent

when a n R ST i s r eceived f or a c onnection, while ot hers j ust return POLLRDNORM. In either case, we call read and if an error has oc curred, i t will return an error. When an existing connection is terminated by the client, we just set the fd member to –1.

---

**Check your progress**

1. Explain w ith f igure s top a nd w ait m ode f or r equest r eply messages between client and server in full duplex network.

2. What are the possible return values of the timeout argument and what do they mean?

---

## 5.7 SUMMARY

In this unit, we study about the five different models in Unix for I/O:

- Non-blocking

- Blocking

- I/O multiplexing

- Signal-driven I/O

- Asynchronous I/O

We s ee t hat t he B locking I/O i s t he p revalently applied de fault. It i s a lso observed t hat t he P OSIX s pecification i s w idely used for de fining t rue asynchronous I/O. The s elect f unction i s us ed for I/O m ultiplexing. T he descriptors, th e ma ximum w aiting time a long with th e ma ximum descriptor number incremented by one are provided to the select function. Readability i s s pecified b y t he calls t o s elect. It i s al so o bserved t hat arrival of out of band da ta is the only exception that arises during socket processes. T he s elect f unction d ictates a l imit o n t he l ength of t ime f or which a block in a function persists. This salient feature can be applied to administer the time limit length for input operations. Similar functionality is also provided b y the poll function. It also describes information related to S TREAM de vices. T hough, t he s elect f unction a s w ell a s t he pol l function is necessary for POSIX but the select function is preferably used in most cases.

## 5.8 TERMINAL QUESTIONS

1. Define I/O multiplexing. Under which circumstances is it used?

2. Compare the different I/O multiplexing models.

3.  What are the five functions used to perform file I/O on a Unix System? Elaborate each function with example.

4.  What is timeout argument? Explain the *timeval* structure.

5.  Describe the steps involved in C language while as signing a descriptor set to another one across the equals sign when the descriptor set is an array of integers.

6.  Differentiate between select and poll functions.

7.  What is the consequence when the second argument provided to shutdown is SHUT_RD?

8.  Describe the conditions under which an application calls shutdown using the argument of SHUT_RDWR as an alternative to simply calling close.

# UNIT-6 :  SOCKET OPTIONS

## Structure

## 6.0    INTRODUCTION

A socket is an endpoint of a connection across a computer network which i s r esponsible t o de liver da ta  packet t o ap propriate p rocess o r thread. It i s a  c ombination of  IP a ddress a nd p ort num ber.  Sockets ar e communication poi nts on t he s ame or  di fferent c omputers t o ex change data. T hese  are s upported b y  Unix,  Windows, M ac, a nd m any ot her operating s ystems. T o be  m ore pr ecise, i t's a  w ay t o t alk t o ot her computers using standard Unix file descriptors. In Unix, every I/O action is done by writing or reading a file descriptor. A file descriptor is just an integer associated with an open file and it can be a network connection, a text file, a terminal, or something else. To a programmer, a socket looks and be haves m uch l ike a  l ow-level f ile d escriptor. T his i s b ecause commands such as read() and write() work with sockets in the same way they do with files and pipes.

## 6.1    OBJECTIVE

To cr eate aw areness ab out t he d ifferent s ocket o ptions av ailable for the development of application.

- Firstly getsockop and setsockopt functions are discussed, then the different s ocket s tates a nd g eneric s ocket o ptions ar e ex plained further.

- This unit also sheds light on socket options for IPV6 and ICMP6. Lastly TCP socket options are specified and the chapter terminates with summary.

## 6.2 GETSOCKOPT AND SETSOCKOPT FUNCTION

There are various ways to control a socket:

**getsockopt**() is used to retrieve options associated with the socket. If an option is to be interpreted by the TCP protocol, protocolLevel is set to the TCP protocol number. The parameters optionValuePtr and optionLengthPtr is used to identify a buffer in which the value(s) for the requested option(s) are to be returned. The socket in use may require the process to have appropriate privileges to use the getsockopt() function. The option_name argument specifies a single option to be retrieved.

It can be one of the following values which have been defined in *<sys/socket.h>*:

### Socket Level Options

The following options are recognized at the socket level:

| *Protocol Level* Options | Data Type | Description |
|---|---|---|
| SO_BINDTODEVICE | string | The device name, as set with tfAddInterface(), will be stored as a null-terminated string in the buffer pointed to by optionValuePtr. |
| SO_DONTROUTE | int | Enable/disable routing bypass for outgoing messages. **Default 0.** |
| SO_ERROR | int | Retrieve the socket error. ***This option is for getsockopt() only!*** |
| SO_KEEPALIVE | int | Enable/disable keep connections alive. **Default 0 (disable).** |
| SO_LINGER | linger | Linger on close if data is present. **Default ON with a linger time of 60 seconds.** |
| SO_OOBINLINE | int | Enable/disable reception of out-of-band data in band. **Default is 0.** |

| | | |
|---|---|---|
| SO_RCVBUF | unsigned long | The buffer size for input. **Default is 8192 bytes.** |
| SO_RCVLOWAT | unsigned long | The l ow water m ark f or receiving in bytes. **Default value is 1.** |
| SO_REUSEADDR | int | Enable this socket option to bind the same port number to multiple sockets us ing di fferent local IP addresses. **Default 0 (disable).** |
| SO_REUSEPORT | int | Enable this socket option to bind the s ame l ocal IP ad dress an d port t o m ultiple s ockets. If multiple U DP s ockets ha ve t he SO_REUSEPORT option s et, then t hose s ockets c an bi nd t o the s ame l ocal IP ad dress, an d local UDP port. **Default 0 (disable).** |
| SO_SNDBUF | unsigned long | The buffer size for output. **Default is 8192 bytes.** |
| SO_SNDLOWAT | unsigned long | The low water mark for sending in bytes. **Default value is 2048.** |
| TM_SO_RCVCOPY | unsigned int | **TCP socket:** fraction u se o f a receive b uffer b elow w hich w e try and a ppend t o a previous receive b uffer i n t he s ocket receive queue. **UDP socket:** fraction u se o f a receive b uffer b elow w hich w e try an d co py t o a n ew r eceive buffer, i f there i s already at least a buffer in the receive queue. **Default value is 4 (25%).** |
| TM_SO_SNDAPPEND | unsigned int | **TCP socket only**. T hreshold i n bytes of s end buf fer be low, which w e t ry and a ppend t o t he previous s end buf fer i n t he T CP send que ue. O nly us ed w ith <u>send</u>(), not w ith <u>tfZeroCopySend</u>(). **Default value is 128 bytes.** |

| | | |
|---|---|---|
| TM_SO_SND_DGRAMS | unsigned int | The num ber o f n on-TCP datagrams that can be queued for send on a socket. **Default is 8 datagrams.** |
| TM_SO_RCV_DGRAMS | unsigned int | The number of non-TCP datagrams that can be queued for receive on a socket. **Default is 8 datagrams.** |
| SO_UNPACKEDDATA | int | **TI C3x and C5x DSP platforms only**. If this option is enabled, all socket data will be sent and received in byte unpacked format. If this option is disabled, all socket data will be sent in a byte packed format, as received from the network. **Default 0 (disable)** |

Table 6.1: Socket level options

## IP Level Options

The following options are recognized at the IP level:

| *protocolLevel* Options | Data Type | Description |
|---|---|---|
| IPO_HDRINCL | int | This is a toggle option used on raw s ockets onl y. If t he va lue is non -zero, it in structs th e Treck s tack t hat t he u ser i s including t he IP he ader w hen sending data. **Default 0** |
| IPO_RCV_TOS | unsigned char | Received IP t ype o f s ervice on the connection (from the last IP datagram arrived on t he connection.) |
| IPO_TOS | unsigned char | IP type of service. **Default 0** |
| IPO_TTL | unsigned char | IP Time To Live in seconds. **Default 64** |
| IPO_SRCADDR | <u>ttUserIpAd dress</u> | Set t he IP s ource a ddress f or the connection. **Default: The first multi-home IP address on the outgoing** |

| | | interface |
|---|---|---|
| IPO_MULTICAST_TTL | unsigned char | Change the default IP TTL for outgoing multicast datagrams. |
| IPO_MULTICAST_IF | in_addr | Specify a configured IP address that will uniquely identify the outgoing interface for multicast datagrams sent on this socket. A zero IP address parameter indicates that we want to reset a previously set outgoing interface for multicast packets sent on that socket. |
| IPO_ADD_MEMBERSHIP | ip_mreq | Add group multicast IP address to given interface (see struct ip_mreq data type). |
| IPO_DROP_MEMBERSHIP | ip_mreq | Delete group multicast IP address to given interface. |
| IP_BLOCK_SOURCE | ip_mreq_source | Block data from a given source to a given multicast group (mute). |
| IP_UNBLOCK_SOURCE | ip_mreq_source | Unblock data from a given source to a given multicast group (un-mute). |
| IP_ADD_SOURCE_MEMBERSHIP | ip_mreq_source | Join a source-specific group. |
| IP_DROP_SOURCE_MEMBERSHIP | ip_mreq_source | Leave a source-specific group. |
| MCAST_JOIN_GROUP | group_req | Add group multicast IP address to given interface. This option also supports IPPROTO_IPV6. |
| MCAST_LEAVE_GROUP | group_req | Delete group multicast IP address to given interface. This option also supports IPPROTO_IPV6. |
| MCAST_BLOCK_SOURCE | group_source_req | Block data from a given source to a given multicast group (mute). This option also supports IPPROTO_IPV6. |
| MCAST_UNBLOCK_SOURCE | group_source_req | Unblock data from a given source to a given multicast |

| | | group ( un-mute). T his opt ion also supports IPPROTO_IPV6. |
|---|---|---|
| MCAST_JOIN_SOURCE _GROUP | group_sour ce_req | Join a s ource-specific g roup. This opt ion a lso s upports IPPROTO_IPV6. |
| MCAST_LEAVE_SOUR CE_GROUP | group_sour ce_req | Leave a source-specific group. This option also supports IPPROTO_IPV6. |
| IP_RCV_TOS | unsigned char | Retrieve the IP header TOS from a packet on a TCP connection, after the TCP connection has been established. |

Table 6.2: IP level options

## 6.3   SOCKET STATES

The following socket options are inherited by a connected TCP socket
from the listening socket:

* SO_DEBUG: It i s a   boolean  option w hich r eports w hether
  debugging information is being recorded.

* SO_ACCEPTCONN: It is a bool ean opt ion w hich r eports s ocket
  listening has been enabled.

* SO_BROADCAST: It is boolean option to report that transmission
  of broadcast messages is being supported by the protocol.

* SO_REUSEADDR: It i s a  bool ean opt ion w hich r eports w hether
  the r ules us ed i n va lidating a ddresses s upplied t o  bind()  should
  allow reuse of local addresses.

* SO_KEEPALIVE: It r eports w hether  connections ar e k ept  active
  with pe riodic t ransmission of m essages.  If t he  connected s ocket
  fails to respond to these messages, the connection shall be broken
  and threads writing to that socket shall be notified with a SIGPIPE
  signal. This opt ion s hall store an int value. This i s also a bool ean
  option.

* SO_LINGER: It r eports w hether t he s ocket l ingers on  close()  if
  data is p resent. If S O_LINGER is s et, th e s ystem s hall b lock th e
  calling t hread dur ing close() until it c an tr ansmit th e d ata or u ntil
  the e nd of  t he i nterval i ndicated b y t he  *l_linger*  member,
  whichever comes first. If SO_LINGER is not specified, and close()
  is i ssued, t he s ystem h andles t he c all in  a  w ay  that a llows th e
  calling thread to  continue as quickly as possible. This option shall
  store a linger structure.

- SO_OOBINLINE: R eports w hether t he s ocket l eaves received out-of-band data (data marked urgent) inline. This option shall store an int value. This is a Boolean option.

- SO_SNDBUF: R eports send buf fer s ize i nformation. T his option shall store an int value.

- SO_RCVBUF: R eports r eceive buf fer s ize i nformation. This option shall store an int value.

- SO_ERROR: R eports i nformation a bout e rror status a nd clears it. This option shall store an int value.

- SO_TYPE: Reports the socket type. This option shall store an int value.

- SO_DONTROUTE: R eports w hether out going m essages bypass the standard routing facilities. The destination shall be on a di rectly-connected ne twork, and m essages a re directed t o t he appropriate ne twork i nterface a ccording t o the de stination a ddress. T he e ffect, i f a ny, of t his opt ion depends on w hat protocol is in use. This option shall store an int value. This is a Boolean option.

- SO_RCVLOWAT: R eports the min imum number of b ytes to process for socket input operations. The default value for SO_RCVLOWAT i s 1 . If S O_RCVLOWAT i s s et t o a larger value, blocking receive calls normally wait until they have r eceived t he s maller of t he l ow w ater m ark v alue o r the r equested a mount. ( They m ay r eturn l ess t han the l ow water mark if an error occurs, a signal is caught, or the type of d ata n ext in t he r eceive q ueue i s d ifferent f rom t hat returned; for example, o ut-of-band d ata.) This opt ion s hall store an int value. N ote th at n ot a ll i mplementations a llow this option to be retrieved.

- SO_RCVTIMEO: R eports t he t imeout va lue for i nput operations. T his opt ion s hall s tore a timeval structure w ith the num ber of s econds a nd m icroseconds s pecifying t he limit on how l ong t o w ait f or a n i nput ope ration t o complete. If a r eceive operation has blocked for this much time without receiving additional data, it shall return with a partial co unt o r *errno* set t o [ EAGAIN] o r [EWOULDBLOCK] i f no d ata w as r eceived. T he d efault for t his opt ion is zero, w hich i ndicates t hat a r eceive operation s hall not t ime out . N ote t hat not a ll implementations allow this option to be retrieved.

## 6.4  GENERIC SOCKET OPTIONS

Protocol-independent c ode ( or not b y any existing pr otocol module) within the kernel are used to handle these generic socket options.

Some o ptions ar e s ocket t ype s pecific. F or ex ample, the SO_BROADCAST socket o ption is c alled "g eneric," w hich is u sed only for datagram sockets.

## SO_BROADCAST Socket Option

This option controls the ability of the process to send broadcast messages. Only datagram sockets support broadcasting and networks that support the concept of a broadcast message (e.g., Ethernet, token ring, etc.).

Applications that doesn't support broadcast mechanism are not allowed to do so because applications have to set this socket option before initializing broadcast. For example, a U DP application m ight t ake the de stination IP address as a command-line a rgument, but the a pplication ne ver i ntended for a us er t o t ype i n a br oadcast a ddress. Rather t han f orcing t he application to try to determine if a given address is a broadcast address or not, t he t est i s i n t he k ernel: If t he de stination a ddress i s a b roadcast address and this socket option is not set, EACCES is returned.

## SO_DEBUG Socket Option

This option is s upported only by TCP. When enabled for a TCP s ocket, the kernel keeps track of detailed information about all the packets sent or received by TCP for the socket. These are kept in a circular buffer_within the kernel that can be examined with the trpt program.

## SO_DONTROUTE Socket Option

This opt ion s pecifies t hat out going p ackets a re t o b ypass t he n ormal routing m echanisms of t he unde rlying pr otocol. The de stination m ust be on a d irectly-connected n etwork, an d m essages ar e d irected t o t he appropriate ne twork i nterface according t o t he de stination a ddress. F or example, i n cas e of IPv4 pa ckets a re routed t hrough uni que l ocal interfaces and if the interface is not found, ENETUNREACH is returned.

The e quivalent of this opt ion c an also be a pplied t o individual da tagrams using t he MSG_DONTROUTE flag with th e send, sendto or sendmsg functions. T his opt ion i s of ten us ed b y r outing da emons (e.g., routed and gated) to bypass the routing table and force a packet to be sent out a particular interface.

## SO_ERROR Socket Option

This opt ion is one t hat can be f etched but c annot be s et. When a n error occurs on a socket, the protocol module in a Berkeley-derived kernel sets a variable namedso_error for that socket to one of the standard Unix Exxx values. This is called the *pending error* for the socket. The process can be immediately notified of the error in one of two ways:

1.  If the process is blocked in a call to select on the socket, for either readability o r w ritability, select returns w ith e ither or bot h conditions set.

2. If t he p rocess i s us ing s ignal-driven I/O, t he SIGIO signal is generated for either the process or the process group.

The pr ocess can t hen obt ain t he va lue of so_error by f etching the SO_ERROR socket opt ion. T he i nteger va lue r eturned by getsockopt is the pending error for the socket. The value of so_error is then reset to 0 by the kernel.

- If so_error is nonz ero w hen t he p rocess calls read and t here i s no data t o r eturn, read returns –1 w ith errno set t o t he v alue of so_error. The value of so_error is then reset to 0. If there is data queued for the s ocket, t hat da ta i s r eturned b y read instead of t he error condition.

- If so_error is nonz ero w hen the pr ocess c alls write, –1 i s r eturned with errno set to the value of so_error and so_error is reset to 0.

## SO_KEEPALIVE Socket Option

When the k eep-alive opt ion i s s et for a T CP s ocket a nd no data ha s been exchanged a cross t he s ocket i n e ither di rection f or t wo hour s, T CP automatically s ends a keep-alive pr obe to t he pe er. T his pr obe i s a T CP segment to which the peer must respond. One of three scenarios results:

1. The peer r esponds w ith the expected A CK. The a pplication i s not notified ( since e verything i s oka y). For f urther t wo hour of inactivity TCP will send a probe.

2. Peer host's c rash or r eboot i s r eported vi a R ST to the l ocal T CP. Sockets remaining errors are set to ECONNRESET and the socket is closed.

3. If p eer d oesn't r esponse t o k eep-alive p robe, B erkeley-derived TCPs s end 8 a dditional pr obes w ith g ap pe riod of 75 s econds. After 11 minutes and 15 seconds of inactivity, TCP will give up.

## SO_LINGER Socket Option

This opt ion s pecifies how t he close function ope rates f or a c onnection-oriented pr otocol (for T CP, but not f or U DP). By de fault, close returns immediately, b ut if th ere is a ny d ata s till r emaining in th e s ocket s end buffer, the system will try to deliver the data to the peer.

The SO_LINGER socket opt ion c an c hange t his de fault. T his opt ion requires t he f ollowing s tructure t o b e p assed ( as t he *optval argument) between t he us er pr ocess a nd t he ke rnel. It i s de fined b y including<sys/socket.h>.

```
struct linger {
int  l_onoff;      /* 0=off, nonzero=on */
 int  l_linger;     /* linger time, POSIX specifies units as seconds
*/ };
```

Calling setsockopt leads t o one of t he f ollowing t hree s cenarios, depending on the values of the two structure members:

1. If l_onoff is 0, t he opt ion i s t urned of f. T he v alue of l_linger is ignored a nd t he previously di scussed T CP de fault applies: close returns immediately.

2. If l_onoff is nonz ero and l_linger is z ero, TCP a borts t he connection when it is closed.

   ❖ In th is c ase, T CP d iscards a ny d ata s till r emaining i n t he socket s end bu ffer a nd sends a n R ST t o t he p eer, not t he normal four-packet connection termination sequence.

   ❖ This s cenario a voids T CP's T IME_WAIT s tate, but l eaves open the possibility of another incarnation of this connection being created within 2MSL seconds and having old duplicate segments f rom th e just-terminated c onnection be ing incorrectly delivered to the new incarnation.

   ❖ Occasional U SENET p ostings a dvocate t he us e of t his feature just to avoid the TIME_WAIT state and to be able to restart a lis tening s erver e ven if c onnections a re s till in u se with the server's well-known port. This should NOT be done and c ould l ead t o da ta corruption, a s de tailed i n RFC1337. Instead, t he SO_REUSEADDR socket opt ion s hould a lways be used in the server before the call to bind. We should make use o f th e T IME_WAIT s tate to le t o ld d uplicate s egments expire in the network rather than trying to avoid it.

   ❖ There ar e cer tain ci rcumstances w hich w arrant u sing t his feature t o s end an ab ortive cl ose. O ne ex ample i s an RS-232 terminal s erver, which mig ht h ang forever in CLOSE_WAIT trying to deliver data to a stuck terminal port, but w ould pr operly reset t he s tuck port i f i t got a n R ST t o discard the pending data.

3. If l_onoff is nonz ero and l_linger is nonz ero, t hen t he ke rnel w ill linger when the socket is closed.

   ❖ In th is s cenario, if th ere is a ny d ata s till r emaining in th e socket send buffer, the process is put to sleep until either:

      1. All t he d ata i s s ent an d ack nowledged b y t he peer TCP, or

      2. The linger time expires.

   ❖ If the socket has been set to nonblocking, it w ill not wait for the close to c omplete, e ven i f t he l inger t ime i s nonz ero. When us ing t his feature of t he SO_LINGER option, it is important f or t he a pplication t o c heck t he r eturn va lue from close, b ecause i f t he l inger t ime ex pires b efore t he remaining d ata i s s ent an d

acknowledged, close returns EWOULDBLOCK and a ny remaining data in the send buffer is discarded.

Given the above three scenarios, consider the situations when a close on a socket returns. Assume that the client writes data to the socket and then calls close.

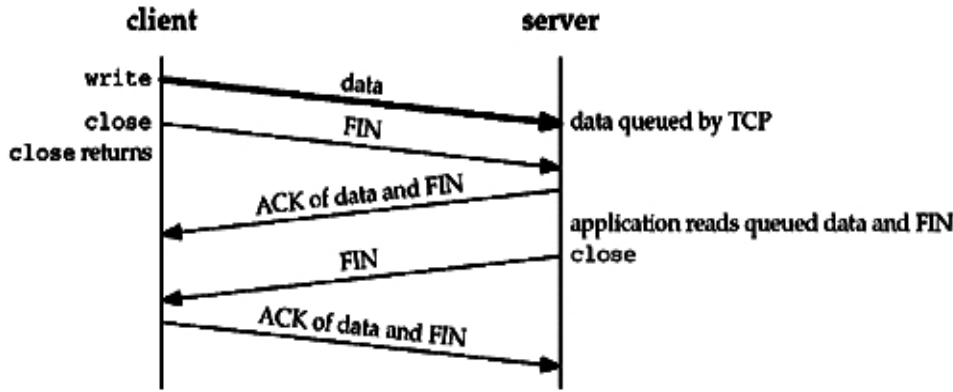## Default operation of close: it returns immediately *



Figure 6.1: Default operation of close

Assume that when the client's data arrives, the server is temporarily busy, so the data is added to the socket receive buffer by its TCP. Similarly, the next segment, the client's FIN, is also added to the socket receive buffer. But by default, the client's close returns immediately. In the scenario shown above, the client's close can return before the server reads the remaining data in its socket receive buffer. Therefore, it is possible for the server host to crash before the server application reads this remaining data, and the client application will never know.

## Close with SO_LINGER socket option set and l_linger a positive value *

The client can set the SO_LINGER socket option, specifying some positive linger time. When this occurs, the client's close does not return until all the client's data and its FIN have been acknowledged by the server TCP, as shown in the figure below.
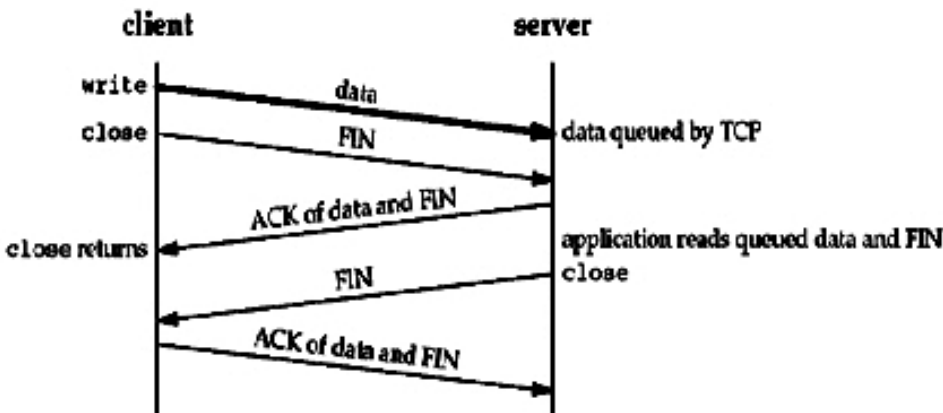


Figure 6.2: close with SO_LINGER option

BCA-E7/115

But this still has the same problem as: The server host can crash before the server application reads its remaining data, and the client application will never know. Worse, the following figure shows what can happen when the SO_LINGER option is set to a value that is too low.
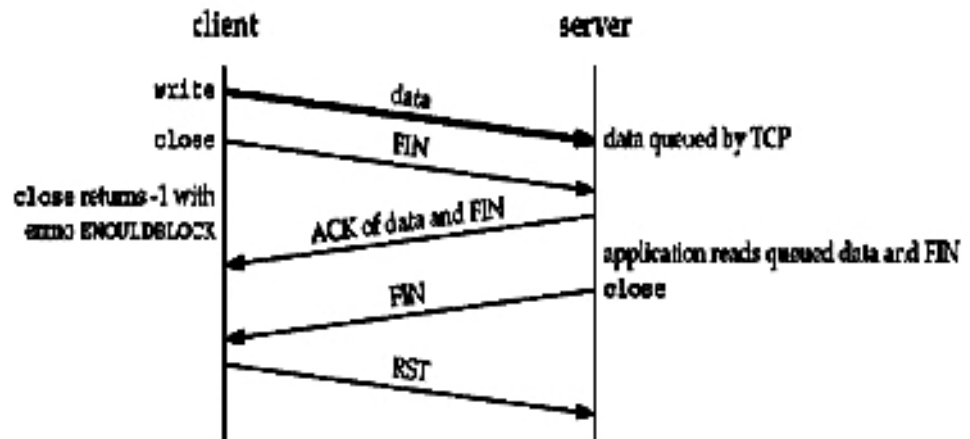


Figure 6.3: close with SO_LINGER option set to very low value

It is important to know that a successful return from close, with the SO_LINGER socket option set, only tells us that the data we sent (and our FIN) have been acknowledged by the peer TCP. This does not tell us whether the peer application has read the data. If we do not set the SO_LINGER socket option, we do not know whether the peer TCP has acknowledged the data.

## 6.4 IPV6 Socket Options

The following options are recognized at the IPv6 level:

| *protocolLevel* Options | Data Type | Description |
|---|---|---|
| IPV6_V6ONLY | int | Force the socket to be IPv6-only. Normally, when running with both TM_USE_IPV4 and TM_USE_IPV6 defined, a socket created with AF_INET6 is able to communicate via both IPv4 and IPv6. Setting this socket option forces the socket to communicate via IPv6 only. |
| IPV6_JOIN_GROUP | ipv6_mreq | Join an IPv6 multicast group. |
| IPV6_LEAVE_GROUP | ipv6_mreq | Leave an IPv6 multicast group. |

| | | |
|---|---|---|
| MCAST_JOIN_GROUP | group_req | Join a n IPv6 m ulticast group. T his opt ion a lso supports IPPROTO_IP. |
| MCAST_LEAVE_GROUP | group_req | Leave a n IPv6 mu lticast group. T his opt ion a lso supports IPPROTO_IP. |
| MCAST_BLOCK_SOURCE | group_sour ce_req | Block da ta f rom a g iven source t o a g iven IPv6 multicast g roup ( mute). This opt ion a lso s upports IPPROTO_IP. |
| MCAST_UNBLOCK_SOUR CE | group_sour ce_req | Unblock da ta from a given source t o a g iven IPv6 multicast g roup ( un-mute). This opt ion a lso s upports IPPROTO_IP. |
| MCAST_JOIN_SOURCE_GR OUP | group_sour ce_req | Join a s ource-specific I Pv6 group. T his opt ion a lso supports IPPROTO_IP. |
| MCAST_LEAVE_SOURCE_ GROUP | group_sour ce_req | Leave a s ource-specific IPv6 group. T his opt ion also supports IPPROTO_IP. |
| IPV6_MULTICAST_HOPS | unsigned int | This opt ion a llows the u ser to s et the h op limit f ield in the I Pv6 header f or multicast p ackets s ent v ia this socket. **Default 1** |
| IPV6_MULTICAST_IF | int | Specify the i nterface i ndex of the outgoing interface for multicast datagrams sent on this s ocket. A n i nterface index of 0 i ndicates that we want t o r eset a p reviously set out going i nterface for multicast p ackets s ent o n this socket. |
| IPV6_UNITCAST_HOPS | int | This opt ion a llows the u ser to s et the h op limit f ield in the IPv6 h eader f or u nicast packets sent via this socket. |

Table 6.3: IPv6 socket options

## 6.5    ICMP6 SOCKET OPTIONS

The ICMP6_FILTER s ocket opt ion c an be us ed b y a R AW application t o f ilter out I CMPv6 m essage t ypes that i t doe s not ne ed t o receive. ICMPv6 provides function comparable to ICMPv4 plus IGMPv4 and ARPv4 functionality. An application might be interested in receiving only a subset of the messages received for ICMPv6.

This option i s e nabled o r di sabled w ith a s etsockopt(). T he opt ion va lue provides a 256 -bit a rray of m essage t ypes t hat s hould be f iltered. T o disable the option, the setsockopt() should be issued with an option length of 0. This causes the TCP/IP protocol stack's default filter to be in effect.

A getsockopt() with this option returns the value set by a setsockopt(). If a setsockopt() has not been done, the TCP/IP protocol stack's default filter is returned.

## 6.6    TCP SOCKET OPTIONS

The following options are recognized at the TCP level:

| Protocol Level Options | Data Type | Description |
|---|---|---|
| TCP_KEEPALIVE | int | Sets t he i dle t ime in seconds for a T CP connection b efore it st arts sen ding k eep alive probes. N ote that k eep al ive probes will be sent only if the SO_KEEPALIVE socket option is enabled. **Default 7,200 seconds.** |
| TCP_MAXRT | int | Sets the amount of time in seconds before the connection is broken once TCP starts retransmitting, or probing a zero window when t he pe er doe s no t r espond. A TCP_MAXRT value of 0 means the system d efault, an d -1 means r etransmit forever. If a p ositive value is specified, it may be rounded up t o t he connection next r etransmission time. Note t hat unless t he TCP_MAXRT v alue i s -1 (transmit f orever), t he connection can also be b roken i f t he num ber of |

| | | maximum r etransmission TM_TCP_MAX_REXMIT ha s be en reached. S ee T M_TCP_MAX_REXMIT below. **Default 0.** Meaning: u se t he sy stem default of T M_TCP_MAX_REXMIT times ne twork c omputed r ound trip time for an established c onnection. For a non-established co nnection, since t here is n o computed r ound t rip time yet, t he connection can be broken when e ither 75 seconds or w hen TM_TCP_MAX_REXMIT t imes d efault network r ound t rip t ime ha ve e lapsed, whichever occurs first). |
|---|---|---|
| TCP_MAXSEG | int | Sets the maximum T CP seg ment si ze sent on the network. N ote t hat t he TCP_MAXSEG v alue i s t he m aximum amount of da ta ( including T CP opt ions, but n ot t he TCP h eader) that can be sent per segment to the peer. This means that the amount of user data sent per segment is the value given by the TCP_MAXSEG option m inus a ny e nabled T CP o ption (for e xample 1 2 b ytes f or a TCP time stamp option). T he T CP_MAXSEG value can be decreased or increased prior to a c onnection establishment, but i t is not r ecommended t o s et it t o a v alue higher than t he IP M TU m inus 40 bytes (for e xample 1460 by tes on E thernet), since t his w ould c ause f ragmentation of TCP segments. Note: set ting the TCP_MAXSEG o ption w ill i nhibit the automatic c omputation of t hat v alue by the system based o n the IP MTU ( which avoids f ragmentation), a nd w ill a lso inhibit P ath MTU D iscovery. A fter the connection h as started, this v alue cannot be c hanged. N ote a lso t hat t he TCP_MAXSEG v alue c annot b e s et below 64 bytes. Default value is IP MTU minus 40 bytes. **Default is IP MTU minus 40 bytes.** |
| TCP_NODELAY | int | Set this option value to a non-zero value, to d isable th e N agle a lgorithm th at buffers the s ent da ta i nside t he TCP. Useful t o a llow cl ient's TCP t o send small p ackets as soon a s p ossible ( like mouse clicks). **Default 0.** |

| TCP_NOPUSH | int | Set this option value to a non-zero value, to force TCP to delay sending any TCP data until a full sized segment is buffered in the TCP buffers. Useful for applications that send continuous big chunks of data like FTP, and know that more data is coming. (Normally the TCP code sends a non full-sized segment, only if it empties the TCP buffer). **Default 0.** |
|---|---|---|
| TCP_STDURG | int | Set this option value to a zero value if the peer is a Berkeley system since Berkeley systems set the urgent data pointer to point to last byte of urgent data+1. **Default 1.** |
| TM_TCP_2MSLTIME | int | Sets the maximum amount of time TCP will wait in the TIME WAIT state, once it has initiated a close of the connection. **Default 60 seconds.** |
| TM_TCP_DELAY_ACK | int | Sets the TCP delay ack time in milliseconds. **Default 200 milliseconds.** |
| TM_TCP_FINWT2TIME | int | Sets the maximum amount of time TCP will wait for the remote side to close, after it initiated a close. **Default 600 seconds.** |
| TM_TCP_KEEPALIVE_CNT | int | Sets the maximum numbers of keep alive probes without any response from the remote, before TCP gives up and aborts the connection. **Default 8.** |
| TM_TCP_KEEPALIVE_INTV | int | Sets the interval between Keep Alive probes in seconds. See TM_TCP_KEEPALIVE_CNT. This value cannot be changed after a connection is established, and cannot be bigger than 120 seconds. **Default 75 seconds.** |
| TM_TCP_MAX_REXMIT | int | Sets the maximum number of retransmissions without any response from the remote, before TCP gives up and aborts the connection. **Default 12.** |
| TM_TCP_PACKET | int | Set this option value to a non-zero value to make TCP behave like a message-oriented protocol (i.e. respect packet boundaries) at the application level in |

| | | |
|---|---|---|
| | | both s end a nd receive d irections of da ta transfer. N ote t hat for t he r eceive direction to r espect pa cket boun daries, the TCP peer which is sending must also implement s imilar functionality in i ts send direction. This is useful as a reliable alternative t o U DP. N ote that p reserving packet boun daries w ith TCP w ill no t work c orrectly i f y ou us e out -of-band data. TM_USE_TCP_PACKET m ust be defined in < trsystem.h> t o use t he TM_TCP_PACKET option. **Default 0.** |
| TM_TCP_PEND_A CCEPT_RECV_W ND | unsigned long | Specify the size (in bytes) of the listening socket's r eceive window. This s ize w ill override the default s ize o r the size specified by setsockopt() with th e SO_RCVBUF flag. Once accept() is called on the listening s ocket, t he window s ize w ill return to the s ize specified by S O_RCVBUF ( or t he default). N ote: This s ize m ay not be larger t han t he de fault w indow s ize t o avoid shrinking of the receive window. |
| TM_TCP_PROBE_ MAX | unsigned long | Sets the maximum window probe timeout interval i n m illiseconds. T he ne twork computed w indow pr obe t imeout i s bound by T M_TCP_PROBE_MIN a nd TM_TCP_PROBE_MAX. **Default 60,000 milliseconds.** |
| TM_TCP_PROBE_ MIN | unsigned long | Sets the minimum window probe timeout interval i n m illiseconds. T he ne twork computed w indow pr obe t imeout i s bound by T M_TCP_PROBE_MIN a nd TM_TCP_PROBE_MAX. **Default 500 milliseconds.** |
| TM_TCP_PURE_A CK_SEGS | int | Option on ly a vailable i f TM_USE_TCP_PURE_ACK is de ined. Sets the number of outstanding un-ACKed segments, before a pure ACK is sent ( even i f the r ecv w indow ha s no t changed.) Default value is zero, in which case the st ack w ill b ehave as if TM_USE_TCP_PURE_ACK had n ot been d efined, a nd w ill on ly A CK e very other s egment pr ovided t hat it i s combined w ith a w indow update, or w ill ACK w hen the delay A CK timer ex pires regardless of the window update. **Default 0.** |

| | | |
|---|---|---|
| TM_TCP_REXMIT_CONTROL | int | Dynamically modify the behavior of the TCP re transmission timer for t he specified so cket. V alid v alues are 1 (Pause), 2 ( Resume), a nd 3 ( Reset). TM_USE_TCP_REXMIT_CONTROL must be defined in <trsystem.h> to make this option available. |
| TM_TCP_RTO_DEF | unsigned long | Sets t he T CP d efault r etransmission timeout value in milliseconds, used when no ne twork r ound t rip t ime ha s be en computed yet. **Default 3,000 milliseconds.** |
| TM_TCP_RTO_MAX | unsigned long | Sets t he maximum retransmission timeout i n m illiseconds. T he ne twork computed r etransmission t imeout i s bound by T M_TCP_RTO_MIN a nd TM_RTO_MAX. **Default 64,000 milliseconds.** |
| TM_TCP_RTO_MIN | unsigned long | Sets the minimum retransmission timeout in m illiseconds. The n etwork c omputed retransmission timeout i s bound by TM_TCP_RTO_MIN a nd TM_TCP_RTO_MAX. **Default 100 milliseconds.** |
| TM_TCP_SEL_ACK | int | Set this option value to a non-zero value to en able s ending th e T CP s elective Acknowledgment option. **Note**: This option can only be changed prior to establishing a TCP connection. **Default 1.** |
| TM_TCP_SLOW_START | int | Set this option value to zero, to disable the TCP slow start algorithm. **Default 1.** |
| TM_TCP_SSL_CLIENT | int | Set t his option t o e nable SSL client negotiation on t his s ocket, opt ionLength must be s izeof(int), a ny non -zero v alue will enable SSL client. |
| TM_TCP_SSL_SERVER | int | Set t his o ption t o en able S SL ser ver negotiation on t his s ocket, opt ionLength must be s izeof(int), a ny non -zero v alue will enable S SL s erver. N ote that, if y ou set t his option for a l istening socket, al l accepted sockets inherit this option value, you don't have to set this option again on an accepted socket. |
| TM_TCP_SSLSESSION | int | Set the S SL s ession I d f or t his s ocket. The op tion l ength m ust be s izeof(int). |

| | | |
|---|---|---|
| | | Note t hat, i f you s et this option f or a listening so cket, a ll acc epted so ckets inherit t his opt ion v alue, y ou don' t h ave to se t this o ption again o n an ac cepted socket |
| TM_TCP_SSL_SE ND_MIN_SIZE | int | Set the SSL send minimum size. If user's send data is less than this value, user data will be qu eued. O ption length must be sizeof(int), a nd o ption va lue c an n ot greater t han 0 xffff. D on't set t his v alue too big.<br>**Default value is defined as macro TM_SSL_SEND_DATA_MIN_SIZE (0)** |
| TM_TCP_SSL_SE ND_MAX_SIZE | int | Set the S SL r ecord m aximum si ze. E ach record w ill a t m ost h ave t hat m uch u ser data e ncapsulated. U ser d ata b igger t han this s ize lim it w ill b e cut into tw o records, O ption l ength m ust be sizeof(int), a nd o ption va lue c an n ot greater than 0x4000 to enable reasonable encapsulate. Don't set t his value too small. (<100 value will be rejected)<br>**Default value is defined as macro TM_SSL_SEND_DATA_MAX_SIZE (8000).** |
| TM_TCP_TS | int | Set this option value to a non-zero value to enable sending the Time stamp option.<br>**Note**: T his option c an onl y be changed prior to establishing a TCP connection.<br>**Default 1.** |
| TM_TCP_WND_S CALE | int | Set this option value to a non-zero value to enable sending the TCP window scale option.<br>**Note**: T his option c an onl y be changed prior to establishing a TCP connection.<br>**Default 1.** |
| TM_TCP_STATE | int | Get the state of the TCP vector associated with the socket.<br>**Note**: Read only value. |
| TM_TCP_USER_P ARAM | ttUserGeneric Union | Use this option t o s et/get user d ata for a specific T CP so cket. T o enable t his feature, un comment t he TM_USE_USER_PARAM macro definition in your <trsystem.h>. |
| TM_TCP_CA_HY BLA | int | Set this option value to 1, to switch to the TCP H ybla C ongestion A voidance Algorithm. T he TCP H ybla a lgorithm |

| | | yields better p erformance for T CP connections w ith a l ong r ound t rip time (such a s on a hi gh-latency terrestrial o r satellite radio link). Set this option value to 0, t o switch ba ck t o t he T CP R eno Congestion Avoidance Algorithm. **Default 0.** |
|---|---|---|
| TM_TCP_PACING | int | Set this option value to 1, to turn on TCP Pacing. W ith TCP P acing t urned on, the stack will attempt to send TCP segments within the congestion window and p eer receive w indow ov er t he R ound T rip Time, instead of sending them all at once. For be tter pe rformance, t his op tion should be turned on, i f the TCP HYBLA algorithm is s witched on. S et this option value to 0, to turn off TCP Pacing. **Default 0.** |
| TM_TCP_CA_WE STOOD | int | Set this option value to 1, to switch to the TCP W estwood+ C ongestion A voidance Algorithm. The TCP W estwood+ algorithm yields better performance on TCP co nnections o ver wireless l ossy links. Set this option value to 0, to switch back t o t he T CP Reno Congestion Avoidance Algorithm. **Default 0.** |

Table 6.4: TCP socket options

---

**Check your progress**

1.    Explain IPV6_MULTICAST_IF option of IPV6 Socket.

2.    What is the use of ICMP6_FILTER socket option?

---

## 6.8 SUMMARY

This unit details the socket concept and its programming tools. The two m ain s ocket f unctions a re di scussed w ith br ief de scription of i ts option p arameter which ar e u sed t o s et s pecific r equirement. Listening socket r etrieve s ome parameters f rom connected TCP t o ge t t he connection s tatus know ledge. Idea of generic s ockets i s a lso i ncluded i n the uni t. S ocket opt ions de scription a t IPv6, ICMP6 a nd T CP l evel a re included with their syntactic representation.

## 6.9    TERMINAL QUESTIONS

1.    Describe the use of socket options with example.

2.    Write a code segment to retrieve output buffer size and set it to 1024 bytes.

3.    How default of close function is modified and why it is needed?

4.    Explain what and why of generic options.

# UNIT-7 : ELEMENTARY UDP SOCKETS

## Structure

## 7.0    INTRODUCTION

The **User Datagram Protocol** (**UDP**) is one of the core members of t he internet pr otocol s uite. UDP i s t he s impler of t he t wo s tandard TCP/IP t ransport pr otocols w here prior c ommunications a re not r equired to set u p transmission c hannels or da ta pa ths. It i s a p rocess-to-process protocol t hat a dds onl y port num ber f or a ddressing, c hecksum f or da ta integrity and length information of data from the upper layer. With UDP, computer a pplications c an s end m essages, t o ot her hos ts on a n internet protocol (IP) network. Although this is an "unreliable" protocol due to no handshaking but unlike TCP, it does not include mechanisms for retrying on transmission failures or data corruption and also it has restrictions on message l ength ( a l ittle unde r 65536 bytes). It i s m ostly ne eded f or applications t hat us e broadcasting or m ulticasting a nd m ay pl ay performance-intensive r oles s uch a s mu ltimedia. UDP is s uitable f or purposes where error checking and correction is either not necessary and it also avoids the overhead of such processing at the network interface level. Time-sensitive a pplications of ten us e U DP be cause dr opping pa ckets i s preferable to waiting for delayed packets, which may not be an option in a real-time system.

## 7.1    OBJECTIVE

To understand the usage, properties and implementation of UDP. After this unit you will come to know about:

- Tasks performed by echo server function and drawbacks of UDP

- What happens when a datagram is lost.

- Lack of support for flow control in UDP

- How to determine an outgoing interface with UDP

## 7.2   ECHO SERVER FUNCTION

UDP i s a  "connectionless" p rotocol  which  enables a p rogram t o use a  s ingle U DP s ocket t o c ommunicate w ith  more t han one  hos t a nd port. U DP por t num bers a re e ntirely i ndependent of  T CP por t nu mbers, though the IANA tries to register the same port number for both UDP and TCP when a  given service is offered through both protocols.  In fact, one of the most important practical differences between TCP and UDP is that there are no message boundaries in a TCP stream, whereas in UDP, every packet (datagram) is effectively a self-contained message. For applications where reliability is not a concern and where all messages are known to fit within the limited size of datagram, this can occasionally make UDP more convenient t o us e t han  TCP. U DP s erver s ocket i s c reated i n m uch t he same way as a TCP server socket. The communication between client and server through UDP protocol is implemented through UDP echo server.

An e cho s erver i s  an a pplication w hich i s us ed t o t est t he c onnection between client and server. This server sends back whatever text the client sent. H owever, c lient s erver i s a n e nvironment  where s erver p rocess t he request sent by client.

In the UDP Echo server, we create a socket and bind to an advertised port number. Then an infinite loop is started to process the client requests for connections. Figure 7.1 shows the working of UDP Echo server.



Figure 7.1: Working of UDP Echo server

The p rocess receives d ata f rom t he client u sing recvfrom() f unction and echoes t he s ame d ata u sing t he s endto()  function. It ha ndles m ultiple clients au tomatically  as U DP i s a  d atagram b ased p rotocol h ence n o exclusive connection is required to a client in this case.

**Drawbacks of UDP:**

TCP ha s e merged as the dom inant protocol used  for t he bul k of i nternet connectivity o  wing t  o s  ervices f  or b  reaking l  arge d  ata s  ets i  nto individual packets, c  hecking   for  and r  esending l  ost pa  ckets a  nd reassembling p  ackets i  nto t he co  rrect s  equence.  But t hese  additional

services co me at a co st in t erms o f ad ditional d ata o verhead, an d d elays called latency.

## 7.3 LOST DATAGRAMS

UDP s ends t he pa ckets ove r l ower bandwidth overhead a nd latency. B ut packets c an be l ost or r eceived out of or der be tween s ender and receiver. UDP i s a n i deal pr otocol f or network applications i n w hich perceived l atency i s cr itical s uch as gaming, voi ce a nd vi deo communications, w hich c an s uffer s ome da ta l oss w ithout a dversely affecting pe rceived qua lity. In s ome c ases, forward e rror correction techniques a re us ed t o i mprove a udio a nd vi deo qua lity in s pite of s ome loss.

UDP c an a lso be us ed i n a pplications th at r equire lo ssless d ata transmission when the application is configured to manage the process of retransmitting l ost packets and co rrectly a rranging r eceived p ackets. T his approach can help to improve the data transfer rate of large files compared with TCP.

UDP client/server is not reliable. If a client datagram is lost, the client will block forever in its call to recvfrom in the function dg_cli. It may wait for a s erver r eply t hat w ill n ever ar rive. S imilarly, i f t he cl ient d atagram arrives at the server but the server's reply is lost, the client will again block forever in its c all to recvfrom. Ho wever, just p lacing a time out o n the recvfrom cannot be the s olution. F or e xample, i f w e do t ime out, w e cannot t ell w hether our datagram n ever m ade i t t o t he s erver, or i f t he server's reply never made it back.

## 7.4 LACK OF FLOW CONTROL WITH UDP

We now e xamine the e ffect of U DP not ha ving a ny flow c ontrol. First, we modify our dg_cli function to send a fixed number of datagrams. It no l onger r eads f rom s tandard i nput. F igure 1 s hows t he ne w ve rsion. This function writes 2,000 1,400-byte UDP datagrams to the server.

We ne xt m odify t he s erver t o r eceive d atagrams a nd c ount t he num ber received. This server no longer echoes datagrams back to the client. Figure 2 shows the new dg_echo function. When we terminate the server with our terminal in terrupt k ey ( SIGINT), i t pr ints t he num ber of r eceived datagrams and terminates.

udpcliserv/dgcliloop1.c

1.#include "unp.h"

2 #define NDG 2000 /* datagrams to send */

3 #define DGLEN 1400 /* length of each datagram */

4 void

```
5 dg_cli(FILE *fp, intsockfd, const SA *pservaddr, socklen_tservlen)
6 {
7     int    i;
8     char   sendline[DGLEN];
9     for (i = 0; i < NDG; i++) {
10        Sendto(sockfd, sendline, DGLEN, 0, pservaddr, servlen);
11    }
12 }
```

udpcliserv/dgecholoop1.c

```
1 #include    "unp.h"
2 static void recvfrom_int(int);
3 static intcount;
4 void
5 dg_echo(intsockfd, SA *pcliaddr, socklen_tclilen)
6 {
7     socklen_tlen;
8     char   mesg[MAXLINE];
9     Signal(SIGINT, recvfrom_int);
10    for ( ; ; ) {
11        len = clilen;
12        Recvfrom(sockfd, mesg, MAXLINE, 0, pcliaddr, &len);
13        count++;
14    }
15 }
16 static void
17 recvfrom_int(intsigno)
18 {
19    printf("\nreceived %d datagrams\n", count);
20    exit(0);
21 }
```

We now run the server on the host freebsd, a slow SPARCStation. We run the client on t he R S/6000 s ystem aix, c onnected di rectly w ith 100M bps

BCA-E7/130

Ethernet. A dditionally, we r un netstat -s on t he s erver, bot h be fore and after, as the statistics that are output tell us how many datagrams were lost.

---

### Check your progress

1. What is an echo server?

2. How does UDP improve the data transfer rate of large files compared with TCP?

---

## 7.5 DETERMINING OUTGOING INTERFACE WITH UDP

A c onnected U DP s ocket c an a lso be us ed t o determine t he out going interface that will be used to a particular destination. This is because of a side ef fect o f t he connect function w hen applied t o a U DP s ocket. T he kernel chooses the local IP ad dress (assuming the process has not al ready called bind to e xplicitly a ssign th is). T his lo cal IP a ddress is c hosen b y searching the r outing t able for t he d estination IP address, an d then us ing the primary IP address for the resulting interface.

udpcliserv/udpcli09.c

```
1 #include   "unp.h"

2 int

3 main(intargc, char **argv)

4 {

5    intsockfd;

6    socklen_tlen;

7    structsockaddr_incliaddr, servaddr;

8    if (argc != 2)

9       err_quit("usage: udpcli<IPaddress>");

10    sockfd = Socket(AF_INET, SOCK_DGRAM, 0);

11    bzero(&servaddr, sizeof(servaddr));

12    servaddr.sin_family = AF_INET;

13    servaddr.sin_port = htons(SERV_PORT);

14    Inet_pton(AF_INET, argv[1], &servaddr.sin_addr);

15    connect(sockfd, (SA *) &servaddr, sizeof(servaddr));

16    len = sizeof(cliaddr);

17    Getsockname(sockfd, (SA *) &cliaddr, &len);
```

18    printf("local address %s\n", Sock_ntop((SA *) &cliaddr, len));

19    exit(0);

20 }

If w e r un t he p rogram on t he m ulti hom ed ho st freebsd, w e ha ve t he following output:

freebsd % udpcli09 206.168.112.96

local address 12.106.32.254:52329

freebsd % udpcli09 192.168.42.2

local address 192.168.42.1:52330

freebsd % udpcli09 127.0.0.1

local address 127.0.0.1:52331

The f irst time w e r un th e p rogram, th e c ommand-line a rgument is a n IP address t hat follows t he de fault route. T he ke rnel a ssigns t he l ocal IP address to the primary address of the interface to which the default route points. T he s econd t ime, t he ar gument i s t he IP ad dress o f a s ystem connected to a second Ethernet interface, so the kernel assigns the local IP address to the primary address of this second interface. Calling connect on a U DP s ocket doe s not send a nything t o t hat ho st; it is e ntirely a lo cal operation that saves the peer's IP address and port. We also see that calling connect on an unbound U DP socket also assigns an ephemeral port to the socket.

## 7.6   SOLVED EXAMPLES

Ques: W hat i s t he l argest l ength t hat w e can p ass t o *sendto* for a UDP/IPv4 socket, th at i s, w hat is th e la rgest a mount o f d ata th at c an fit into a UDP/IPv4 datagram?

Solution:

The l argest IPv4 da tagram i s 65,535 b ytes, l imited b y t he 16 -bit to tal length field. The IP header requires 20 bytes and the UDP header requires 8 b ytes, l eaving a m aximum of 65,507 b ytes f or us er da ta. W ith IPv6 without j umbogram s upport, t he s ize of t he I Pv6 he ader i s 40 b ytes, leaving a m aximum of 65,487 b ytes f or us er d ata. The n ew v ersion of dg_cli has b een u sed. If you forget to s et the s end buf fer s ize, B erkeley-derived kernels return an error of EMSGSIZE from sendto, s ince the size of t he s ocket s end buf fer i s normally l ess t han r equired for a m aximum-sized UDP datagram. But if we set the client's socket buffer sizes and run the c lient pr ogram, not hing i s r eturned b y the s erver. W e c an ve rify t hat the c lient's da tagram i s s ent t o the s erver b y r unning t cpdump, but i f w e put a printf in the server, its call to recvfrom does not return the datagram. The problem is that the server's UDP socket receive buffer is smaller than the da tagram w e a re s ending, s o t he da tagram i s di scarded a nd no t

delivered t o t he s ocket. O n a FreeBSD s ystem, w e can v erify t his b y running netstat -s and looking at the "dropped due to full socket buffers" counter before and after our big datagram is received. The final solution is to modify the server, setting its socket send and receive buffer sizes.

## 7.7   SUMMARY

UDP is an unreliable transport layer protocol. It serves processes where error checking and correction is not necessary and processes that are time s ensitive, th at is , r eal ti me s ystem. Being c onnectionless, it enables to use single socket to communicate with more than one host and port. Every packet is a self-contained message in UDP. An echo server is an application used to test the connection between client and server. If a client datagram is lost it is blocked forever. U DP does not support flow control, but i s c an be used f or de termining o utgoing i nterface. T his protocol i s i deal f or ne twork a pplications l ike gaming, voi ce a nd vi deo communication that can suffer some data loss without adversely affecting the quality.

## 7.8   TERMINAL QUESTIONS

1.   Explain with examples the drawbacks of UDP.

2.   "If a client datagram is lost, the client will block forever in its call to recvfrom in the function dg_cli". Explain.

3.   State some real life examples of where UDP is used.

4.   Write a program to implement Echo server function.

# UNIT-8 : NAME AND ADDRESS CONVERSION

## Structure

## 8.0 INTRODUCTION

There i s s ome  IP a ddress a ttached w ith c orresponding dom ain name s erver ( DNS). D NS l ookup, N SLOOKUP or   IP l ookup a re t he process to find the IP address by searching the DNS until a match found. The D omain N ame S ystem al so s pecifies t he t echnical f unctionality o f the database service that is at its core. A DNS name server is a server that stores the DNS records for a domain; a DNS name server responds with answers to queries against its database. In a nutshell, you tell it what the human readable address is for a site and it will give you the IP address. There are some special IP addresses such as 127.0.0.1 which is default IP address of every computer. No m atter w hich c omputer you us e, i t w ill always have an IP address of 127.0.0.1 and a name 'localhost'. In addition, a c omputer c an ha ve m ore t han one  IP a ddress. In or der t o c onnect t o other c omputers, i t w ill ha ve a n  IP  address t hat i s know n t o ot her computers.

## 8.1 OBJECTIVE

At the end of this unit, we will able to know the working of DNS.

* The i mportance  and w orking of   gethost b y n ame f unction i s explained.

* The different resolver options are discussed.

- Functions of IPV6 and its support is mentioned.

- Use of Uname f unction a nd ot her i mportant ne tworking information is explained.

## 8.2 DNS

DNS i s hi erarchical na ming c onvention w hich c ontains information about services or any other resource connected to the network. It defines the DNS protocol, a detailed specification of the data structures and d ata communication e xchanges. T he Internet m aintains t wo principal namespaces, t he d omain n ame h ierarchy[ and t he Internet Protocol address spaces.

Most imp ortantly, it tr anslates mo re r eadily u sed d omain n ames to th e numerical IP addresses needed for the purpose of locating and identifying that resource. It provides worldwide directory service created in 1983 by Paul Mockapetris. The Domain Name System delegates the responsibility of assigning domain names and mapping those names to Internet resources by designating authoritative name servers for each domain.

There is often confusion about a host name and a domain name. A domain name is the name that is purchased from a registrar. It will be something like hc idata.com or hc idata.co.uk. N ote t hat t here i s no " www" a t t he beginning of a domain name. A domain name can be subdivided into sub-domains - for example www.hcidata.com. Once you own a domain, there is no reasonable limit to the number of the sub-domains you can create. In fact m any s ub-domains can be al located t o t he s ame h ost m achine. A ny requests for a sub-domain (e.g. www.hcidata.com) are converted to an IP address by DNS and the IP address is used to route the request through the network until it reaches the host machine.

In the early years of the internet, each sub-domain would have a unique IP address so it was common for a host machine to have only one sub domain name. Network a dministrators m ay de legate a uthority over sub-domains of t heir a llocated na me s pace t o ot her na me s ervers. T his mechanism pr ovides di stributed a nd f ault t olerant s ervice a nd w as designed to avoid a single large central database.

Nowadays, t he co mmon p ractice i s t o h ave m any sub-domains w ith t he same IP address. It is also common for the domain name to get converted into the IP address of the host machine that runs the www sub domain. For example, a hos t m achine t hat c onverts hos t na mes t o IP a ddress us ing DNS m ay b e cal led d ns.hcidata.com an d a h ost m achine t hat i s a w eb server may be called www.hcidata.com.

### IP address to Country

IP a ddresses a re a llocated b y r egional or ganizations. Therefore, it is relatively easy to work out the country in which an IP is likely to reside. When an IP is allocated to a company they are expected to be used in the

country in which the organization resides. But, there is nothing to stop a company allocating an IP to a machine in another country. A company is allocated a range of IP addresses X.Y.Z.0 to X.Y.Z.255 for use in England. This company has a private network with a branch office in New York. So, it uses most of the IP address in England but uses some of them in the United States. So, we cannot guarantee that the country is 100% correct when converting an IP address, but we would expect it to be correct at least 90% of the time.

## 8.3 GETHOST BY NAME FUNCTION

The gethostbyname function retrieves host information corresponding to a host name. This function has been deprecated by the introduction of the getaddrinfo function. Developers creating Windows Sockets 2 applications are advised to use the getaddrinfo function instead of gethostbyname.

struct hostent* FAR gethostbyname( _In_ const char *name);

### Return value

If no error occurs, gethostbyname returns a pointer to the hostent structure described above. Otherwise, it returns a null pointer and a specific error number. The gethostbyname function does not check the size of the *name* parameter before passing the buffer which may result heap corruption.

## 8.4 RESOLVER OPTION

The OptionsResolver component helps you configure objects with option arrays. It supports default values, option constraints and lazy options. The *resolver* is a set of routines in the C library that provide access to the Internet Domain Name System (DNS). The OptionsResolver component helps you configure objects with option arrays. It supports default values, option constraints and lazy options. The resolver configuration file is designed to be human readable format which contains a list of keywords with values that provide various type of resolver option.

The different configuration options are:

### nameserver- Name server IP address

Internet address of a name server that the resolver should query

Resolver query IP address from the name server. If there are multiple servers, the resolver queries them in order. If no nameserver entries are present, the default is the name server on the local machine.

### domain -Local domain name

Short names are used relative to the local domain. If no domain entry is present, the domain is determined form the local hostname returned by gethostname(). The domain part is taken to be everything after the first '.'.

The root domain is assumed if the hostname does not contain a domain part.

## search -Search list for host-name lookup

The search list is normally determined from the local domain name. However, by default, it contains only the local domain name. This may be changed by listing the desired domain search path following the search keyword with space and tabs separating the names. This process may be slow and may generate network traffic if the servers for the listed domains are not local. Queries will time out if no server is available for one of the domains.The search list is currently limited to six domains with a total of 265 characters.

## shortlist

Sorted address are returned by gethostbyname() through this option. A shortlist is specified by IP-address-netmask pairs. The IP address and optional network pairs are separated by slashes.

## options

It allows certain internal resolver variables to be modified. The syntax is options option where option is as follows:

## debug

It sets RES_DEBUG.

## ndots: n

It sets a threshold for the number of dots which must appear in a name given to res_query before an initial absolute query will be made. The default value for n is 1. It implies that if there are any dots in a name, the name will be tried first as an absolute name before any search list elements are appended to it.

## timeout: n

It sets the amount of time the resolver will wait for a response from a remote name server. It is measured in seconds.

## attempts: n

It sets the number of times the resolver will send a query to its name server before giving up.

## rotate

It makes round robin selection of name servers by spreading the query load among all listed servers.

## no-check-names

It di sables the m odern BIND checking of i ncoming hos tnames a nd m ail names f or i nvalid c haracters s uch a s unde rscore, non -ASCII or c ontrol characters.

### int6

This ha s t he e ffect of t rying a A AAA que ry b efore a n A que ry i nside gethostname() function. It maps IPv4 responses in IPv6 "tunneled form".

### ip6-bytestring

It c auses r everse IPv6 lookups t o be made us ing t he bi t-label f ormat described in RFC 2673.

### ip6-dotint/no-ip6-dotint

When t his opt ion i s c lear, r everse IPv6 l ookups a re m ade i n t he i p6.int zone. W hen t his opt ion i s s et, r everse IPv6 l ookups a re m ade i n t he ip6.arpa z one. r everse IPv6 l ookups a re m ade i n t he *ip6.arpa* zone b y default. This option is set by default.

### ends0

It enables support for the DNS extension described in RFC 2671.

### single-request

Sometime D NS s erver c annot ha ndle t hese que ries pr operly a nd m ake a requests t ime out . T his option di sables t he behaviour and m akes glibc perform the IPv6 and IPv4 requests sequentially.

### single-request-reopen

The r esolver us es t he s ame s ocket f or A and AAAA requests. S ome hardware doe s m istake t o s end ba ck onl y one r eply. T he c lient s its a nd waits for second reply. By turning this option ON, it closes the socket and opens a new one before sending the second request.

---

### Check your progress

1. Explain how DNS can be used in recursive way?

2. What are the return values returned by gethostbyname function?

---

## 8.5   FUNCTION AND IPV6 SUPPORT

Internet Protocol V ersion 6 ( IPv6) is a network layer protocol that enables d ata communications o ver a p acket s witched n etwork. P acket switching i nvolves t he s ending and r eceiving o f data i n pa ckets be tween two node s in a n etwork. The w orking s tandard for the IPv6 protocol w as published b y t he Internet E ngineering T ask F orce ( IETF) i n 1998.J apan

and Korea were acknowledged as having the first public deployments of IPv6

IPv6 and IPv4 share a similar architecture. The majority of transport layer protocols that function with IPv4 will also function with the IPv6 protocol. Most application layer protocols are expected to be interoperable with IPv6 as well. A main advantage of IPv6 is increased address space. The 128-bit length of IPv6 addresses is a significant gain over the 32-bit length of IPv4 addresses, allowing for an almost limitless number of unique IP addresses. The size of the IPv6 address space makes it less vulnerable to malicious activities such as IP scanning. IPv6 packets can support a larger payload than IPv4 packets resulting in increased throughput and transport efficiency. Notable exception of File Transfer Protocol (FTP).

### IPv6 functions

IBM is implementing IPv6 on i5/OS® over several software releases. IPv6 functions are transparent to existing TCP/IP applications and coexist with IPv4 functions.

These are the main i5/OS features that are affected by IPv6:

If you configure IPv6, you are sending IPv6 packets over an IPv6 network. Creating an IPv6 local area network for a scenario that describes a situation in which you configure IPv6 on your network.

The Start and Stop menu items on the TCP/IP Configuration folder are removed. IPv6 can be started and stopped in the same way as IPv4, with STRTCP and ENDTCP commands. IPv6 cannot be started or stopped independent of IPv4.

The Configure IPv6 wizard is removed from iSeries Navigator. The line configuration options in the wizard are replaced by actions on individual lines in the **Lines** folder. Similarly, you can use a new wizard to create IPv6 interfaces.

## 8.6   UNAME FUNCTION

This function is used to get name and information about current kernel. This is a system call, and the operating system presumably knows its name, release and version. It also knows what hardware it runs on. So, four of the fields of the struct are meaningful. On the other hand, the field *nodename* is meaningless: it gives the name of the present machine in some undefined network, but typically machines are in more than one network and have several names. Moreover, the kernel has no way of knowing about such things, so it has to be told what to answer here. The same holds for the additional *domainname* field.

It returns system information in the structure pointed to by *buf*.

#include <sys/utsname.h>

int uname(struct utsname *buf*);

The *utsname* struct is defined in *<sys/utsname.h>*:

```
struct utsname {
char sysname[];    /* Operating system name*/
char nodename[];   /* Name within "some implementation-defined
                            network" */
 char release[];   /* Operating system release*/
 char version[];    /* Operating system version */
 char machine[];    /* Hardware identifier */
#ifdef _GNU_SOURCE
 char domainname[]; /* NIS or YP domain name */
 #endif
 };
```

On success, zero is returned. On error, -1 is returned, and *errno* is set appropriately.

The length of the fields in the struct varies. Some operating systems or libraries use a hardcoded 9 or 33 or 65 or 257. Other systems use SYS_NMLN or _SYS_NMLN or UTSLEN or _UTSNAME_LENGTH. Clearly, it is a bad idea to use any of these constants; just use sizeof(...). Often 257 is chosen in order to have room for an internet hostname.

## 8.7    Other Networking Information

When looking at networking basics, understanding the way a network operates is the first step to understanding routing and switching. The network operates by connecting computers and peripherals using two pieces of equipment; switches and routers. Switches and routers, essential networking basics, enable the devices that are connected to your network to communicate

### Networking Basics: Switches

- Switches are used to connect multiple devices on the same network within a building or campus. For example, a switch can connect your computers, printers and servers, creating a network of shared resources. The switch, one aspect of your networking basics, would serve as a controller, allowing the various devices to share information and talk to each other. Through information sharing and resource allocation, switches save you money and increase productivity.

There are two basic types of switches to choose from as part of your networking basics: managed and unmanaged.

❖ An unmanaged switch works out of the box and does not allow you to make changes. Home-networking equipment typically offers unmanaged switches.

❖ A managed switch allows you access to program it. This provides greater flexibility to your networking basics because the switch can be monitored and adjusted locally or remotely to give you control over network traffic, and who has access to your network.

• **Routers**, the second valuable component of your networking basics, are used to tie multiple networks together. For example, you would use a router to connect your networked computers to the Internet and thereby share an Internet connection among many users. The router will act as a dispatcher, choosing the best route for your information to travel so that you receive it quickly. Routers analyse the data being sent over a network, change how it is packaged, and send it to another network, or over a different type of network. They connect your business to the outside world, protect your information from security threats, and can even decide which computers get priority over others. Depending on your business and your networking plans, you can choose from routers that include different capabilities. These can include networking basics such as:

❖ **Firewall:** Specialized software that examines incoming data and protects your business network against attacks.

❖ **Virtual Private Network (VPN):** A way to allow remote employees to safely access your network remotely.

❖ **IP Phone network:** Combine your company's computer and telephone network, using voice and conferencing technology, to simplify and unify your communications.

---

**Check your progress**

1. How does IPV6 improve throughput and transport efficiency of a network?

2. What is the use of uname function?

## 8.8   SOLVED EXAMPLE

Q. Modify following pr ogram t o c all getnameinfo i nstead of s ock_ntop. What flags should you pass to getnameinfo?

names/daytimetcpcli1.c

1 #include "unp.h"

2 int

3 main (int argc, char **argv)

4 {

5 int sockfd, n;

6 char recvline [MAXLINE + 1];

7 struct sockaddr_in servaddr;

8 struct in_addr **pptr;

9 struct in_addr *inetaddrp [2];

10 struct in_addr inetaddr;

11 struct hostent *hp;

12 struct servent *sp;

13 if (argc ! = 3)

14 err_quit ("usage: daytimetcpclil ");

15 if ( (hp = gethostbyname (argv [1]) ) == NULL) {

16 if (inet_aton (argv [1], &inetaddr) == 0) {

17 err_quit ("hostname error for %s: %s", argv [1],

18 hstrerror (h_errno) );

19 } else {

20 inetaddrp [0] = &inetaddr;

21 inetaddrp [1] = NULL;

22 pptr = inetaddrp;

23 }

24 } else {

25 pptr = (struct in_addr **) hp->h_addr_list;

26 }

27 if ( (sp = getservbyname (argv [2], "tcp") ) == NULL)

28 err_quit ("getservbyname error for %s", argv [2] );

29 for ( ; *pptr != NULL; pptr++) {

30 sockfd = Socket (AF_INET, SOCK_STREAM, 0) ;

31 bzero (&servaddr, sizeof (servaddr) ) ;

32 servaddr.sin_family = AF_INET;

33 servaddr.sin_port = sp->s_port;

34 memcpy (&servaddr.sin_addr, *pptr, sizeof (struct in_addr) ) ;

35 printf ("trying %s\n", Sock_ntop ( (SA *) &servaddr, sizeof (servaddr) ) ) ;

36 if (connect (sockfd, (SA *) &servaddr, sizeof (servaddr) ) == 0)

37 break; /* success */

38 err_ret ("connect error");

39 close (sockfd) ;

40 }

41 if (*pptr == NULL)

42 err_quit ("unable to connect");

43 while ( (n = Read (sockfd, recvline, MAXLINE) ) > 0) {

44 recvline [n] = 0; /* null terminate */

45 Fputs (recvline, stdout);

46 }

47 exit (0);

48 }

Solution:

Following modifications are made in the code given above.

1. We first allocate arrays to hold the hostname and service name as follows:

   char host[NI_MAXHOST], serv[NI_MAXSERV];

2. After accept returns, we call getnameinfo instead of sock_ntop as follows:

   if (getnameinfo(cliaddr, len, host, NI_MAXHOST, serv, NI_MAXSERV, NI_NUMERICHOST | NI_NUMERICSERV) == 0)

   printf("connection from %s.%s\n", host, serv);

**3.** Since this is a server, we specify the NI_NUMERICHOST and NI_NUMERICSERV flags to avoid a DNS query and a lookup of /etc/services.

## 8.9   SUMMARY

The Domain Name System specifies the technical functionality of the database service that is at its core. A DNS name server is a server that stores the DNS records for a domain; a DNS name server responds with answers to queries against its database. The Internet maintains two principal namespaces, the domain name hierarchy[ and the Internet Protocol address spaces. The *resolver* is a set of routines in the C library that provide access to the Internet Domain Name System (DNS). Internet Protocol Version 6 ( IPv6) is a network layer protocol that enables data communications over a packet switched network. The 128-bit length of IPv6 addresses is a significant gain over the 32-bit length of IPv4 addresses, allowing for an almost limitless number of unique IP addresses. The Uname function is used to get name and information about current kernel. It also knows what hardware it runs on. Switches and routers enable the devices that are connected to your network to communicate. Switches are used to connect multiple devices on the same network providing information sharing and resource allocation that in turn saves your money and increases productivity. Routers are used to tie multiple networks together, choosing the best route for your information, connect your business to the outside world, protect your information from security threats, and can even decide which computers get priority over others.

## 8.10   TERMINAL QUESTIONS

**1.** State the similarities and differences between IPv4 and IPv6.

**2.** Write a program for choosing the best route for some hypothetical network.

**3.** How load balancing is achieved using DNS?

**Uttar Pradesh Rajarshi Tandon
Open University**

# Bachelor of Computer Application

## BCA-E7
### Network Programming

**Block**

# 3

## DAEMON PROCESSES, ADVANCE I/O FUNCTIONS AND UNIX DOMAIN PROTOCOLS

# Course Design Committee

| | |
|---|---|
| **Dr. Ashutosh Gupta** | Chairman |
| Director (In-charge) | |
| School of Compruter and Information Science, UPRTOU Prayagraj | |
| **Prof. R. S. Yadav** | Member |
| Department of Computer Science and Engineering | |
| MNNIT-Allahabad, Prayagraj | |
| **Ms Marisha** | Member |
| Assistant Professor (Computer Science), | |
| School of Science UPRTOU Prayagraj | |
| **Mr. Manoj Kumar Balwant** | Member |
| Assistant Professor, (Computer Science) | |
| School of Sciences UPRTOU Prayagraj | |

# Course Preparation Committee

| | |
|---|---|
| **Dr. Prabhat Kumar** | Author (Block 1,2) |
| Assistant Professor, Department of IT | |
| NIT Patna | |
| **Dr. Prabhat Ranjan** | Author (Block 3,4) |
| Assistant Professor, Department of Computer Science | |
| Central University of South Bihar | |
| **Dr. Rajiv Mishra** | Editor |
| Associate Professor, Department of CSE | |
| IIT Patna | |
| **Dr. Ashutosh Gupta** (Director in Charge) | |
| School of Computer & Information Sciences, | |
| UPRTOU Prayagraj | |
| **Mr. Manoj Kumar Balwant** | Coordinator |
| Assistant Professor, (Computer Science) | |
| School of sciences UPRTOU Prayagraj | |

# BLOCK INTRODUCTION

The objective of this course is to introduce the basic concept about the network programming as well as provides a mix of practical experience and a depth of understanding. The network programming course address today's most crucial standards, implementations and techniques. The aim is to provide an extensive variety of topics on this subject with appropriate examples. The course is organized into following blocks:

Block 3 describes the daemon processes, advance I/O functions and UNIX domain protocols.

# UNIT-9 : DAEMON PROCESSES

**Structure**

## 9.1    INTRODUCTION

In t his uni t, we w ill le arn about daemons an d characteristics o f daemon processes. Further on we will look on how to log messages using *syslog* facility. Then da emon pr oviding i nternet s ervices i s di scussed. Then we will also have a look on *daemon_init* function, *inetd* daemon and *daemon_inetd* functions.

## 9.2    OBJECTIVES

At the end of this unit we will have knowledge about: -

- Daemons and their characteristics.
- Ways to start a daemon.
- *syslogd* Daemon and *syslog* function.
- *daemon_init* function, *inetd* daemon and *daemon_inetd* functions.

## 9.3    DAEMON

A daemon is a process that runs in the background as a background process instead of being under the direct control of an interactive user. In other w ords, it is not associated w ith controlling te rminal or lo gin s hell. Unix systems typically have many processes that are daemons, running in the background, performing different administrative tasks.

Generally, in UNIX system, the name of the daemon process end with the letter *d*. As for example, *syslogd* daemon, *inetd* daemon, *sshd* daemon, etc.

*System* daemons have the following characteristics: -

- Started once when the system is initialized.

- Runs until the system is shut down.

- During the service time, spends most of their time waiting for some event to occur.

- Frequently spawn other processes to handle service requests.

**Ways to start a daemon: -**

a. Many daemons are started by the system in itialization scripts. These scripts are mainly in the directory / etc. or in a directory whose name begins with /etc/rc.

b. Many network servers are started by *inetd* superserver.

c. The execution of programs on a regular basis is performed by the *cron* daemon, and programs that it invokes run as daemons.

d. The execution of a program at one time in the future is specified by the at command. The cron daemon normally initiates these programs when their time arrives, so these programs run as daemons.

e. Daemons can be started from user terminals, either in the foreground or in the background.

## 9.4  *syslogd* DAEMON

Many versions of UNIX provide a general-purpose logging facility called *syslog*. Individual programs that need to have information logged send the information to *syslog*. In order to handle these logs status *syslogd* daemon comes into play. Purpose of *syslogd* daemon is to log system messages. It reads the log message and does what the configuration file

(normally */etc. / syslog.conf*) s pecifies t o do w ith that message. If t he daemon receives the SIGHUP signal, it rereads its configuration file.

Berkeley-derived imp lementations o f s yslogd perform t he f ollowing actions on startup:

1. The c onfiguration f ile, nor mally / etc/syslog.conf, i s r ead, specifying w hat t o do with e ach t ype o f l og message t hat t he daemon can r eceive. T hese m essages c an b e ap pended t o a f ile written to a s pecific u ser, or forwarded to the s yslogd da emon on another host.

2. A U nix dom ain s ocket i s c reated a nd bound to t he pa thname /var/run/log.

3. A UDP s ocket i s created a nd bound t o por t 514 (the s yslog service).

4. The pa thname / dev/klog i s ope ned. A ny error m essages f rom within the kernel appear as input on this device.

## 9.5 *syslog* FUNCTION

Since d aemon d oesn't h ave a controlling te rminal, it n eeds s ome way t o out put m essages when s omething ha ppens l ike nor mal informational messages or emergency messages that need to be handled by an administrator. So, there comes the role of *syslog* function.

Structure of syslog function:

#include <syslog.h>

void syslog (int priority, const char *message, …);

Here, priority is combination of *level* (0 to 7) and *facility* (to identify the type of process sending the message). Log messages have a level between 0 and 7, which shown in table 9.1. These are ordered values. If no level is specified by the sender, LOG_NOTICE is the default**.**

| Level | Value | Description |
|---|---|---|
| LOG_EMERG | 0 | System is unusable (highest priority) |
| LOG_ALERT | 1 | Action must be taken immediately |
| LOG_CRIT | 2 | Critical conditions |
| LOG_ERR | 3 | Error conditions |
| LOG_WARNING | 4 | Warning conditions |
| LOG_NOTICE | 5 | Normal but significant condition(default) |
| LOG_INFO | 6 | Informational |
| LOG_DEBUG | 7 | Debug-level messages (lowest priority) |

Table 9.1: Level of log messages

Log m essages al so contain a f acility t o i dentify the t ype of pr ocesses sending t he m essages. We s how t he di fferent values i n table 9 .2. I f no facility is specified, LOG_USER is the default.

| Facility | Description |
|---|---|
| LOG_AUTH | Security/authorization messages |
| LOG_AUTHPRIV | Security/authorization messages(private) |
| LOG_CRON | Cron daemon |
| LOG_DAEMON | System daemons |
| LOG_FTP | FTP daemon |
| LOG_KERN | Kernel messages |
| LOG_LOCAL0 | Local use |
| LOG_LOCAL1 | Local use |
| LOG_LOCAL2 | Local use |
| LOG_LOCAL3 | Local use |
| LOG_LOCAL4 | Local use |
| LOG_LOCAL5 | Local use |
| LOG_LOCAL6 | Local use |
| LOG_LOCAL7 | Local use |
| LOG_LPR | Line printer system |
| LOG_MAIL | Mail system |
| LOG_NEWS | Network news system |
| LOG_STSLOG | Messages g enerated internally b y syslogd |
| LOG_USER | Random user-level messages(default) |
| LOG_UUCP | UUCP system |

Table 9.2: Facility of log messages

The p urpose o f f acility and l evel i s t o al low al l m essages f rom a given facility to be handled the same in the /etc/syslog.conf file or to allow all messages of a given level to be handled the same.

<table>
<tr><td>

***Check Your Progress***

1.    *Can you define the steps performed by syslogd Daemon?*

2.    *Can you define the different levels of log message?*

</td></tr>
</table>

# 9.6   *DAEMON_INIT* FUNCTION

daemon_init function is used to demonize a process i.e. to start an arbitrary program and run it as a daemon. This function should be suitable for use on a ll va riants o f UNIX b ut some o ffer C l ibrary function cal led daemon that provides similar feature.

The program below shows a function named daemon_init that can call to daemonize the process.

```
#include "unp.h"
#include    <syslog.h>
#define  MAXFD  64
extern int daemon_proc;
int
daemon_init(const char *pname, int facility)
{
    int  i;
    pid_t  pid;
    if ((pid=Fork ()) <0)
      return (-1);
    else if (pid)
        _exit (0); /*
    /* child 1 continues…*/
  if(setsid()   < 0)
      return (-1);
  signal (SIGHUP,SIG_IGN);
  if ((pid =Fork ()) < 0)
      return (-1);
  else if (pid)
```

```
    _exit (0);

/* child 2 continues…*/

daemon_proc = 1;

chdir("/")

/*close off file descriptors */

for (i=0; i<MAXFD; i++)

  close(i);

/* redirect stdin, stdout and stderr  to /dev/null */

open ("/dev/null", O_RDONLY);

open ("/dev/null", O_RDWR);

open ("/dev/null", O_RDWR);

openlog (pname,  LOG_PID, facility);

return (0);

}
```

In the program, the daemon_init function first call fork and then the parent terminates, and child continues. If the process starts as a shell command in the foreground, when the parent terminates, the shell thinks the command is done. This automatically runs the child process in the background. Also, the child inherits the process group ID from the parent but gets its own process ID. This guarantees that the child is not a process group leader, which is required for the next call to setsid.

The setsid is a POSIX function that creates a new session. The process becomes the session leader of the new session. The process becomes group leader of a new process group and has no controlling terminal. Ignore SIGHUP and call fork again. When this function returns, the parent is the first child and it terminates, leaving the second child running. The purpose of this second fork is to guarantee that the daemon cannot automatically acquire a controlling terminal should it open a terminal device in the future. The calling fork in a second time, guarantee that the second child is no longer a session leader, so it cannot acquire a controlling terminal.

Then set flag for error functions. Set the global daemon_proc to nonzero. Then change the working directory to the root directory, although some daemons might have a reason to change to some other directory. After that close, any open descriptors that are inherited from the process that executed the daemon.  After that redirect stdin, stdout, and stderr to /dev/null for standard input, standard output, and standard error. Then open log is called. The first argument is from the caller and is normally the name of the program (e.g., argv[0]). The process ID should be added to each log message. This facility is also specified by the caller.

## 9.7  *INETD* DAEMON

inetd r efers t o i nternet s ervice da emon. i netd  daemon i s a superserver ( service d ispatcher) d aemon o n m any  UNIX  systems th at provide internet s ervices. This d aemon i s u sed b y servers t hat u se ei ther TCP or UDP.

This *inetd* process establishes itself as a daemon using the techniques that we described with our daemon_init function. It then reads and processes its c onfiguration f ile,  typically / etc/inetc.conf.  This f ile s pecifies t he services that the super server is to handle, and what to do when a service request arrives. The table 9.3 shows the fields in inetd.conf file.

| Field | Description |
|---|---|
| service-name | Must be in /etc/services |
| socket-type | Stream(TCP) or dgram (UDP) |
| Protocol | Must be in /etc/protocols either tcp or udp |
| wait-flag | Typically, nowait for TCP or wait for UDP |
| login-name | From /etc/passwd: typically root |
| server-program | Full pathname to exec |
| server-program-arguments | Arguments for exec |

Table 9.3: Fields in inetd.conf file

When a  TCP packet o r UDP packet ar rives w ith a p articular d estination port number, *inetd* launches the appropriate server program to handle the connection.

The steps performed by inetd shown in figure 9.1. The steps are as follows: -

- On s tartup, i t r eads t he */etc/inetd.conf* and  creates a s ocket o f appropriate type for all the services specified in the file.

- *bind* is called for the socket to specify port and IP address for the server.

- *listen* is called for TCP (not needed for datagram sockets).

- After cr eation o f s ocket *select* is c alled to  w ait f or a ny o f th e socket to become readable.

- When the *select* returns that a socket is readable, *accept* is called to accept the new connection (only for TCP connection).

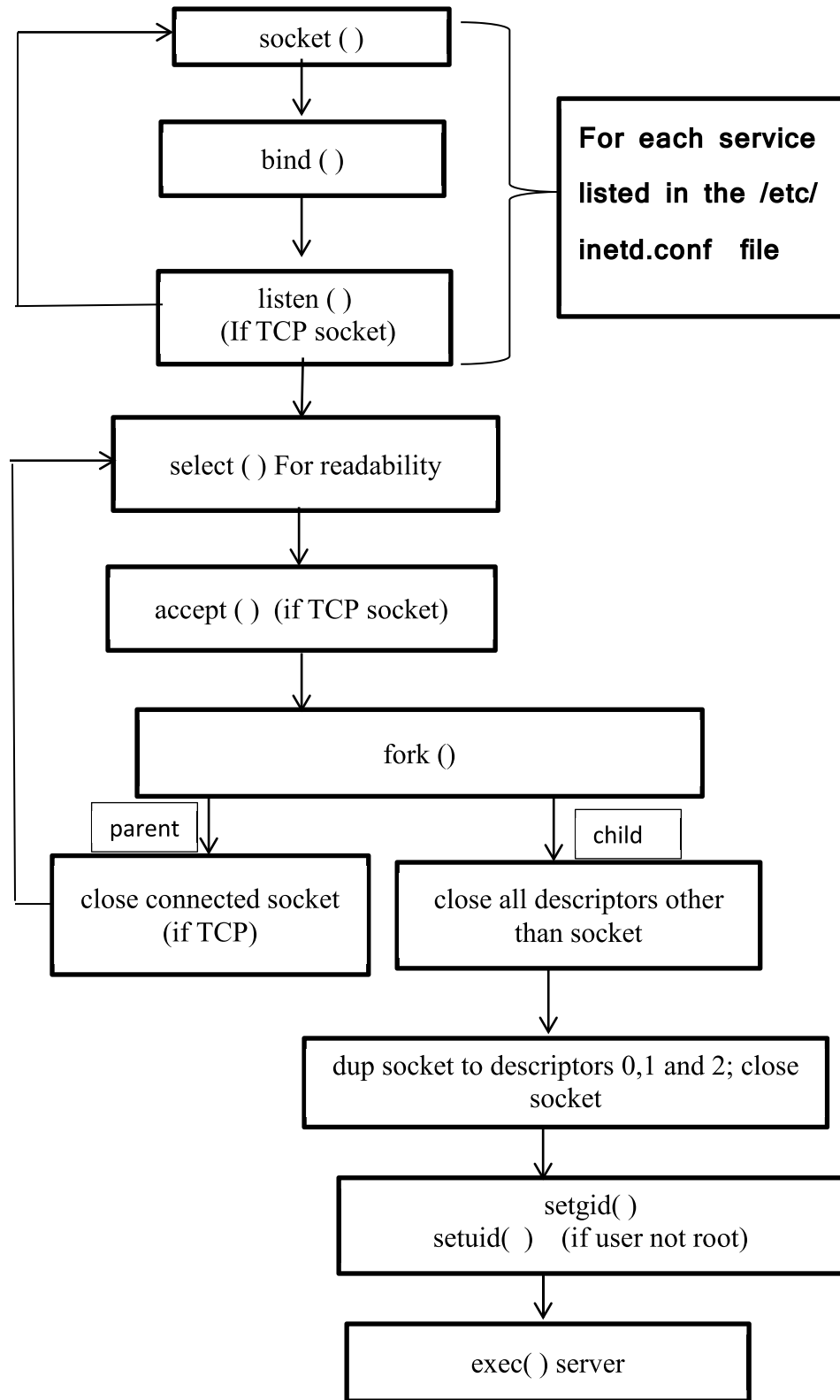- The *inetd* daemon *forks* and the child process handle the service request.

```
        ┌────────────────────┐
    ┌──→│      socket ( )     │───┐
    │   └────────────────────┘   │     ┌──────────────────────┐
    │            │               │     │ For each service     │
    │            ↓               │     │ listed in the /etc/   │
    │   ┌────────────────────┐   │     │ inetd.conf  file     │
    │   │       bind ( )      │   ├────→│                      │
    │   └────────────────────┘   │     └──────────────────────┘
    │            │               │
    │            ↓               │
    │   ┌────────────────────┐   │
    └───│     listen ( )      │   │
        │    (If TCP socket)  │───┘
        └────────────────────┘
                 │
                 ↓
        ┌────────────────────┐
    ┌──→│ select ( ) For readability │
    │   └────────────────────┘
    │            │
    │            ↓
    │   ┌────────────────────┐
    │   │ accept ( )  (if TCP socket) │
    │   └────────────────────┘
    │            │
    │            ↓
    │   ┌──────────────────────────────────┐
    │   │             fork ()              │
    │   └──────────────────────────────────┘
    │     parent │              │ child
    │            ↓              ↓
    │   ┌──────────────┐  ┌──────────────────┐
    └───│ close connected │ │ close all descriptors other │
        │ socket (if TCP) │ │     than socket   │
        └──────────────┘  └──────────────────┘
                                    │
                                    ↓
                          ┌──────────────────────────┐
                          │ dup socket to descriptors 0,1 and 2; close socket │
                          └──────────────────────────┘
                                    │
                                    ↓
                          ┌──────────────────────────┐
                          │ setgid( )                │
                          │ setuid( )  (if user not root) │
                          └──────────────────────────┘
                                    │
                                    ↓
                          ┌──────────────────────────┐
                          │     exec( ) server        │
                          └──────────────────────────┘
```

Figure 9.1: Steps performed by inetd

## 9.8  *DAEMON_INETD* FUNCTION

It demonizes process run b y *inetd*. This function is tr ivial compared t han da emon_init be cause a ll of t he daemonization s teps a re performed by inetd when it starts.

**The program below shows a function named daemon_inetd**

```
#include "unp.h"

#include  <syslog.h>

extern int daemon_proc;

void

daemon_inetd(const char *pname,int facility)

{

   daemon_proc =1;

   openlog (pname, LOG_PID,facility)

}
```

The daemonization a re performed b y i netd w hen i t s tarts. T he daemon_proc f lag f or e rror f unctions a nd o pen log is c alled. T he f irst argument is from the caller and is normally the name of the program. The process ID s hould b e a dded t o ea ch l og m essage. T his facility is a lso specified by the caller.

## 9.9  SUMMARY

Daemons a re t he pr ocesses r unning i n ba ckground and a re independent of control from all terminals. All outputs from a daemon are normally s ent t o *syslog* daemon b y c alling t he *syslog* function. Start o f daemon r equires a few steps. *daemon_init* handles d etails o f s tarting a daemon. Many UNIX servers that provide internet services are started b y the *inetd* daemon.

## 9.10  TERMINAL QUESTIONS

1. Define daemon and its characteristics.

2. Write ways to start a daemon.

3. Define *syslogd* daemon and *syslog* function.

4. Explain function of *inetd* daemon.

5. Explain steps of *inetd* daemon.

6. Explain *daemon_init* function with an example.

# UNIT-10 : ADVANCE I/O FUNCTIONS

**Structure**

## 10.1  INTRODUCTION

In t his unit, we w ill c over a va riety of f unctions a nd t echniques that i s cat egorized as " Advance I/O". First, we w ill s ee t hree w ays o f setting a time out o n I/O o peration involved i n s ocket. N ext co mes t hree variations on t he *read* and *write* functions. We will s tudy about ancillary data. Then we will also have a look on how to determine the amount of data in the socket receive buffer and how to use the C standard I/O library with sockets.

## 10.2  OBJECTIVE

At the end of this unit we will get to know: -

- Three ways to place a timeout on I/O operation involving a socket.

- Format of variations of *read* and *write.*

- About ancillary data

- How to determine the amount of data in the socket receive buffer?

- How to use the C standard I/O library with sockets?

## 10.3  SOCKET TIMEOUTS

Sockets involve some I/O operation (like *read, write, etc.*). So, a timeout can be placed on these I/O operations.

Following are three ways to place the timeout: -

- Call *alarm*, which generates the SIGALRM signal when the specified time has expired.

- Block waiting for I/O in *select*, which has a time limit built in instead of blocking in a call to read or write.

- Use the newer *SO_RCVTIMEO* and *SO_SNDTIMEO* socket options.

## 10.4  *RECV* AND *SEND* FUNCTIONS

*recv* function is used to receive messages from a socket and may be used to receive data on a socket. It is normally used only on a connected socket.

*send* function is used to transmit a message to another socket. It is normally used only on a connected socket.

These two functions are similar to the standard *read* and *write* functions, but one additional argument is required. Header file for these operations is <sys/socket.h>. Format of *read* and *write* functions are given below:

#include <sys/socket.h>

ssize_t recv(int *sockfd*, void *\*buff*, size_t *nbytes*, int *flags*);

ssize_t send (int *sockfd*, void const *\*buff*, size_t *nbytes*, int *flags*);

Both return: number of bytes read or write if OK, -1 on error

Here, first three arguments are same as the first three arguments to *read* and *write*. *recv* read *nbytes* bytes from socket file descriptor *sockfd* into buffer *buff*. *send* function sends *nbytes* to the socket file descriptor *sockfd* from buffer *buff* starting at *buff*.

The *flags* argument is either 0 or is formed by logically OR'ing one or more of constant shown in table 10.1 below: -

| *Flags* | Description | recv | Send |
|---|---|---|---|
| MSG_DONTROUTE | Bypass routing table lookup | | * |
| MSG_DONTWAIT | Only this operation is nonblocking | * | * |
| MSG_OOB | Send or receive out-of-band data | * | * |
| MSG_PEEK | Peek at incoming message | * | |
| MSG_WAITALL | Wait for all the data | * | |

Table 10.1: flags for I/O function

## 10.5 *READV* AND *WRITEV* FUNCTIONS

These two functions are similar to read and write but, *readv and writev* let us read into or write from one or more buffers with a single function call. These operations are called *scatter read* and *gather write.*

*readv* function *iovcnt* blocks from the file associated with the file descriptor *fields* into the multiple buffers described by *iov*.

*writev* function writes at most *iovcnt* blocks described by the *iovec* to the file associated with the file descriptor *filedes*.

Format of readv and writev functions are given below:

> #include <sys/uio.h>
>
> ssize_t readv(int *filedes*, const struct *iovec *iov*, int *iovcnt*);
>
> ssize_t writev(int *filedes*, const struct *iovec *iov*, int *iovcnt*);
>
> Both return: number of bytes read or write if OK, -1 on error

Here, iovec is a structure defined to denote buffer starting address and its size.

> struct iovec{
>
>     void *iov_base;    /*starting address of buffer*/
>
>     size_t iov_len;/*size of buffer*/
>
> };

Header file for these operations is <sys/uio.h>.

---

### Check Your Progress

1.  *Explain the recv( ) and send( ) function with syntax.*

2.  *Write the difference between readv( ) and writev( ) function.*

---

## 10.6 *RECVMSG* AND *SENDMSG* FUNCTIONS

These two functions are the most general of all the I/O functions. We could replace all calls to *read, readv, recv, and recvfrom* with calls to *recvmsg.* Similarly, all calls to the various output functions could be replaced with calls to *sendmsg.* Header file to be included is <sys/socket.h>

Format of recvmsg and sendmsg functions are given below:

> #include <sys/socket.h>
>
> ssize_t recvmsg(int *sockfd*, struct msghdr *msg*, int *flags*);
>
> ssize_t sendmsg(int *sockfd*, struct msghdr *msg*, int *flags*);

msghdr structure is defined as: -

```
struct msghdr{
        void    *msg_name;                  /*protocol address*/
        socklen_t       msg_namelen; /*size of protocol address*/
        struct  iovec   *msg_iov;               /*scatter/gather
array*/
        int     msg_iovlen;                 /*# e    lements in
msg_iov*/
        void    *msg_control;           /*ancillary d   ata    (cmsghdr
struct) */
        socklen_t       msg_controllen;      /*length o  f a  ncillary
data*/
        int     msg_flags;                  /*flags r   eturned b  y
recvmsg ()*/
        };
```

Here, *msg_name* and *msg_namelen* members are used when the socket is not c onnected. *msg_iov* and *msg_iovlen* members s pecify t he array of input or out put buf fers ( the a rray of *iovec* structures). *msg_control* and *msg_controllen* members s pecify t he l ocation a nd s ize of t he opt ional ancillary data. *msg_flags* member is used only by *recvmsg* while ignored by *sendmsg.*

Summary of t he flags t hat a re e xamined b y t he kernel for t he i nput a nd output f unctions, a s w ell a s t he *msg_flags* that might be r eturned b y *recvmsg* is shown in table 10.2 below:

| 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 | | | * |
| MSG_NOTIFICATION | | | * |

Table 10.2: Summary of input and output flags by various I/O functions

The first four flags are only examined and never returned; the next two are both examined and returned; and the last four are only returned.

## 10.7 ANCILLARY DATA

Control messages or control information is also called as ancillary data. Ancillary Data can be sent and received using the *msg_control* and *msg_controllen* members of the *msghdr* structure with the *sendmsg* and *recvmsg* functions. Summary of the various uses of ancillary data is shown in table 10.3 below.

| Protocol | cmag_level | cmag_type | Description |
|---|---|---|---|
| IPv4 | IPPROTO_IP | IP_RECVDSTAD DR | Receive destination address with UDP datagram |
| | | IP_RECVIF | Receives interface index with UDP datagram |
| IPv6 | IPPROTO_IPv6 | IPv6_DSTOPTS | Specify destination options |
| | | IPv6_HOPLIMIT | Specify hop limit |
| | | IPv6_HOPOPTS | Specify hop-by-hop options |
| | | IPv6_NEXTHOP | Specify next-hop address |
| | | IPv6_PKTINFO | Specify packet information |
| | | IPv6_RTHDR | Specify routing header |
| | | IPv6_TCLASS | Specify traffic class |
| Unix domain | SQL_SOCKET | SCM_RIGHTS | Send/receive descriptors |
| | | SCM_CREDS | Send/receive user credentials |

Table 10.3: Summary of uses of ancillary data

Ancillary data consists of one or more *ancillary data objects*, each one beginning with a *cmsghdr* structure, defined by including*<sys/socket.h>*.

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

*The following five macros are defined by including the <sys/socket.h> header to simplify processing of the ancillary data: -*

```
#include <sys/socket.h>
#include < sys/param.h>/* f or A LIGN m acro on m  any
```
implementations */
```
struct cmsghdr *CMSG_FIRSTHDR(struct msghdr *mhdrptr) ;
```

> Returns: poi nter t o f irst c msghdr s tructure o r N ULL i f no ancillary data

```
struct cmsghdr *CMSG_NXTHDR(struct msghdr *mhdrptr, struct
cmsghdr *cmsgptr) ;
```

> Returns: poi nter t o ne xt cmsghdr s tructure or N ULL i f no more ancillary data objects

```
unsigned char *CMSG_DATA(struct cmsghdr *cmsgptr) ;
```

> Returns: p ointer to f irst b yte o f d ata a ssociated w ith cmsghdr structure

```
unsigned int CMSG_LEN(unsigned int length) ;
```

> Returns: va lue t o s tore in c msg_len g iven t he a mount of
data

```
unsigned int CMSG_SPACE(unsigned int length) ;
```

> Returns: to tal s ize o f an a ncillary data o bject given th e amount of data

---

### *Check Your Progress*

1.  *Write a brief note on flags that are examined by the kernel for the input and output functions.*

2.  *What are the various uses of ancillary data?*

---

## 10.8  HOW MUCH DATA IS QUEUED?

There are times when w e want to see how much data is queued to be r ead on a  s ocket, without r eading t he  data. There a re f ollowing techniques has to covered.

* If the goal is not to block in the kernel because we have something else to do w hen nothing i s ready t o be  read, non blocking I/O can be used.

* If we w ant to e xamine th e d ata b ut s till le ave i t o n th e  receive queue for s ome ot her pa rt of our pr ocess t o r ead, we c an us e the *MSG_PEEK* flag.

* Some implementations support the *FIONREAD* command of *ioctl.*

## 10.9  SOCKETS AND STANDARD I/O

One of the methods of performing I/O is the *standard I/O library*. It is specified by the ANSI C standard. The standard I/O library handles some of the details such as automatically buffering the input and output streams.

The standard I/O library can be used with sockets, but there are a few things to consider: -

- A standard I/O stream can be created from any descriptor by calling the *fdopen* function. Similarly, given a standard I/O stream, we can obtain the corresponding descriptor by calling *fileno*.

- TCP and UDP sockets are full-duplex. Standard I/O streams can also be full-duplex.

- The easiest way to handle this read-write problem is to open two standard I/O streams for a given socket: one for reading and one for writing.

  Standard I/O uses three types of buffering: -

- *Fully buffered*: I/O takes place only when the buffer is full, or the process calls *fflush* or *exit*.

- *Line buffered*: I/O takes place only when a new line is encountered, or the process calls *fflush* or *exit.*

- *Unbuffered*: I/O takes place each time a standard I/O output function is called.

Most UNIX implementations of the standard I/O library use the following rules:

- Standard error is always unbuffered.

- Standard input and standard output are fully buffered, unless they refer to a terminal device, in which case, they are line buffered.

- All other streams are fully buffered unless they refer to a terminal device, in which case, they are line buffered.

## 10.10  SUMMARY

There are three main ways to set a time limit on a socket operation:

**(i)**    Use the *alarm* function and the *SIGALRM* signal,

**(ii)**   Use the time limit that is provided by *select* and

**(iii)**  Use the newer *SO_RCVTIMEO* and *SO_SNDTIMEO* socket options.

*recvmsg* and *sendmsg* are t he most g eneral o f a ll th e I/O f unctions provided. We know the various uses of ancillary data.

## 10.11  TERMINAL QUESTIONS

1. Write t he t hree w ays of pe rforming I/O ope rations involving sockets.

2. Write syntax/format of (i) readv and writev function, (ii) recvmsg and sendmsg function.

3. What are the three types of buffering used by standard I/O?

4. What are the different methods to check queued data?

5. What are the different rules used by most UNIX implementations of the standard I/O library?

6. What ar e d ifferent macros a re de fined b y i ncluding t he <sys/socket.h> header to simplify processing of the ancillary data?

# UNIT-11 : UNIX DOMAIN PROTOCOLS

## Structure

## 11.1  INTRODUCTION

The UNIX domain protocols are not an actual protocol suite, but a way of performing client/server communication on a single host using the same API that is used for clients and servers on different hosts. The UNIX domain protocols are an alternative to the interprocess communication (IPC) when the client and server are on the same host. It is important to note that the protocol addresses used to identify clients and servers in the UNIX domain are pathnames within the normal file system.

**Two types of sockets provided in UNIX domain: -**

- Stream sockets (similar to TCP).

- Datagram sockets (similar to UDP).

**Reason for using UNIX domain sockets: -**

- UNIX domain sockets are often twice as fast as a TCP socket when both peers are on same host.

- UNIX domain sockets are used when passing descriptors between processes on the same host.

- Newer implementations of UNIX domain sockets provide the client's credentials (user ID and group IDs) to the server, which can provide additional security checking.

Also, we will learn about UNIX domain socket address structure, socket pair function and socket function are discussed to provide the insight view. We will see programs of UNIX domain stream/datagram client/server. Then there is discussion on topics descriptors passing and receiving sender credentials.

## 11.2 OBJECTIVES

At the end of this unit we get to know about: -

- What are UNIX domain protocols?

- Unix domain socket address structure

- Format of socket pair and socket functions

- UNIX domain stream client/server and UNIX domain datagram client/server

- Descriptors passing and receiving sender credentials

## 11.3 UNIX DOMAIN SOCKET ADDRESS STRUCTURE

UNIX domain socket address structure is defined by including the <*sys/un.h*>header.

       struct sockaddr_un{

             sa_family_t          sun_family;
*/\*AF_LOCAL\*/*

             char          sun_path[104];     */\*null-terminated pathname\*/*

       };

(Unix domain socket address structure: socketaddr_un )

The pathname stored in the sun_path array must be null-terminated. The macro SUN_LEN is provided and it takes a pointer to a sockaddr_un structure and returns the length of the structure, including the number of non-null bytes in the pathname.

***The below program creates a unix domain socket, binds a pathname to it and then calls get sockname and prints the bound pathname.***

    #include "unp.h"

    int

    main (int argc, char **argv)

    {

```
int sockfd;

socklen_t len;

struct sockaddr_un addr1, addr2;

if (argc != 2)

err_quit("usage: unixbind <pathname>");

sockfd = Socket(AF_LOCAL, SOCK_STREAM, 0);

unlink(argv[1]); /* OK if this fails */

bzero(&addr1, sizeof(addr1));

addr1.sun_family = AF_LOCAL;

strncpy(addr1.sun_path, argv[1], sizeof(addr1.sun_path) - 1);

Bind(sockfd, (SA *) &addr1, SUN_LEN(&addr1));

len = sizeof(addr2);

Getsockname(sockfd, (SA *) &addr2, &len);

printf("bound name = %s, returned len = %d\n", addr2.sun_path,
len);

exit(0);

}
```

In the program, the pathname that bind to the socket is the command-line argument. But the bind will fail if the pathname already exists in the filesystem. Therefore, it calls unlink to delete the pathname, in case it already exists. If it does not exist, unlink returns an error, which ignore bind and then getsockname. Copy the command-line argument using strncpy, to avoid overflowing the structure if the pathname is too long. Since initialize the structure to zero and then subtract one from the size of the sun_path array, the pathname is null-terminated. After that bind is called and use the macro SUN_LEN to calculate the length argument for the function. Then call getsockname to fetch the name that was just bound and print the result.

## 11.4 *SOCKETPAIR* FUNCTION

This function creates two sockets that are then connected together. This function applies only to UNIX domain sockets. It includes header *<sys/socket.h>*

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

int socketpair(int *family*, int *type*, int *protocol*, int *sockfd[2]*);

Here,          family -          AF_LOCAL

protocol-        0

type -   SOCK_STREAM or SOCK_DGRAM

and two s ocket d escriptor t hat are r eturned are sockfd[0] and sockfd[1].

The result of socketpair with a type of SOCK_STREAM is called a stream pipe. It is similar to a regular Unix pipe, but a stream pipe is full-duplex; that is both descriptors can be read and written.

---

***Check Your Progress:***

*1.* *Write the difference between Unix pipe and stream pipe.*

*2.* *Can you create a function for socket pair?*

---

# 11.5 SOCKET FUNCTIONS

While using UNIX domain sockets there exists several differences and r estrictions i n t he s ocket f unctions. B elow i s t he l ist of P OSIX requirements w hen a pplicable a nd n ote th at n ot all imp lementations a re currently at this level.

- The default file access permissions for a pathname created by bind should be 0777, modified by the current unmask value.

- The pathname associated with a UNIX domain socket should be an absolute pathname, not a relative pathname.

- The p athname s pecified in a cal l t o c onnect m ust be a pa thname that i s c urrently bound to a n ope n UNIX domain s ocket of t he same type (stream or datagram).

- The p ermission te sting associated w ith th e *connect* of a UNIX domain socket is the same as if *open* had been called for write-only access to the pathname.

- UNIX domain stream sockets provide a byte stream interface to the process with no record boundaries.

- If s ocket's qu eue i s f ull t hen E CONNREFUSED i s r eturned i n response to a call to *connect* for UNIX domain stream socket.

- UNIX domain datagram sockets are similar to UDP sockets.

- Unlike U DP s ockets, s ending a da tagram on a n unbound UNIX domain datagram socket does not bind a pathname to the socket.

## 11.6 UNIX DOMAIN STREAM CLIENT/SERVER

UNIX domain stream client/server uses stream socket. Most of the steps are similar to TCP echo client/server. But some modifications have been down n l ike associating p athname t o UNIX. The s teps f or cr eating UNIX domain stream server are following: -

- Call s ocket ( ) - A c all to socket ( ) with t he p roper a rguments creates the UNIX socket. Here, we will pass SOCK_STREAM as second argument to create a stream socket.

- We first unlink the pathname, in case it e xists from an e arlier run of the server, and then initialize the socket address structure before calling bind (). Rest of steps is same as of TCP echo client/server.

- Call bind ( )- We get a socket descriptor from the c all to socket(), now bind that to an address in the UNIX domain.

- Call listen ( ) - This i nstructs th e s ocket to lis ten f or in coming connections from client programs.

- Call accept ( ) - This will accept a connection from a client.

- Close the connection.

In o rder t o cr eate UNIX domain s tream c lient s ome mo difications h ave been made to TCP client/socket. The steps are given below: -

- The socket address structure to contain the server's address is now a *sockaddr_un* structure.

- Call s ocket ( ) -The f irst a rgument t o s ocket i s A F_LOCAL a nd second one is SOCK_STREAM.

- Call connect ( ) – To connect with server.

***The program below shows the unix domain stream protocol echo server.***

```
#include   "unp.h"
int
main (int argc, char **argv)
{
int listenfd, connfd;
pid_t childpid;
socklen_t clilen;
struct sockaddr_un cliaddr, servaddr;
void sig_chld(int);
listenfd = Socket(AF_LOCAL, SOCK_STREAM, 0);
```

```
unlink(UNIXSTR_PATH);

bzero(&servaddr, sizeof(servaddr));

servaddr.sun_family=AF_LOCAL;

strcpy(servaddr.sun_path, UNIXSTR_PATH);

Bind(listenfd, (SA *) &servaddr, sizeof(servaddr));

Listen(listenfd, LISTENQ);

Signal(SIGCHLD, sig_chld);

for ( ; ; ) {

clilen = sizeof(cliaddr);

if ( (connfd = accept(listenfd, (SA *) &cliaddr, &clilen)) < 0) {

if (errno == EINTR)

continue; /* back to for() */

else

err_sys("accept error");

}

if ( (childpid = Fork()) == 0) { /* child process */

Close(listenfd); /* close listening socket */

str_echo(connfd); /* process request */

exit(0);

}

Close(connfd);

}

}
```

The program of the server is use the Unix domain stream protocol instead of T CP. The d atatype of t he t wo s ocket a ddress s tructures i s no w sockaddr_un. The f irst argument t o s ocket i s AF_LOCAL, t o cr eate a Unix domain stream socket. The constant UNIXSTR_PATH is defined in unp.h to be /tmp/unix.str. First unlink the pathname, in case it exists from an earlier run of the server and then initialize the socket address structure before cal ling bind. A n e rror fro m unlink is acc eptable. The s tream protocol echo server program in bind call it specify the size of the socket address structure (the third argument) as the total size of the sockaddr_un structure, not j ust t he num ber of b ytes oc cupied b y t he pa thname. B oth lengths are valid since the pathname must be null-terminated.

**The program below shows the Unix domain stream protocol echo client**

```
#include "unp.h"

int

main(int argc, char **argv)

{

int sockfd;

struct sockaddr_un servaddr;

sockfd = Socket(AF_LOCAL, SOCK_STREAM, 0);

bzero(&servaddr, sizeof(servaddr));

servaddr.sun_family = AF_LOCAL;

strcpy(servaddr.sun_path, UNIXSTR_PATH);

Connect(sockfd, (SA *) &servaddr, sizeof(servaddr));

Str_cli(atdin,sockfd);

Exit(0);

}
```

The s ocket a ddress s tructure t o c ontain t he s erver's a ddress i s a sockaddr_un structure. The first argument to socket is AF_LOCAL. The code to fill in the socket address structure is identical to the code shown in the pr evious pr ogram for t he server: Initialize th e s tructure to 0, s et th e family to AF_LOCAL, and copy the pathname into the sun_path member.

# 11.7 UNIX DOMAIN DATAGRAM CLIENT /SERVER

UNIX domain d atagram c lient/server requires mo dification to UDP e cho c lient/server. Important poi nts i n creating UNIX domain datagram echo server are following: -

- The da tatype of t he two s ocket a ddress s tructures i s now *sockaddr_un*.

- The f irst ar gument t o s ocket i s A F_LOCAL, t o cr eate a UNIX domain datagram socket.

- We first unlink the pathname, in case it e xists from an earlier run of the server, and then initialize the socket address structure before calling bind ().

- Others are same as UDP echo server.

Similarly, in UNIX domain datagram echo client, some modifications have been done.

- The socket address structure to contain the server's address is now a *sockaddr_un* structure. We also allocate one of these structures to contain the client's address.

- The first argument to socket is AF_LOCAL.

- Unlike our UDP client, when using the UNIX domain datagram protocol, we must explicitly bind a pathname to our socket so that the server has a pathname to which it can send its reply. Other is same as UDP echo client.

**The program below shows the Unix domain datagram protocol echo server.**

```
#include "unp.h"

int

main(int argc, char **argv)

{

int sockfd;

struct sockaddr_un servaddr, cliaddr;

sockfd = Socket(AF_LOCAL, SOCK_DGRAM, 0);

unlink(UNIXDG_PATH);

bzero(&servaddr, sizeof(servaddr));

servaddr.sun_family = AF_LOCAL;

strcpy(servaddr.sun_path, UNIXDG_PATH);

Bind(sockfd, (SA *) &servaddr, sizeof(servaddr));

dg_echo(sockfd, (SA *) &cliaddr, sizeof(cliaddr));

}
```

In the program, the datatype of the two socket address structures is now sockaddr_un. The first argument to socket is AF_LOCAL, to create a Unix domain datagram socket. The constant UNIXDG_PATH is defined in unp.h to be /tmp/unix.dg. First unlink the pathname, in case it exists from an earlier run of the server, and then initialize the socket address structure before calling bind. An error from unlink is acceptable. The dg_echo function is used.

***The program below shows the Unix domain datagram protocol echo client***

#include "unp.h"

```
int

main(int argc, char **argv)

{

int sockfd;

 struct sockaddr_un cliaddr, servaddr;

 sockfd = Socket(AF_LOCAL, SOCK_DGRAM, 0);

bzero(&cliaddr,sizeof(cliaddr));

cliaddr.sun_family=AF_LOCAL;

strcpy(cliaddr.sun_path,tmpnam(NULL));

Bind(sockfd, (SA *) &cliaddr, sizeof(cliaddr));

bzero(&servaddr, sizeof(servaddr));

servaddr.sun_family=AF_LOCAL;

strcpy(servaddr.sun_path,UNIXDG_PATH);

dg_cli(stdin,sockfd,(SA*) &servaddr,sizeof(servaddr));

exit(0);

}
```

In the program, the socket address structure to contain the server's address is now a sockaddr_un structure. Also allocate one of these structures to contain the client's address. The first argument to socket is AF_LOCAL. Unlike our UDP client, when using the Unix domain datagram protocol, then must explicitly bind a pathname to socket so that the server has a pathname to which it can send its reply. The code to fill in the socket address structure with the server's well-known pathname is identical to the code shown earlier for the server. The function dg_cli is the used.

---

### Check Your Progress

1.  *Write the difference between Unix domain datagram and Unix domain stream client/server.*

---

## 11.8 PASSING DESCRIPTORS

Steps involved in passing a descriptor between two processes are as follows: -

*   Create UNIX domain sockets either a stream socket or a datagram socket and connect them for communication between a server and a client.

*   One process opens a descriptor. Any type of descriptor can be exchanged.

- Sender builds an *msghdr* structure containing the descriptor to be passed, and calls *sendmsg* with the structure across one of the UNIX domain sockets.

- Reciever calls *recvmsg* to receive the descriptor from the other UNIX domain socket.

Client and server must have an application protocol so they know when the descriptor is to be passed.

## 11.9   RECEIVING SENDER CREDENTIALS

When a client and server communicate using UNIX domain protocols, the server often needs a way to know exactly who the client is, to validate that the client has permission to ask for the service being requested.

FreeBSD passes credentials in a *cmsgcred* structure, which is defined by including the *<sys/socket.h>* header.

*structcmsgcred {*

    *pid_t cmcred_pid;*        */\* PID of sending process \*/*

    *uid_t cmcred_uid;*        */\* real UID of sending process \*/*

    *uid_t cmcred_euid;*   */\* effective UID of sending process \*/*

    *gid_t cmcred_gid;*        */\* real GID of sending process \*/*

    *short cmcred_ngroups;*    */\* number of groups \*/*

    *gid_t cmcred_groups[CMGROUP_MAX];*
        */\* groups \*/*

    *};*

**The program below shows the read_cred function that reads and returns sender's credentials.**

```
#include "unp.h"

#defineCONTROL_LEN(sizeof(structcmsghdr)+sizeof(struct
cmsgcred))

ssize_t

read_cred(int fd, void *ptr, size_t nbytes,  struct cmsgcred
*cmsgcredptr)

{
```

```c
struct msghdr msg;

struct iovec iov[1];

char control[CONTROL_LEN];

int n;

msg.msg_name = NULL;

msg.msg_namelen = 0;

iov[0].iov_base = ptr;

iov[0].iov_len = nbytes;

msg.msg_iov = iov;

msg.msg_iovlen = 1;

msg.msg_control = control;

msg.msg_controllen = sizeof(control);

msg.msg_flags = 0;

if ( (n = recvmsg(fd, &msg, 0)) < 0)

return (n);

cmsgcredptr->cmcred_ngroups = 0;  / * i ndicates no c redentials
returned */

if (cmsgcredptr && msg.msg_controllen > 0) {

struct cmsghdr *cmptr = (struct cmsghdr *) control;

if (cmptr->cmsg_len < CONTROL_LEN)

err_quit("control length = %d", cmptr->cmsg_len);

if (cmptr->cmsg_level != SOL_SOCKET)

err_quit("control level != SOL_SOCKET");

if (cmptr->cmsg_type != SCM_CREDS)

err_quit("control type != SCM_CREDS");

memcpy(cmsgcredptr, C    MSG_DATA(cmptr), s    izeof(struct
cmsgcred));

}

return (n);

}
```

In t he program, the f irst th ree a rguments a re id entical to r ead, w ith th e
fourth argument being a pointer to a cmsgcred structure that will be filled
in. If credentials were returned, the length, level, and type of the ancillary

data are verified, and the resulting structure is copied back to the caller. If no credentials were returned, then set the structure to 0. Since the number of groups (cmcred_ngroups) is always 1 or more, the value of 0 indicates to the caller that no credentials were returned by the kernel. The main function for echo server, str_echo function is called by the child after the parent has accepted a new client connection and called fork. If credentials were returned, they are printed. The further code reads buffers from the client and writes them back to the client.

***The program below shows the str_echo function that asks for client credentials***

```
#include "unp.h"

ssize_t read_cred(int, void *, size_t, struct cmsgcred *);

void

str_echo(int sockfd)

{

ssize_t n;

int i;

char buf[MAXLINE];

struct cmsgcred cred;

again:

while ( (n = read_cred(sockfd, buf, MAXLINE, &cred))>0) {

if (cred.cmcred_ngroups == 0) {

printf("(no credentials returned)\n");

} else {

printf("PID of sender = %d\n", cred.cmcred_pid);

printf("real user ID = %d\n", cred.cmcred_uid);

printf("real group ID = %d\n", cred.cmcred_gid);

printf("effective user ID = %d\n", cred.cmcred_euid);

printf("%d groups:", cred.cmcred_ngroups - 1);

for (i = 1; i < cred.cmcred_ngroups; i++)

printf(" %d", cred.cmcred_groups[i]);

printf("\n");

}

writen(sockfd, buf, n);
```

```
        }

        if (n < 0 && errno == EINTR)

        goto again

        else if(n<0)

           err_sys("str_echo:read error");

        }
```

In t he program, the m ain function for e cho s erver i s s tr_echo function. This function i s c alled by t he c hild after t he pa rent ha s a ccepted a ne w client c onnection a nd called fork.   If credentials w ere returned, t hey ar e printed. Further code reads buffers from the client and writes them back to the client. Here client is to pass an empty cmsgcred structure that will be filled in when it calls sendmsg.

## 11.10   SUMMARY

       UNIX domain sockets are an alternative to IPC when the client and server are on the same host. The advantage in using UNIX domain sockets over s ome form of IPC i s that t he A PI i s nearly i dentical to a n etworked client/server. We m odified our T CP and U DP echo c lients and s ervers to use the UNIX domain protocols and the only major difference had to bind a pa thname t o t he U DP c lient's s ocket, s o t hat  the  UDP s erver h ad somewhere to send the replies. Descriptor passing is a powerful technique between cl ients an d servers o n the s ame h ost an d i t t akes p lace ac ross a UNIX domain socket.

## 11.11   TERMINAL QUESTIONS

1.    Explain UNIX domain protocol.

2.    Write two types of socket provided in UNIX domain socket.

3.    Write structure of UNIX domain socket address.

4.    What is use of Bind ( ) system call?

5.    Write a program to show read_cred function that reads and returns sender's credentials.

6.    Write a program to show the str_echo function that asks for client credentials.

**Bachelor of Computer Application**

**BCA-E7**
**Network Programming**

**Uttar Pradesh Rajarshi Tandon**
**Open University**

**Block**

# 4

**Broadcast, Multicast and Inter Process Communication**

# Course Design Committee

| | |
|---|---|
| **Dr. Ashutosh Gupta** | Chairman |
| Director (In-charge) | |
| School of Compruter and Information Science, UPRTOU Prayagraj | |
| **Prof. R. S. Yadav** | Member |
| Department of Computer Science and Engineering | |
| MNNIT-Allahabad, Prayagraj | |
| **Ms Marisha** | Member |
| Assistant Professor (Computer Science), | |
| School of Science UPRTOU Prayagraj | |
| **Mr. Manoj Kumar Balwant** | Member |
| Assistant Professor, (Computer Science) | |
| School of Sciences UPRTOU Prayagraj | |

# Course Preparation Committee

| | |
|---|---|
| **Dr. Prabhat Kumar** | Author (Block 1,2) |
| Assistant Professor, Department of IT | |
| NIT Patna | |
| **Dr. Prabhat Ranjan** | Author (Block 3,4) |
| Assistant Professor, Department of Computer Science | |
| Central University of South Bihar | |
| **Dr. Rajiv Mishra** | Editor |
| Associate Professor, Department of CSE | |
| IIT Patna | |
| **Dr. Ashutosh Gupta** (Director in Charge) | |
| School of Computer & Information Sciences, | |
| UPRTOU Prayagraj | |
| **Mr. Manoj Kumar Balwant** | Coordinator |
| Assistant Professor, (Computer Science) | |
| School of sciences UPRTOU Prayagraj | |

# BLOCK INTRODUCTION

The objective of this course is to introduce the basic concept about the network programming as well as provides a mix of practical experience and a depth of understanding. The network programming course address today's most crucial standards, implementations and techniques. The aim is to provide an extensive variety of topics on this subject with appropriate examples. The course is organized into following blocks:

Block 4 covers broadcasting, multicast, inter process communication and remote login.

# UNIT 12 : BROADCASTING

## Structure

## 12.1  INTRODUCTION

In th is unit, we w ill l earn about; br oadcasting and i ts us es, broadcast a ddress, di fference be tween uni cast a nd br oadcast addressing, *dg_cli* function using broadcasting and race conditions.

There ar e f our types of a ddressing; U nicast, Anycast, M ulticast a nd Broadcast.  Unicasting i s pr ocessing t alking to e xactly one another process, f or e xample T CP.  A ny casting i s a dded i n IPv6 a ddressing architecture. M ulticasting s upport i s opt ional i n IPv4 but m andatory i n IPv6.  Broadcasting i s not a vailable in I Pv6.  A ny I Pv6 a pplication th at uses br oad c asting m ust be r ecorded IPv6 t o use m ulticasting. B road casting an d m ulticasting require d ata gram t ransport s uch as U DP o r r aw IP, they cannot work with TCP.

## 12.2 OBJECTIVES

At the end of this unit, you should be able to: -

- Know what is broadcasting, its uses.

- How to write broadcast address.

- Able to differentiate between unicast and broadcast address.

- Able to write dg_cli function that broadcasts.

- Know race conditions.

## 12.3 BROADCAST ADDRESSES

Broadcasting refers to transferring a message to all the recipients. Broadcasting require datagram transport like UDP or raw IP, it cannot work with TCP.

---

### *Uses*

- To locate a server on the local subnet when the server is assumed to be on the local subnet but its unicast IP address is not known. This is sometimes called *resource discovery*.

- To minimize the network traffic on a LAN when there are multiple clients communicating with a single server.

---

If we denote an IPv4 address as { *subnetid*, *hostid*}, where *subnetid* represents the bits that are covered by the network mask (or the CIDR prefix) and *hostid* represents the bits that are not covered, then we have **two types of broadcast addresses**. We denote a field containing all one bits as –1.

**a)** *Subnet-directed broadcast address - <subnetid, -1>*: This addresses all the interfaces on the specified subnet. For example, if we have the subnet 192.168.42/24, then 192.168.42.255 would be the subnet-directed broadcast address for all interfaces on the192.168.42/24 subnet.

**b)** *Limited broadcast address -<-1, -1, -1>* or 255.255.255.255: Datagrams destined to this address must never be forwarded by a router.

---

### *Check Your Progress*

*1.* *Can you explain the subnetid and hostid?*

*2.* *Can you explain types of broadcast addresses?*

---

## 12.4 UNICAST VERSUS BROADCAST

Unicast is the term used to describe communication where a piece of information is sent from one point to another point. In this case there is just one sender, and one receiver. Unicast uses IP delivery methods such as Transmission Control Protocol (TCP) and User Datagram Protocol (UDP), which are session-based protocols. When a Windows Media Player client connects using unicast to a Windows Media server, that client has a direct relationship to the server. Each unicast client that connects to the server takes up additional bandwidth. For example, if you have 10 clients all playing 100-kilobits per second (Kbps) streams, those

clients as a group are taking up 1,000 K bps. If you have only one client playing the 100 Kbps stream, only 100 Kbps is being used.

Broadcast is the term used to describe communication where a piece of information is sent from one point to all other point. In this case there is just one sender, but the information is sent to all connected receivers.

## 12.5  DG_CLI FUNCTION USING BROADCASTING

The *dg_cli* function is used to perform most of the client processing in UDP echo client. In order to broadcast to the standard UDP daytime server and printing all replies we make some modifications to *dg_cli function*. In *main ( )* function we change the destination port number to 13.

> s*ervaddr.sin_port = htons(13);*

The working of *dg_cli* function is as follows: -

- Allocate room for server's address, set socket option.

- Read line; send to socket, read all replies.

- Print each received reply.

**The dg_cli function that broadcasts as shown below:**

> *#include "unp.h"*
>
> *static void recvfrom_alarm(int);*
>
> *void*
>
> *dg_cli(FILE *fp, int sockfd, const SA *pservaddr, socklen_t servlen)*
>
> *{*
>
>> *int n;*
>>
>> *const int on = 1;*
>>
>> *char sendline[MAXLINE], recvline[MAXLINE + 1];*
>
>> *socklen_t len;*
>
>> *struct sockaddr *preply_addr;*
>>
>> *preply_addr = Malloc(servlen);*
>>
>> *Setsockopt(sockfd, SOL_SOCKET, SO_BROADCAST, &on, sizeof(on));*
>>
>> *Signal(SIGALRM, recvfrom_alarm);*
>>
>> *while (Fgets(sendline, MAXLINE, fp) != NULL) {*
>>
>>> *Sendto(sockfd, sendline, strlen(sendline), 0, pservaddr, servlen);*

```
                        alarm(5);

                        for ( ; ; ) {

                                len = servlen;

                                n = recvfrom(sockfd, recvline, MAXLINE, 0,
                                preply_addr, &len);

                                if (n < 0) {

                                        if (errno == EINTR)

                                                break; /* waited long enough
for replies */

                                        else

                                                err_sys("recvfrom error");

                                } else {

                                        recvline[n] = 0; /* null terminate */

                                        printf("from %s: %s",

                                        Sock_ntop_host  (preply_addr,  len),
recvline);

                                        }

                                }

                        }

                        free(preply_addr);

                }

                static void

                recvfrom_alarm(int signo)

                {

                        return; /* just interrupt the recvfrom() */

                }
```

The dg_cli function sets the SO_BROADCAST socket option and prints all the replies received within five seconds. In the program, the malloc allocates room for the server's address to be returned by recvfrom. The SO_BROADCAST socket option is set and a signal handler is installed for SIGALRM.    The next two steps, fgets and sendto of this function are sending a broadcast datagram, receive multiple replies, call recvfrom in a loop and print all the replies received within five seconds. After five seconds, SIGALRM is generated, signal handler is called, and recvfrom returns the error EINTR. For each reply received, in the program call sock_ntop_host, which in the case of IPv4 returns a string containing the

dotted-decimal IP address of t he s erver. T his i s pr inted a long w ith th e server's reply.

A race condition is a situation which occurs usually when multiple processes are a ccessing data t hat i s s hared among t hem, b ut t he co rrect outcome depends on the execution order of the processes.

Race conditions are always a co ncern with threads programming since so much data is shared among all the threads (e.g., all the global variables). Race conditions of a different type often exist when dealing with signals. The pr oblem oc curs because a s ignal can n ormally b e d elivered at an y time while o ur p rogram is executing. POSIX allows us to bl ock a signal from being delivered, but this is often of little use while we are performing I/O operations.

A race condition exists in above program (*dg_cli* function using broadcast) if force the condition to occur as follows: Change the argument to alarm from 5 t o 1, a nd add sleep (1) immediately before the printf. When make these changes to the function and then type the first line of input, the line is sent as a broadcast and set the alarm for one second in the future. The block in the call to recvfrom, and the first reply then arrives for our socket, probably within a few milliseconds. The reply is returned by recvfrom, but we then g o to s leep for one s econd. A dditional r eplies are r eceived, an d they are pl aced i nto our sockets receive b uffer. But while w e are as leep, the alarm t imer ex pires an d t he S IGALRM s ignal i s g enerated: s ignal handler i s c alled, a nd i t just r eturns a nd i nterrupts the s leep i n w hich w e are b locked. Then l oop around a nd read the que ued r eplies w ith a one - second p ause e ach t ime w e p rint a r eply. W hen w e h ave read al l the replies, w e b lock a gain in th e c all to r ecvfrom, b ut th e timer is n ot running. T hus, w e w ill block f orever i n r ecvfrom. T he f undamental problem i s that our intent i s for our s ignal ha ndler to i nterrupt a bl ocked recvfrom, but t he s ignal c an b e de livered at a ny time, and w e c an be executing anywhere in the infinite for loop when the signal is delivered.

## 12.7  SUMMARY

In this unit, we have covered about broadcast addresses, difference between uni cast a nd br oadcast, *dg_cli* function u sing b roadcasting an d also studied race conditions.

Broadcasting sends datagram that all hosts on the attached subnet receive. While unicasting sends datagram to a single intended host.

There are two ways to write broadcast address

(i)    subnet-directed broadcast address and

**(ii)**   limited broadcast address.

## 12.8   TERMINAL QUESTIONS

1.   Explain broadcast address and its uses.

2.   Explain difference between unicast and broadcast?

3.   Explain dg_cli function using broadcasting.

4.   What is race condition?

5.   Write a dg_cli function with race condition.

# UNIT 13 : MULTICASTING

**Structure**

## 13.1  INTRODUCTION

In this unit, we will learn details about multicasting. We will study multicasting on  LAN and its difference with broadcasting. Then we get knowledge about multicasting on WAN.  Then we see difficulties related to multicasting on WAN and discuss solutions to it in terms of source specific multicast. Then we will have a look over multicast socket options and *mcast_join* and related functions. We will see *dg_cli* function using multicasting. Then there will be discussion on topics Receiving IP Multicast Infrastructure Session Announcements, Sending and Receiving, and Simple Network Time Protocol (SNTP).

Multicast address identifies a set of IP interfaces. A multicast datagram should be received by only those interfaces interested in the datagram, that is, by the interfaces on the host running applications wishing to participate

in the multicast group. Multicasting is used on a LAN or across a WAN. Indeed, applications multicast across a subset of internet on a daily basis.

## 13.2   OBJECTIVES

At the end of this unit, you should be able to know about: -

- Multicast and multicast address.
- Difference between Multicasting versus Broadcasting on a LAN
- Multicasting on a WAN and Source-Specific multicast
- Multicast Socket Options
- mcast_join and related functions and dg_cli Function using multicasting
- Receiving IP Multicast Infrastructure Session Announcements
- Sending and Receiving multicast datagrams
- Simple Network Time Protocol (SNTP)

## 13.3  MULTICAST ADDRESSES

A multicast address is a logical identifier for a group of hosts in a computer network that are available to process datagrams or frames intended to be multicast for a designated network service. Multicast addresses identify a set of IP interfaces. IP multicast address ranges and uses are shown below in table 13.1.

| Range Start Address | Range End Address | Description |
|---|---|---|
| 224.0.0.0 | 224.0.0.255 | Reserved for special "well-known" multicast addresses. |
| 224.0.1.0 | 238.255.255.255 | Globally-scoped (Internet-wide) multicast addresses. |
| 239.0.0.0 | 239.255.255.255 | Administratively-scoped (local) multicast addresses. |

Table 13.1 : IP Multicast Address Ranges and Uses

***IPv4 Multicast addresses***

In IPv4, class D addresses ranging from 224.0.0.0 to 239.255.255.255 are the multicast addresses. The lower order 28 bits of class D address form the multicast *groupID* and the 32-bit address is called the *group address.*

Mapping IPv4 multicast address to Ethernet address involves copy of low-order 23 bi ts of m ulticast a ddress t o l ow-order 23 -bits of t he E thernet address. T he hi gh o rder 24 bi ts of t he E thernet a ddress a re always *01:00:5e a*nd the next bit is always 0.

### IPv6 Multicast Addresses

The high-order byte of an IPv6 multicast address has the value *ff.*

The mapping from a 16-byte IPv6 multicast address into a 6-byte Ethernet address involves c opy o f l ow-order 32 bi ts of t he g roup a ddress i nto t he low-order 32 bi ts of t he E thernet address. T he h igh-order 2 b ytes o f t he Ethernet address are 33:33.

The table 13.2 and figure 13.1 shown the IPv6 Multicast Address Format.

| Field Name | Size (bits) | Description |
|---|---|---|
| *(Indicator)* | 8 | The f irst e ight bi ts a re a lways " 1111 1111" t o indicate a multicast address. |
| *Flags* | 4 | Four bi ts a re reserved for flags t hat c an be us ed t o indicate the nature of certain multicast addresses. At the p resent time, the f irst t hree o f t hese ar e u nused and s et t o z ero. T he f ourth i s t he *"T"* ( *Transient*) flag. If le ft as zero, this marks the multicast ad dress as a p ermanently-assigned, "well-known" mu lticast address. If s et to o ne, th is me ans th is is a *transient* multicast a ddress, me aning th at it is not permanently assigned. |
| *Scope ID* | 4 | Four bi ts a re us ed t o de fine t he s cope of the multicast a ddress; 1 6 d ifferent v alues from 0 to 15 are possible. <br><br> **Scope ID Value**     **Multicast Address Scope** <br> 0                        Reserved <br> 1                  Node-Local Scope <br> 2                  Link-Local Scope <br> 5                  Site-Local Scope <br> 8                  Organization-Local Scope <br> 14                Global Scope <br> 15                Reserved |
| *Group ID* | 112 | *Group ID:* Defines a p articular group w ithin e ach scope level. |

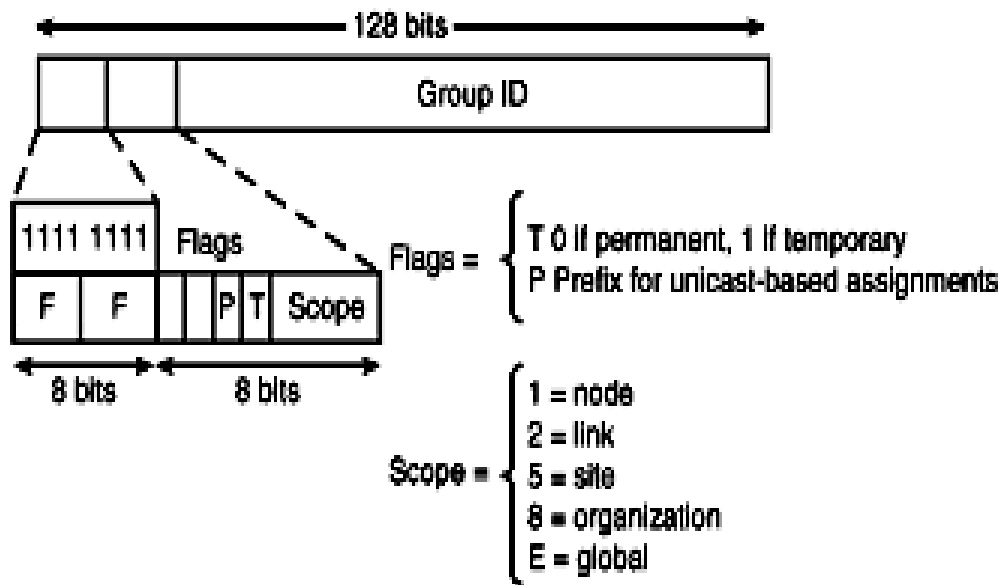Table 13.2 : IPv6 Multicast Address Format

Figure 13.1 : IPv6 Multicast Address Format

## 13.4    MULTICASTING VERSUS BROADCASTING ON A LAN

Broadcast is a term used to describe communication where a piece of information is sent from one point to all other points. In network case, there i s j ust one s ender, but the i nformation i s s ent t o a ll c onnected receivers. Broadcast i s mostly u sed i n l ocal sub ne tworks. In or der t o transmit b roadcast p acket, t he d estination M AC a ddress i s s et t o FF:FF:FF:FF:FF:FF and all such packets will be received by other NICs.

Multicast i s a t erm u sed t o d escribe communication w here a p iece o f information is sent from a source host to a group of destination hosts. The notion of gr oup i s e ssential t o t he c oncept of multicasting. A m ulticast group a ddress i s d efined. A ll t he hos ts t hat ha ve j oined t his group w ill receive messages sent to this multicast group address.

## 13.5   MULTICASTING ON A WAN

We know W AN pr ovides l ong di stance t ransmission of da ta, image, a udio a nd vi deo i nformation ove r l arge g eographical areas t hat may comprise a country, a continent, or even the whole world. WAN i.e wide area network is a combination of LANs connected through routers.

In order to have multicasting on a WAN we need to have multicast routers for connecting LAN. M ulticast r outers c ommunicate e ach ot her us ing multicast r outing pr otocol ( MRP). G roup of hos ts be longing t o di fferent LANs ma y f orm a mu lticast group. A n ew h ost c an jo in th e mu lticast group b y s ending a n IGMP t o a ny a ttached multicast router which t hen exchanges this information with other multicast routers using MRP.

Suppose a host on a LAN want to send a message to a multicast group on a WAN. It will multicast the message to its LAN. Other hosts on this LAN belonging t o t he required m ulticast g roup w ill r eceive th is me ssage. Multicast router attached to this LAN will also receive the message. This multicast r outer w ill th en s end th e me ssage to another multicast r outer attached to it. A ll mu lticast r outer w ill th en mu lticast o n its r espective LAN and multicast routers attached to it. Intended recipient on LAN will then receive message from multicast.

Multicasting on a WAN has been difficult to deploy for several reasons.

- The biggest problem is that the MRP needs to get the data from all the senders, which may be located anywhere in the network, to all receivers, which may similarly be located anywhere.

- Another l arge p roblem i s m ulticast ad dress al location. T here ar e not e nough IPv4 m ulticast a ddresses t o s tatically assign t hem t o everyone who wants one, as is done with unicast addresses.

# 13.6  SOURCE-SPECIFIC MULTICAST

Source-specific mu lticast (SSM) i s a m ethod of delivering multicast packets in which the only packets that are delivered to a r eceiver ar e t hose o riginating f rom a s pecific source address r equested by the receiver.

SSM combines t he group a ddress w ith a s ystem's s ource a ddress, w hich solves the multicasting problems in WAN by the following ways:

- The r eceivers s upply t he s ender's s ource a ddress to t he r outers a s part of joining the group.

- It r edefines th e id entifier f rom s imply being a mu lticast g roup address t o be ing a combination of a uni cast s ource a nd m ulticast destination.

IP version 4 ( IPv4) addresses in the 232/8 (232.0.0.0 to 232.255.255.255) range a re d esignated as s ource-specific mu lticast (SSM) d estination addresses and are reserved for use by source-    specific applications and protocols. F or IP ve rsion 6 ( IPv6), t he address p refix F F3x: :/ 32 i s reserved fo r s ource-specific mu lticast use. T his doc ument de fines a n extension to the Internet network service that applies to datagrams sent to SSM ad dresses an d d efines the hos t a nd router requirements t o s upport this extension.

---

### *Check Your Progress*

1. *Can you explain the difference between multicast and broadcast?*

2. *What are the advantages of SSM ?*

## 13.7  MULTICAST SOCKET OPTION

The API support for traditional multicasting requires only five new socket options. Source-filtering support, which is required for SSM, adds four more. The following are the multicast socket options: -

- IP_MULTICAST_IF – Specify d efault i nterface f or o utgoing multicasts.

- IP_MULTICAST_TTL – Specify TTL for outgoing multicasts.

- IP_MULTICAST_LOOP – Enable or disable loopback of outgoing multicasts.

- IPV6_MULTICAST_IF – Specify d efault i nterface f or out going multicasts.

- IPV6_MULTICAST_HOPS – Specify hop l imit f or out going multicasts.

- IPV6_MULTICAST_LOOP – Enable or di sable l oopback of outgoing multicasts.

Working with multicast sockets and UNIX (FreeBSD) as follows:

1.  Sending socket: - In general, there's nothing s pecial you ne ed to do on the sending end. The key is simply to send to a multicast IP (group) address. Tips:

    ❖ Use socket ( ) with AF _INET and S OCK_DGRAM arguments as normal.

    ❖ Use bind () to associate this socket with a l ocal address and port.

    ❖ Do n ot a ttempt to a ssociate th e s ocket w ith a mu lticast destination address using connect ().

    ❖ Use sendto () for sending data.

2.  Receiving socket: - Receiving i s ne arly t he s ame, but w ith one additional system call setsockopt ().

    ❖ Use socket ( ) with AF _INET and S OCK_DGRAM arguments as normal.

    ❖ Use setsockopt () with the IP_ADD_MEMBERSHIP option. This tells the system to receive packets on the network whose destination is the group address (but not its own).

## 13.8  *mcast_join* and RELATED FUNCTIONS

The mu lticast s ocket o ptions f or IPv4 a re s imilar to th e mu lticast socket opt ions f or IPv6, t here are e nough di fferences t hat pr otocol-

independent c ode using mu lticasting b ecomes complicated w ith lo ts o f #ifdefs. A b etter s olution is to h ide th e d ifferences w ithin th e f ollowing eight functions:

| #include "unp.h" |
|---|

| int mc ast_join(int *sockfd*, c onst s truct s ockaddr * *grp*, s ocklen_t *grplen*, const char \**ifname*, u_int *ifindex*); |
|---|

| int mcast_leave(int *sockfd*, const struct sockaddr \**grp*, socklen_t *grplen*); |
|---|

| int m cast_block_source(int *sockfd*, const s truct s ockaddr * *src*, socklen_t *srclen*, const struct sockaddr \**grp*, socklen_t *grplen*); |
|---|

| int m cast_unblock_source(int *sockfd*, c onst struct s ockaddr * *src*, socklen_t *srclen*, const struct sockaddr \**grp*, socklen_t *grplen*); |
|---|

| int m cast_join_source_group(int *sockfd*, c onst s truct s ockaddr * *src*, socklen_t *srclen*, c onst struct s ockaddr \**grp*, s ocklen_t *grplen*, c onst c har \**ifname*, u_int *ifindex*); |
|---|

| int m cast_leave_source_group(int *sockfd*, c onst s truct s ockaddr * *src*, socklen_t *srclen*, const struct sockaddr \**grp*, socklen_t *grplen*); |
|---|

| int mcast_set_if(int *sockfd*, const char \**ifname*, u_int *ifindex*); |
|---|

| int mcast_set_loop(int *sockfd*, int *flag*); |
|---|

| int mcast_set_ttl(int *sockfd*, int *ttl*); |
|---|

| All above return: 0 if OK, −1 on error |
|---|

| int mcast_get_if(int *sockfd*); |
|---|

| Returns: non-negative interface index if OK, −1 on error |
|---|

| int mcast_get_loop(int *sockfd*); |
|---|

| Returns: current loopback flag if OK, −1 on error |
|---|

| int mcast_get_ttl(int *sockfd*); |
|---|

| Returns: current TTL or hop limit if OK, −1 on error |
|---|

- *mcast_join* joins the any-source multicast group whose IP address is within the socket address structure pointed to by *grp,* and whose length is specified by *grplen.*

- *mcast_leave* leaves t he m ulticast g roup w hose IP a ddress i s contained within the socket address structure pointed to by *grp.*

- *mcast_block_source* blocks r eception on t he given s ocket of t he source and group whose IP address are contained within the socket address s tructures poi nted t o b y *src and grp,* respectively, an d whose lengths are specified by *srclen* and *grplen.*

- *mcast_unblock_source* unblocks reception of traffic from the given source to the given group.

- *mcast_join_source_group* joins t he s ource-specific g roup w here the source and group IP addresses are contained within the socket address s tructures poi nted t o b y *src and grp,* respectively, an d whose lengths are specified by *srclen* and *grplen.*

- *mcast_leave_source_group* leaves the source-specific group whose source a nd group IP a ddresses a re c ontained w ithin t he s ocket address s tructures poi nted t o b y *src and grp,* respectively, an d whose lengths are specified by *srclen* and *grplen.*

- *mcast_set_if* sets the default interface index for outgoing multicast datagrams.

- *mcast_set_loop* sets t he l oopback opt ion t o e ither 0 or 1, a nd *mcast_set_ttl* sets either the IPv4 TTL or the IPv6 hop limit.

***The program below shows the first third of mcast_join function. The program shows how straightforward the protocol-independent API can be.***

```
#include    "unp.h"

#include    <net/if.h>

int

mcast_join(int sockfd, const SA *grp, socklen_t grplen,

       const char *ifname, u_int ifindex)

{

#ifdef MCAST_JOIN_GROUP

   struct group_req req;

   if (ifindex > 0) {
```

```
            req.gr_interface = ifindex;
        } else if (ifname != NULL) {
            if ( (req.gr_interface = if_nametoindex(ifname)) == 0) {
                errno = ENXIO;      /* i/f name not found */
                return (-1);
            }
        } else
            req.gr_interface = 0;
        if (grplen > sizeof(req.gr_group)) {
            errno = EINVAL;
            return -1;
        }
        memcpy(&req.gr_group, grp, grplen);
        return (setsockopt(sockfd, family_to_level(grp->sa_family),
                    MCAST_JOIN_GROUP, &req, sizeof(req)));
#else
```

In the program, the caller is supplied an index, and then use it directly. Otherwise, if the caller supplied an interface name, the index is obtained by calling if_nametoindex. Otherwise, the interface is set to 0, telling the kernel to choose the interface. The caller's socket address is copied directly into the request's group field. Recall that the group field is a sockaddr_storage, so it is big enough to handle any socket address type the system supports. However, to guard against buffer overruns caused by sloppy coding, check the sockaddr size and return EINVAL if it is too large. The setsockopt performs the join. The level argument to setsockopt is determined using the family of the group address and family_to_level function. Some systems support a mismatch between level and the socket's address family, for instance using IPPROTO_IP with MCAST_JOIN_GROUP, even with an AF_INET6 socket, but not all do, so it turns the address family into the appropriate level.

**The program below shows the second third of mcast_join, which handles IPv4 sockets.**

```
switch (grp->sa_family) {
case AF_INET:{
        struct ip_mreq mreq;
        struct ifreq ifreq;
```

```
memcpy(&mreq.imr_multiaddr,
       &((const struct sockaddr_in *) grp)->sin_addr,
       sizeof(struct in_addr));
if (ifindex > 0) {
    if (if_indextoname(ifindex, ifreq.ifr_name) == NULL) {
        errno = ENXIO; /* i/f index not found */
        return (-1);
    }
    goto doioctl;
} else if (ifname != NULL) {
    strncpy(ifreq.ifr_name, ifname, IFNAMSIZ);
doioctl:
    if (ioctl(sockfd, SIOCGIFADDR, &ifreq) < 0)
        return (-1);
    memcpy(&mreq.imr_interface,
           &((struct sockaddr_in *) &ifreq.ifr_addr)->sin_addr,
           sizeof(struct in_addr));
} else
    mreq.imr_interface.s_addr = htonl(INADDR_ANY);

return (setsockopt(sockfd, IPPROTO_IP, IP_ADD_MEMBERSHIP,
                   &mreq, sizeof(mreq)));
```

In the program, the IPv4 multicast address in the socket address structure is c opied i nto a n i p_mreq structure. If a n i ndex is specified, if_indextoname i s c alled, s toring t he na me into i freq structure. If t his succeeds then branch a head is to is sue th e io ctl. The c aller's n ame i s copied into an i freq structure, and an i octl of SIOCGIFADDR returns the unicast address associated with this name. Upon success, the IPv4 address is c opied i nto t he imr_interface member of t he ip_mreq s tructure. If a n index is not specified and a name is not specified, the interface is set to the wildcard address, telling the kernel to choose the interface. The setsockopt performs the join.

***The final portion of the function, which handles IPv6 sockets, is shown below.***

```
#ifdef IPV6
case AF_INET6:{
    struct ipv6_mreq mreq6;
    memcpy(&mreq6.ipv6mr_multiaddr,
           &((const struct sockaddr_in6 *) grp) ->sin6_addr,
```

```
        sizeof(struct in6_addr));
        if (ifindex > 0) {
        mreq6.ipv6mr_interface = ifindex;
        } else if (ifname != NULL) {
        if ( (mreq6.ipv6mr_interface = if_nametoindex(ifname)) == 0) {
        errno = ENXIO; /* i/f name not found */
        return (-1);
}
} else
        mreq6.ipv6mr_interface = 0;
        return          (setsockopt(sockfd,           IPPROTO_IPV6,
        IPV6_JOIN_GROUP,
        &mreq6, sizeof(mreq6)));
}
#endif
        default:
        errno = EAFNOSUPPORT;
        return (-1);
}
#endif
}
```

In the program, first the IPv6 multicast address is copied from the socket address structure into the ipv6_mreq structure. If an index was specified, it is stored in the ipv6mr_interface member; if a name was specified, the index is obtained by calling if_nametoindex; otherwise, the interface index is set to 0 for setsockopt, telling the kernel to choose the interface. The group is joined.

---

### Check Your Progress

1.  *Can you explain the use of mcast_join and mcast_leave.*

2.  *Can you explain the use of mcast_join_source_group*

---

## 13.9  dg_cli FUNCTION USING MULTICASTING

Modify dg_cli function by removing the call to *setsockopt.* Run a modified UDP echo server that joins the all-hosts group, and then run our program specifying the all hosts group as the destination address. We get a response from both the system on the subnet. They are each running the multicast echo server. Each reply is unicast because the source address of

the request which is used by each server as the destination address of the reply is a unicast address.

# 13.10  RECEIVING IP MULTICAST INFRASTRUCTURE SESSION ANNOUNCEMENTS

The IP multicast infrastructure is the portion of the Internet with inter-domain multicast enabled. Multicast is not enabled on the entire Internet.

In order to receive a multimedia conference on the IP multicast infrastructure, a site needs to know only the multicast address of the conference and the UDP ports for the conference's data streams. The *Session Announcement Protocol (SAP),* describes the way this is done (the packet headers and frequency with which these announcements are multicast to the IP multicast infrastructure) and the *Session Description Protocol (SDP),* describes the contents of these announcements (how the multicast addresses and UDP port numbers are specified). A site wishing to announce a session on the IP multicast infrastructure periodically sends a multicast packet containing a description of the session to a well-known multicast group and UDP port. Sites on the IP multicast infrastructure run a program named *sdr* to receive these announcements.

***The below shows the main program to receive SAP/SDP announcements***

```
#include "unp.h"

#define SAP_NAME "sap.mcast.net" /* default group name and port */

#define SAP_PORT  "9875"

void loop(int, socklen_t);

 int

main(int argc, char **argv)

 {

        int sockfd;

        const int on = 1;

        socklen_t salen;

        struct sockaddr *sa;

        if (argc == 1)

        sockfd =  Udp_client(SAP_NAME, SAP_PORT, ( void **)  &sa, &salen);

        else if (argc == 4)
```

```
    sockfd = Udp_client(argv[1], argv[2], (void **) &sa, &salen);

    else

    err_quit("usage: m   ysdr <m   cast-addr> <p   ort#> <i   nterface-
    name>");

    Setsockopt(sockfd, S  OL_SOCKET, S  O_REUSEADDR, &  on,
    sizeof(on));

    Bind (sockfd, sa, salen);

    Mcast_join(sockfd, sa, salen, (argc == 4) ? argv[3] : NULL, 0);

    loop (sockfd, salen); /* receive and print */

    exit (0);

}
```

In the program, the multicast address assigned for SAP announcements is 224.2.127.254 and its name is sap.mcast.net. All the well-known multicast appears in the DNS under the mcast.net hierarchy. The well-known UDP port is 9875. In the program udp_ client function is call to look up the name and port, and it fills in the appropriate socket address structure. In the program set the S O_REUSEADDR socket option to allow multiple instances of this program to run on a host, and bind the port to the socket. By binding the multicast address to the socket, prevent the socket from receiving any other UDP datagrams that may be received for the port. After that mcast_join function is call to join the group. If the interface name is specified as a command-line argument, it is passed to function; otherwise, the kernel chooses the interface on which the group is joined. Lastly call loop function to read and print all the announcements.

# 13.11  SENDING AND RECEIVING

The program that sends and receives multicast datagrams consists of two parts. The first part sends a multicast datagram to a specific group every five seconds and the datagram contains the sender's hostname and process ID. The second part is an infinite loop that joins the multicast group to which the first part is sending and prints every received datagram (containing the hostname and process ID of the sender). This allows us to start the program on multiple hosts on a LAN and easily see which host is receiving datagrams from which sender.

***The program below shows the main function of the program.***

#include "unp.h"

void recv_all(int, socklen_t);

```c
void send_all(int, SA *, socklen_t);

int

main (int argc, char **argv)

{

int sendfd, recvfd;

const int on = 1;

socklen_t salen;

struct sockaddr *sasend, *sarecv;

if (argc != 3)

err_quit("usage: sendrecv <IP-multicast-address> <port#>");

sendfd = Udp_client(argv[1], argv[2], (void **) &sasend, &salen);

recvfd = Socket(sasend->sa_family, SOCK_DGRAM, 0);

Setsockopt(recvfd, SOL_SOCKET, SO_REUSEADDR, & on, sizeof(on));

sarecv = Malloc(salen);

memcpy(sarecv, sasend, salen);

Bind (recvfd, sarecv, salen);

Mcast_join(recvfd, sasend, salen, NULL, 0);

Mcast_set_loop(sendfd, 0);

if (Fork () == 0)

recv_all (recvfd, salen); /* child -> receives */

send_all (sendfd, sasend, salen); /* parent -> sends */
```

In the program two sockets is created, one for sending and one for receiving. The receiving socket is to bind the multicast group and port. Then the receiving socket is to join the multicast group. The sending socket will send datagrams to this same multicast address and port. But if we try to use a single socket for sending and receiving, the source protocol address is 239.255.1.2:8888 from the bind (using netstat notation) and the destination protocol address for the sendto is also 239.255.1.2: 8888. However, now the source protocol address that is bound to the socket becomes the source IP address of the UDP datagram, and RFC 1122 forbids an IP datagram from having a source IP address that is a multicast

address or a broadcast address. The udp_client function creates the sending socket, processing the two command-line arguments that specify the multicast address and port number. This function also returns a socket address structure that is ready for calls to sendto along with the length of this socket address structure. Then create the receiving socket using the same address family that was used for the sending socket. Then set the SO_REUSEADDR socket option to allow multiple instances of this program to run at the same time on a host. Then allocate room for a socket address structure for this socket, copy its contents from the sending socket address structure, and bind the multicast address and port to the receiving socket. After that call mcast_join function to join the multicast group on the receiving socket and mcast_set_loop function to disable the loopback feature on the sending socket. For the join, specify the interface name as a null pointer and the interface index as 0, telling the kernel to choose the interface. Lastly the fork and then the child is the receive loop and the parent is the send loop.

***The program below shows the send a multicast datagram every five seconds.***

```
#include "unp.h"
#include <sys/utsname.h>
#define SENDRATE 5 /* send one datagram every five seconds */
void
send_all(int sendfd, SA *sadest, socklen_t salen)
{
        char line[MAXLINE]; /* hostname and process ID */
        struct utsname myname;
      if (uname(&myname) < 0)
      err_sys("uname error");;
      snprintf(line,   sizeof(line), " %s, %d  \n", myname.nodename,
      getpid());
      for ( ; ; ) {
      Sendto(sendfd, line, strlen(line), 0, sadest, salen);
      sleep(SENDRATE);
}
}
```

In the program send_all function, which sends one multicast datagram every five seconds. The main function passes as arguments the socket descriptor, a pointer to a socket address structure containing the multicast

destination and port, and the structure's length. In the program obtain the hostname from the uname function and build the output line containing it and t he pr ocess ID. S end da tagram, t hen go t o s leep. Then s end a datagram and then sleep for five seconds.

***The program below shows that receive all multicast datagrams for a group we have joined***.

```
#include "unp.h"
void
recv_all(int recvfd, socklen_t salen)
{
        int n;
        char line[MAXLINE + 1];
        socklen_t len;
        struct sockaddr *safrom;
        safrom = Malloc(salen);
        for ( ; ; ) {
        len = salen;
        n = Recvfrom(recvfd, line, MAXLINE, 0, safrom, &len);
        line[n] = 0; /* null terminate */
        printf("from %s: %s", Sock_ntop(safrom, len), line);
}
}
```

In t he program recv_all f unction, w hich i s t he infinite r eceive l oop. A socket ad dress s tructure i s al located t o r eceive t he s ender's p rotocol address f or ea ch c all t o r ecvfrom. E ach d atagram i s r ead b y r ecvfrom, null-terminated, and printed.

---

### Check Your Progress

1. *Write a program to send and receive the multicast datagram with two sockets.*

2. *What happens if we create one socket for both sending and receiving.*

---

## 13.12   SIMPLE NETWORK TIME PROTOCOL (SNTP)

Simple Network Time Protocol (SNTP) is a simplified version of Network Time Protocol (NTP) that is used to synchronize computer clocks

on a network. This simplified version of NTP is generally used when full implementation of NTP is not needed

SNTP i s a si mplified a ccess s trategy for s ervers a nd c lients u sing N TP. SNTP s ynchronizes a c omputer's s ystem time with a s erver th at h as already been synchronized by a source such as a radio, satellite receiver or modem.

SNTP supports unicast, multicast and anycast operating modes. In unicast mode, t he cl ient s ends a r equest t o a d edicated s erver b y r eferencing i ts unicast ad dress. O nce a r eply i s r eceived f rom t he s erver, t he cl ient determines the time, roundtrip delay and local clock offset in reference to the server. In multicast mode, the server sends an unsolicited message to a dedicated IPv4 or IPv6 l ocal br oadcast a ddress. G enerally, a m ulticast client d oes n ot s end an y requests t o t he s ervice b ecause o f t he s ervice disruption c aused b y u nknown a nd unt rusted m ulticast s ervers. T he disruption c an be a voided t hrough a n a ccess c ontrol m echanism t hat allows a client to select a designated server he or she knows and trusts.

NTP is a sophisticated protocol for synchronizing clocks across a WAN or a LAN, and c an o ften achieve m illisecond ac curacy. SNTP, a s implified but pr otocol-compatible version i ntended f or hosts t hat do not ne ed t he complexity o f a c omplete N TP imp lementation. It is c ommon for a few hosts on a LAN t o s ynchronize t heir c locks a cross t he Internet t o o ther NTP hos ts and t hen redistribute th is time o n th e LAN u sing e ither broadcasting or multicasting.

*The below program shows the NTP packet format and definitions*

#define JAN_1970        2208988800UL /* 1970 - 1900 in seconds */

struct l_fixedpt { /* 64-bit fixed-point */

uint32_t int_part;

uint32_t fraction;

};

struct s_fixedpt { /* 32-bit fixed-point */

uint16_t int_part;

uint16_t fraction;

};

struct ntpdata { /* NTP header */

u_char status;

u_char stratum;

u_char ppoll;

int precision:8;

```
        struct s_fixedpt distance;

        struct s_fixedpt dispersion;

        uint32_t refid;

        struct l_fixedpt reftime;

        struct l_fixedpt org;

        struct l_fixedpt rec;

        struct l_fixedpt xmt;

        };
```

#define VERSION_MASK 0x38

#define MODE_MASK 0x07

#define MODE_CLIENT 3

#define MODE_SERVER 4

#define MODE_BROADCAST 5

In the program, the l_fixedpt defines the 64-bit fixed-point values used by NTP for timestamps and s_fixedpt defines the 32-bit fixed-point values that are also used by NTP. The ntpdata structure is the 48-byte NTP packet format.

## 13.13   SUMMARY

A mu lticast a pplication s tarts b y jo ining th e mu lticast g roup assigned to the application. This tells the IP layer to join the group, which in turns tells the datalink layer to receive multicast frames that are sent to the corresponding hardware layer multicast address.

Multicasting on a WAN requires multicast-capable routers and a multicast routing protocol. Until all the routers on the Internet are multicast-capable, multicast is o nly a vailable to a subset of Internet us ers. The te rm "IP multicast infrastructure" is use to describe the set of all multicast-capable systems on the Internet.

Nine socket options provide the API for multicasting:

- Join an any-source multicast group on an interface

- Leave a multicast group

- Block a source from a joined group

- Unblock a blocked source

- Join a source-specific multicast group on an interface

- Leave a source-specific multicast group

- Set the default interface for outgoing multicasts

- Set the TTL or hop limit for outgoing multicasts

- Enable or disable loopback of multicasts

The first six are for receiving, and the last three are for sending.

## 13.14  TERMINAL QUESTIONS

1.  Explain multicast address.

2.  Write a short note on multicast socket option.

3.  Explain the use of mcast_block_source.and mcast_unblock_source

4.  Explain source specific multicast.

5.  Write about SNTP in terms of multicast.

6.  Write a program to show NTP packet format and definitions

# UNIT-14 : INTER PROCESS COMMUNICATION

## Structure

## 14.1 INTRODUCTION

In this unit, the different methods of IPC will be discussed. In this unit, we will learn about; File and record locking, Pipes, FIFOs streams and messages, name spaces, system IPC, message queues and semaphores.

A process is an active operating system entity which executes programs. Normally, a process, like a specialist, does one particular job (well). In real life, there are complex workflows and we, often have multiple processes collaborating to accomplish certain objectives. In order to work together, processes need to exchange data. So, we have various inter process communication (IPC) mechanisms.

The figure 14.1 shows the IPC between two processes on a single system. The information between the two processes going through the kernel. The figure 14.2 shows the IPC between two processes on different system.
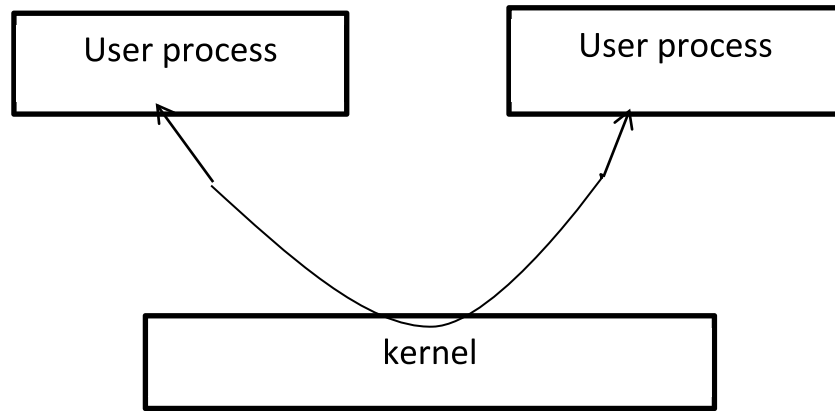
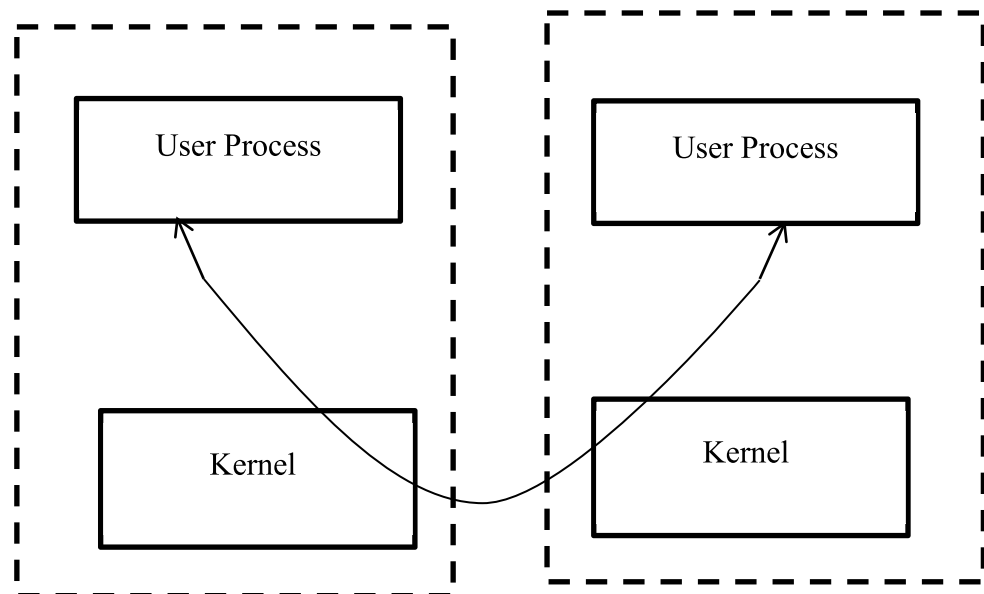Figure 14.1: IPC between two processes on a single system



Figure 14.2: IPC between two processes on different system

## 14.2 OBJECTIVE

At the end of this unit, you should be able to know about: -

- Purpose of File and record locking

- How to use Pipes and FIFOs.

- Streams and messages, Message queues

- Name spaces, system IPC

- What is the use of Semaphores?

## 14.3  FILE AND RECORD LOCKING

When multiple process wants to share resource, it is essential that some form of mutual exclusion be provide so that only one process at a time can access the resource. The example is line printer daemon. The process that places a job on the print queue (to be printed at a later time by another process) has to assign a unique sequence number to each print job. Each process that needs to assign a sequence number goes through three steps:

- It reads the sequence number file

- It uses the number

- It increments the number and writes it back

The problem is that in the time it takes a single process can perform the same three steps; another process can perform the same three steps. The need is for a process to be locked so no other process can access the same file until the first process is done.

In file locking locks an entire file, while record locking allows a process to lock a specified portion of a file. Used to ensure that a process has exclusive access to a file before using it

#include int lockf (int fd, int function, long size);

fd---file descripter (not a file pointer)

size--- define the record size or lock area: [ offset, offset + size]. size=0 means the rest

of the file. Use lseek() to move the current offset. When the offset position is

set to the beginning and size=0 then lock the whole file.

## Function:

F_ULOCK---unlock a previous lock

F_LOCK ---lock a region(blocking)

F_TLOCK ---Test and lock a region(nonblocking)

F_TEST ---Test a region to see if it is locked.

## Example:

Use F_TLOCK instead of F_TEST and F_LOCK.

If (lockf(fd, F_TEST, size)==0) /* If the region is locked, -1 is returned and the

process is in sleep state*/

Re=lockf(fd, F_LOCK, size); /*a small chance that another process locks between

TEST and LOCK*/

rc=lockf(fd, F _TLOCK, s ize) / * T est + l ock d one as an atomic operation, If

unsuccessful, l ockf() returns −1 and the calling

process c ontinues t o do ot her things*/

*The following are the two types of Linux file locking:*

1.    Advisory locking

2.    Mandatory locking

## 1.    Advisory Locking

Advisory l ocking r equires c ooperation f rom t he pa rticipating processes. S uppose pr ocess " A" a cquires a WRITE l ock, a nd i t started w riting in to th e file, a nd p rocess " B", without tr ying to acquire a l ock, i t can o pen the file and w rite i nto i t. H ere p rocess "B" is the non-cooperating process. If process "B", tries to acquire a l ock, t hen i t m eans t his p rocess i s co-operating t o ensure t he "serialization". Advisory locking w ill w ork, onl y i f t he participating processes ar e cooperative. A dvisory l ocking sometimes also called as "unenforced" locking.

Posix r ecord l ocking i s c alled a dvisory l ocking. T his m eans t he kernel m aintains correct know ledge of all f iles t hat ha ve be en locked b y each pr ocess, but i t doe s not pr event a pr ocess f rom writing to a file th at is read-locked b y a nother pr ocess. S imilarly, the kernel does not prevent a process from reading from a file that is w rite-locked b y a nother pr ocess. A p rocess c an i gnore a n advisory lock and write to a file that is read-locked, or read from a file th at is w rite-locked, a ssuming t he p rocess ha s a dequate permissions t o r ead or w rite t he f ile. A dvisory locks a re f ine f or cooperating processes.

## 2.    Mandatory Locking

Mandatory l ocking d oesn't r equire c ooperation f rom t he participating pr ocesses. Mandatory l ocking causes t he k ernel t o check every open, read, and write to verify that the calling process isn't violating a lock on the given file.

## 14.4  PIPES

A pi pe pr ovides a on e-way f low o f da ta.  Two p rocesses can  be joined b y t he pi pe s ymbol ( |) on t he s hell c ommand l ine. T he s tandard output of  t he f irst pr ocess be comes t he s tandard i nput f or t he s econd process. For example,

$ ls -ls | more

Example a pipe provides a one-way flow of data.

int pipe (int * filedes);

int pipefd[2]; /* pipefd[0] is opened for reading; pipefd[1] is opened for writing */

***The program below shows how to create and use a pipe:***

```
main ()
{
int pipefd[2], n;
 char buff[100];
if (pipe(pipefd) < 0 ) err_sys("pipe error");
printf("read fd = %d, write fd = %d\n", pipefd[0], pipefd[1]);
if (write(pipefd[1], "hello world\n", 12) != 12) err_sys("write
error");
if ((n=read(pipefd[0], buf f, s izeof(buff)))  < = 0)  err_sys("read
error");
write (1, buff, n); /*fd=1=stdout*/
}
```

**Result:**　　　hello world
　　　　　　　read fd=3, write df =4

## Properties of Pipe:

* Pipes do not have a name. For this reason, the processes must share a pa rent pr ocess. T his i s t he m ain dr awback t o pi pes. H owever, pipes are treated as file descriptors, so the pipes remain open even after fork and exec.

* Pipes do not  di stinguish between m essages; t hey just read a fixed number of bytes. Newline (\n) can be used to separate messages. A structure w ith a  le ngth field c an  be us ed f or m essage c ontaining binary data.

* Pipes c an a lso be  us ed t o g et t he out put of  a  c ommand or  t o provide input to a command
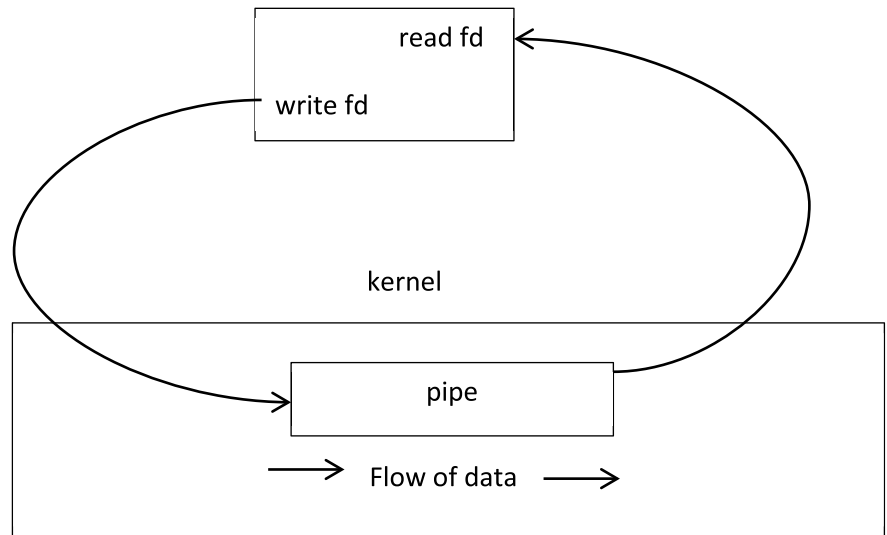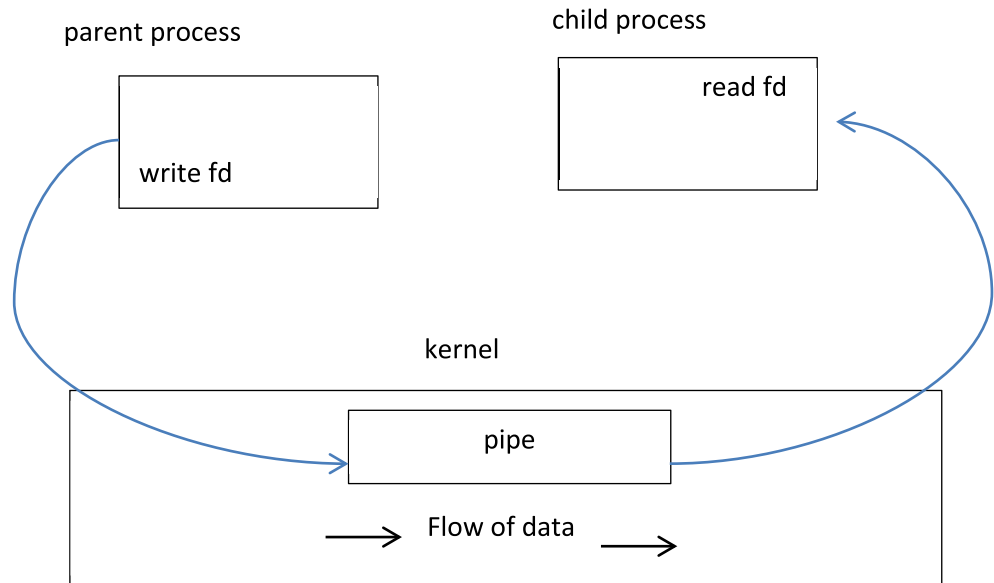
Figure 14.3: Pipe in a single process



Figure 14.4: Pipe between two processes

The figure 14.3 shows the pipe in a single process and figure 14.4 shows the pipe between two processes.

One major feature of pipe is that the data flowing through the communication medium is transient, that is, data once read from the read descriptor cannot be read again. Also, if we write data continuously into the write descriptor, then we will be able to read the data only in the order in which the data was written. One can experiment with that by doing successive writes or reads to the respective descriptors.

## 14.5   FIFO

A F IFO is s imilar to a p ipe. A F IFO (First in First Out) is a one-way flow of data. FIFOs have a name, so unrelated processes can share the FIFO. F IFO i s a na med pi pe. T his is t he main difference be tween pi pes and F IFOs. Another m ajor d ifference b etween FIFOs and pi pes i s t hat FIFOs last th roughout th e lif e-cycle o f th e s ystem, w hile p ipes la st o nly during the life-cycle of the process in which they were created. To make it more clearly, FIFOs exist be yond t he l ife of t he process. S ince t hey are identified b y the file s ystem, th ey remain in the h ierarchy u ntil e xplicitly removed us ing unl ink, b ut pi pes a re i nherited on ly b y r elated pr ocesses, that is, processes which are descendants of a single process.

**Create:** A FIFO is created by the mkfifo function:

#include <sys/types.h>

#include <sys/stat.h>

int mkfifo(const char *pathname, mode_t mode);

pathname – a UNIX pathname (path and filename).

mode – the file permission bits.

FIFO c an al so b e created by the mknod s ystem call, e.g., m knod("fifo1", S_IFIFO|0666, 0) is same as mkfifo("fifo1", 0666).

**Open:** mkfifo tries to create a new FIFO. If the FIFO already exists, then an EEXIST error is returned. To open an existing FIFO, use open (), fopen () or freopen ()

**Close:** to close an open FIFO, use close (). To delete a cr eated F IFO, use unlink ().

The table 14.1 shows the Effect of O_NDELAY flag on pipes and FIFOs. A pipe or FIFO follows these rules for reading and writing:

- A read requesting less data than is in the pipe or FIFO returns only the requested amount of data.

- If a process asks to read more data than is currently available in the pipe F IFO, O nly the da ta a vailable i s r eturned. The pr ocess m ust be prepared to handle a return value from read that is less than the requested amount.

- If there is no da ta in the pipe or FIFO, and if no pr ocesses have it open for w riting, a r ead return zero, s ignifying t he e nd of file. If the r eader h as s pecified O _NDELAY, i t can not t ell i f a r eturn value of zero means there is no data currently available or if there are no writers left.

| Condition | Normal | O_NDELAY set |
|---|---|---|
| Open FIFO, read-only with no process having the FIFO open for writing | Wait until a process opens the FIFO for writing | Return immediately, no error |
| Open FIFO, write-only with no process having the FIFO open for reading | Wait until a process opens the FIFO for reading | Return an error immediately, errno set to ENXIO |
| read pipe or FIFO, no data | Wait until there is data in the pipe or FIFO, or until no processes have it open for writing; return a value of zero if no processes have it open for writing, otherwise return the count of data | Return immediately, return value of zero |
| Write, pipe or FIFO is full | Wait until space is available, then write data | return immediately, return value of zero |

Table 14.1: Effect of O_NDELAY flag on pipes and FIFOs.

- If a process writes less than the capacity of a pipe (which is at least 4096 bytes) the write is guaranteed to be atomic. This means that if two processes each write to a pipe or FIFO at about the same time, either all the data from the first process is written, followed by all the data from the second process, or vice versa. The system does not mix the data from the two processes-i.e., part of the data from one process, followed by part of the data from the other process. If, however, the write specifies more data than the pipe can hold, there is no guarantee that the write operation is atomic.

- If a process writes to a pipe or FIFO, but there are no processes in existence that have it open for reading, the SIGPIPE signal is

generated, and the write returns zero with errno set to EPIPE. If the process h as not called s ignal t o ha ndle t he S IGPIPE not ification, the d efault a ction is to te rminate th e S IGPIPE s ignal, o r if it handles the signal and returns from its signal handler.

# 14.6  STREAMS AND MESSAGES

A S TREAM is a general, flexible p rogramming model f or U NIX system communication services. STREAMS define standard interfaces for character input/output (I/O) within the kernel, and between the kernel and the rest of the UNIX system. The mechanism consists of a s et of s ystem calls, kernel resources, and kernel routines.

A S TREAM e nables you t o c reate m odules t o pr ovide s tandard da ta communications s ervices a nd t hen m anipulate t he m odules on a stream. From t he a pplication l evel, m odules c an be dy namically s elected an d interconnected. N o ke rnel pr ogramming, compiling, and l ink e diting a re required to create the interconnection.

A S TREAM pr ovides a n e ffective e nvironment f or ke rnel s ervices and drivers r equiring m odularity. S TREAMS parallel the la yering mo del found in networking protocols. For example, STREAMS are suitable for:

- Implementing network protocols

- Developing character device drivers

- Developing network controllers (for example, for an Ethernet card)

- I/O terminal services

In S TREAMS, a ll in formation i s e xchanged v ia m essages i.e., bot h da ta and c ontrol m essages of va rious p riorities. A m ulti-component m essage structure is used to reduce the overhead of

1. Memory-to-memory copying i.e., via reference counting

2. Encapsulation/de-encapsulation i.e., via composite messages.

Messages m ay be que ued a t S TREAM m odules. Many Unix pr ocesses that n eed t o i mpose a m essage s tructure o n t op o f a s tream b ased IPC facility. More s tructured m essage c an a lso be bui lt, a nd t his i s w hat the Unix message queue form of IPC does. We can also add more structure to either a pipe or FIFO. We define a message in mesg.h header file as

    /*
    *Definition of "our" message.

```
* You may have to change the 4096 to a smaller value, if message
*queues on   your s ystem w ere  configured w ith "m sgmax" l ess
*than 4096.
 */
# define MAXMESGDATA       (4096-16)
                                                    /* w  e don' t   want
sizeof(Mesg) > 4096 */
#define MESGHDRSIZE      (sizeof(Mesg) – MAXMESGDATA)
                                                    /* l      ength of
mesg_len and mesg_type*/
typedef struct {
int mesg_len;  /*#bytesin mesg_data, can be 0 or > 0 */
long mesg_type; /* message type, must be > 0 */
char mesg_data [MAXMESGDATA];
} Mesg;
```

***Check Your Progress:***

1. *Can you write a program that create FIFO in which it writes first then read?*

2. *How stream and message are useful in Unix?*

## 14.7  NAME SPACES

The set of possible names for a given type of IPC is called its name space. The n ame  space i s i mportant be cause a ll forms of  IPC ot her t han plain pipes, the name is how the client and server connected to exchange message. The table 14.2 shows the list of available name space below.

| IPC type | Name Space | Identification |
|---|---|---|
| pipe | No name | File descriptor |
| fifo | Path name | File descriptor |
| message queue | Key_t key | identifier |
| shared memory | Key_t key | identifier |
| semaphore | Key_t key | identifier |
| socket-unix domain | Path name | File descriptor |
| socket-other domains | Domain depndent | File descriptor |

Table 14.2 : List of available name space

# 14.8   SYSTEM IPC

The three types of IPC

- Message queues

- Semaphores

- Shared memory

These are collectively referred as "system V IPC"

Linux s upports t hree t ypes of i nterprocess communication m echanisms that first appeared in Unix System V (1983). These are message queues, semaphores an d s hared memory. T hese S ystem V IPC m echanisms al l share co mmon au thentication m ethods. P rocesses m ay ac cess t hese resources only by passing a unique reference identifier to the kernel via system c alls. A ccess t o t hese S ystem V IPC o bjects i s ch ecked u sing access p ermissions; much like ac cesses to files are checked. The access rights to the S ystem V IPC object is set by the creator of the object via system calls. The object's reference identifier is used by each mechanism as an index into a table of resources. It is not a straight forward index but requires s ome m anipulation t o g enerate t he index. All L inux d ata structures representing S ystem V I PC objects in th e s ystem in clude an ipc_perm s tructure which c ontains t he ow ner and c reator pr ocess's us er and group identifiers. The access mode for this object (owner, group and other) and the IPC object's key. The key is used as a way of locating the System V IPC o bject's r eference i dentifier. Two s ets o f k eys are supported: public and private. If the key is public then any process in the system, subject to rights checking, can find the reference identifier for the System V IPC object. System V IPC objects can never be referenced with a key, only by their reference identifier.

A summary of their system calls is shown in table 14.3.

|  | Message queue | Semaphore | Shared memory |
|---|---|---|---|
| Include file | <sys/msg.h> | <sys/sem.h> | <sys/shm.h> |
| System c all to c reate or open | msgget | semget | shmget |
| System c all f or control operations | msgctl | semctl | shmctl |
| System c alls f or IPC operations | msgsnd | semop | shmat |
|  | msgrcv |  | shmdt |

Table 14.3: Summary of system V IPC system calls

The value returned by msgget is the message queue identifier, msqid, or -1 if an error occurred.

## 14.9  MESSAGE QUEUES

Message q ueues al low o ne o r m ore p rocesses t o w rite m essages, which w ill be read b y o ne or m ore reading pr ocesses. Linux m aintains a list o f m essage q ueues, in the m sgque vector; each el ement o f w hich points t o an msqid_ds da ta s tructure t hat f ully de scribes t he m essage queue. When message queues are created a new msqid_ds data structure is allocated f rom s ystem memory and in serted in to th e v ector. For e very message query in the system, the kernel maintains the following structure of information

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

#include<sys/ipc.h>

struct msqid_ds{
        struct ipc_perm msg_perm;

        struct msg     *msg_first;

        struct msg     *msg_last;

        ushort    msg_cbytes;

        ushort    msg_qnum;

        ushort    msg_qbytes;

        ushort    msg_lspid;

        ushort     msg_lrpid;

        time_t     msg_stime;

        time_t    msg_rtime;

        time_t    msg_ctime;
    };
```

A new message query is created or an existing message queue is accessed with the msgget system call

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

#include    <sys/ipc.h>

#include    <sys/msg.h>

int msgget (key_t key,int msgflag);
```

The msgflag value is a combination of constants shown in table 14.4.

| Numeric | Symbolic | Description |
|---|---|---|
| 0400 | MSG_R | Read by owner |
| 0200 | MSG_W | Write by owner |
| 0040 | MSG_R >> 3 | Read by group |
| 0020 | MSG_W >>3 | Write by group |
| 0004 | MSG_R >>6 | Read by world |
| 0002 | MSG_W>>6 | Write by world |
|  | IPC_CREAT | Create new entry |
|  | IPC_EXCL | Create new entry |

Table 14.4: msgflag values for msgget system call.

Each m sqid_ds data s tructure contains a n i pc_perm da ta s tructure a nd pointers to the messages entered onto this queue. In addition, Linux keeps queue modification times such as the last time that this queue was written to a nd s o on. T he m sqid_ds a lso contain two w ait q ueues; one f or t he writers to the queue and one for the readers of the message queue. Each time a process attempts to write a message to the write queue its effective user a nd group i dentifiers a re c ompared with t he m ode i n t his que ue's ipc_perm d ata s tructure. If t he pr ocess c an w rite t o t he que ue t hen t he message may be copied from the process's address space into an msg data structure and put at the end of this message queue. Each message is tagged with an ap plication s pecific t ype, a greed b etween t he co operating processes.

However, t here m ay be no r oom f or t he m essage a s Linux r estricts t he number and length of messages that can be written. In this case the process will be added to this message queue's write wait queue and the scheduler will be called to select a new process to run. It will be woken up when one or more messages have been read from this message queue. Reading from the q ueue i s a s imilar p rocess. A gain, t he p rocesses acc ess r ights t o t he write queue are checked. A reading process may choose to either get the first m essage i n t he q ueue r egardless o f i ts t ype o r s elect m essages w ith particular t ypes. If n o m essages m atch these cr iteria the reading p rocess will be added to the message queue's read wait queue and the scheduler run. W hen a ne w m essage i s w ritten t o t he que ue t his pr ocess w ill be woken up and run again.

## 14.10  SEMAPHORES

Semaphores a re s ynchronization pr imitive. T he m ain us e of semaphores is to synchronize the access to shared memory segments. In its simplest form, a s emaphore i s a l ocation i n m emory whose v alue can b e tested an d s et b y m ore t han o ne p rocess. S emaphores can be us ed t o

implement critical regions, areas of critical code that only one process at a time should be executing.

***The following information is related to semaphore:***

1. The semaphore is stored in the kernel:

   a. Allows atomic operations on the semaphore.

   b. Processes are prevented from indirectly modifying the value.

2. A process acquires the semaphore if it has a value of zero. The value of the semaphore is then incremented to 1. When a process releases the semaphore, the value of the semaphore is decremented.

3. If the semaphore has non-zero value when a process tries to acquire it, that process blocks.

4. In 2 and 3, the semaphore acts as a customer counter. In most cases, it is a resource counter.

5. When a process waits for a semaphore, the kernel puts the process "to sleep" until the semaphore is available. This is better (more efficient) than busy waiting such as TEST & SET.

6. The kernel maintains information on each semaphore internally, using a data structure struct semis_ds that keeps track of permission, number of semaphores, etc.

7. Apparently, a semaphore in Unix is not a single binary value, but a set of nonnegative integer values.

8. There are 3 (logical) types of semaphores:

   - Binary semaphore – have a value of 0 or 1. Similar to a mutex lock. 0 means locked; 1 means unlocked.

   - Counting semaphore – has a value $\geq 0$. Used for counting resources, like the producer-consumer example. Note that value =0 is similar to a lock (resource not available).

   - Set of counting semaphores – one or more semaphores, each of which is a counting semaphore.

9. There are 2 basic operations performed with semaphores:

   - Wait – waits until the semaphore is > 0, then decrements it.

   - Post – increments the semaphore, which wakes waiting processes.

**Operations on a semaphore are performed using:**

int semop(int *semid*, struct sembuf *\*opsptr*, unsigned int *nops*)

*semid* — value returned by semget.

*nops* — # of operations to perform, or the number of elements in the *opsptr* array.

*opsptr* — points to an array of one or more operations. Each operation is defined as:

struct sembuf { ushort sem_num; /* semaphore #, numbered from 0, 1, 2 … */

short sem_op; /* semaphore operation */

short sem_flg; / *operations flags, such as 0, IPC_NOWAIT for nonblocking call,

or SEM_UNDO to have the semaphore automatically released when the process is terminated prematurely. */

};

sem_op = 0 – wait until the semaphore is 0. IPC_NOWAIT causes an error if semval≠0.

sem_op > 0 – increment the semaphore value: semval + sem_op, (acquire)

sem_op < 0 – wait until the semaphore value≥|sem_op| and decrement the semaphore value: semval - |sem_op|, (release)

**Example: How to write lock/unlock (somewhat like P/V operations)**

#include <sys/types.h>

#include <sys/ipc.h>

#include <sys/sem.h>

#define SEMKEY 123456L /* key value for semget() */

#define PERMS 0666

static struct sembuf op_lock[2]= { 0, 0, 0, / * wait for sem #0 to become 0 */

0, 1, SEM_UNDO /* then increment sem #0 by 1 */ };

static struct sembuf op_unlock[1]= { 0, -1, ( IPC_NOWAIT | SEM_UNDO)

/* decrement sem #0 by 1 (sets it to 0) */ };

int semid = -1; /* semaphore id. Only the first time will create a semaphore. */

my_lock( )

{ if (semid <0) {

if ( ( semid=semget(SEMKEY, 1, IPC_CREAT | PERMS )) < 0 ) printf("semget error"); }

```
        if (semop(semid, &op_lock[0], 2) < 0) printf("semop lock error");
}
my_unlock( )
{
        if ( semop(semid, & op_unlock[0], 1)  <  0)  pr intf("semop unl ock
error");
}
```

The s emaphore has us ed for s ynchronization.  T he bi nary s emaphore has
created, a s ingle s emaphore va lue that is e ither z ero or one.  For l ocking
the semaphore call *semop( )* to do operations automatically.  First, wait for
the semaphore value to become zero, and then increment the value to one.
This i s  an ex ample w here m ultiple  semaphore ope ration  must be  don e
atomically by the kernel. If it to ok two system calls to do this, one to test
the va lue a nd w ait f or i t t o be come z ero, a nd a nother t o i ncrement t he
value, the operations would not work. For unlocking the resource, s*emop(*
*)* will cal l t o d ecrement the s emaphore v alue. S ince w e h ave l ock o n the
resource, we know that the semaphore value is one before the call, so the
call cannot wait.

## 14.11  SUMMARY

        IPC ha s t raditionally be en a  m assive a rea i n U NIX.   In t his  unit,
we have covered record locking and file locking since the sharing of single
file b etween m ultiple p rocesses i s a   common o ccurrence.  V arious  IPC
techniques l ike P IPES, F IFO, M essage que ues, s emaphores a nd s hared
memory are covered.

## 14.12  TERMINAL QUESTIONS

1.      What i s a s ignal ge nerated for t he w riter of a pi pe of F IFO w hen
        the ot her e nd di sappears, a nd f or t he r eader of  a P IPE or  FIFO
        when its writer disappears?

2.      What ha ppens w ith t he c lient s erver e xample us ing m essage
        queues if the file to be copied is a binary file?

3.      What happens to the version that uses the popen function if the file
        is a binary file?

4.      What is the use of semaphore?

5.      Can you design a message in mesg.h header file?

6.      List few of the available name spaces.

7.      Write a program to lock and unlock a semaphore.

# UNIT-15 : REMOTE LOGIN

## Structure

## 15.1 INTRODUCTION

In this unit, we will learn details about rlogin (remote login). We will study Terminal line disciplines. Then we get knowledge about Pseudo- Terminal. Then we see the terminal modes and control terminal. Lastly, we will discuss the transparent issues in RPC.

rlogin (remote login) is a UNIX command that allows an authorized user to login to other UNIX machines (hosts) on a network and to interact as if the user were physically at the host computer. Once logged in to the host, the user can do anything that the host has given permission for, such as read, edit, or delete files.

## 15.2   OBJECTIVES

In this unit, we will understand the following:

- Terminal line disciplines

- Pseudo- Terminal and Terminal modes

- Control Terminal

- rlogin overview

- RPC Transparency Issues

## 15.3  TERMINAL LINE DISCIPLINES

Terminal drivers are complicated by the line discipline associated with their Terminal. It is assumed to be a full duplex device so that the input path and output path are separate. The line discipline is within the kernel, somewhere between the actual device driver and the user process. The terminal line discipline is just a module that is pushed onto a stream on top of the actual terminal device driver. Figure 15.1 shows a normal interactive shell showing terminal line discipline.

There are several functions that can be done by a line discipline mode.

- Echo the characters entered into lines

- Assemble the characters entered into lines, so that a process reading from the terminal receives complete lines.

- Edit the lines that are input. UNIX allows you to erase the preceding character and also to kill the entire line being input and start over with a new line.

- Generate signals when certain terminal keys are entered. The SIGINT and SIGQUIT signals can be generated this way, for example.

- Process flow control characters. For Example, when you press the control –S key, the output to the terminal is stopped. The restart the output, the Control-Q key is entered.

- Allow you to enter an end –of-file character.
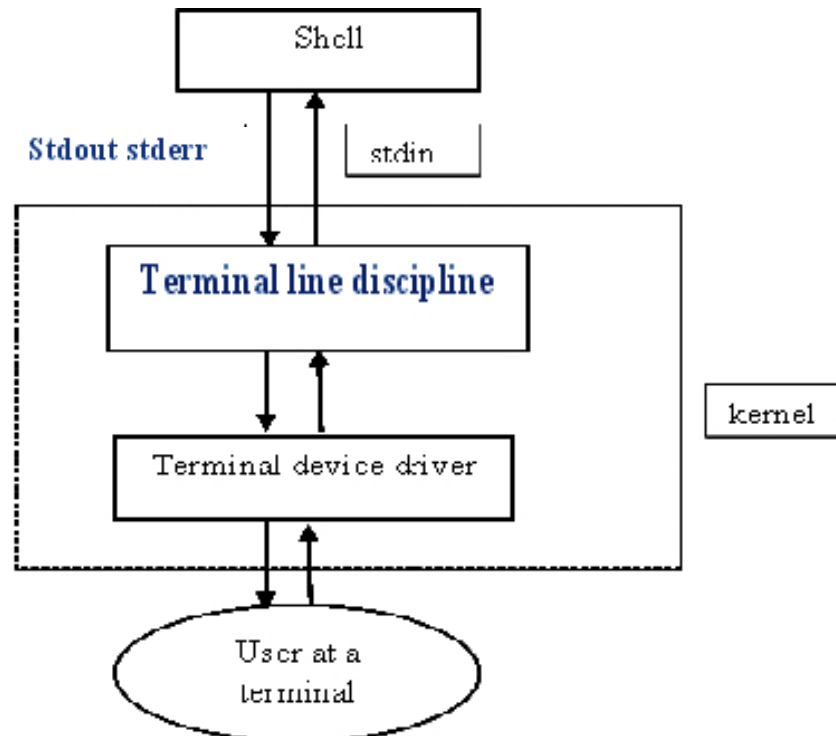
- Do character conversions

Figure 15.1: Normal interactive shell showing terminal line discipline

There are many versions of the line discipline modules. For example, BSD supplies five modules.

- The "old d iscipline" th at is s imilar to th e v ersions 7 U NIX terminal handler.

- The new discipline is a superset of the old discipline.

- It provides the features needed for job control along with enhanced editing capabilities.

- The Berknet line discipline.

- The s erial Line Internet P rotocol c an be us ed t o t ransfer IP datagrams across serial lines.

---

### Check Your Progress

1.  *Can you explain the different function in terminal line discipline mode?*

---

## 15.4    PSEUDO- TERMINAL

A pseudo-terminal is pair of devices. One half is called the master and the other half is called the slave. A process opens a p air o f p seudo-terminal de vices and gets t wo f ile de scriptors. T he s lave po rtion of pseudo-terminal d evices gets t wo file de scriptors. The s lave po rtion of a pseudo-terminal presents an interface to the user process that looks like a terminal device.

A ps eudo-terminal is mainly u sed to ma ke a p rocess b elieve th at it interacts w ith a t erminal a lthough it a ctually in teracts w ith o ne o r mo re processes. The f igure 15.2 s hows t he Pseudo-terminals as t hey ar e u sed by script.
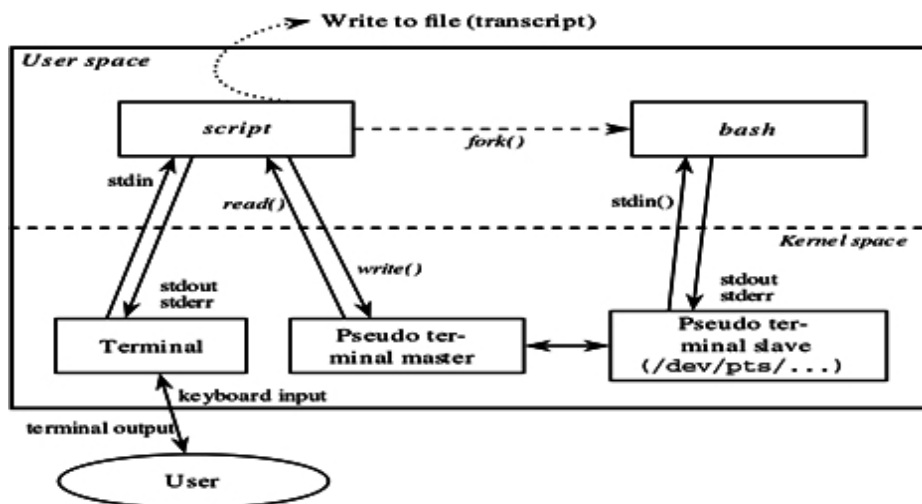


Figure 15.2 : Pseudo-terminals

## 15.5  TERMINAL MODES

In terminal models we are considering only standard terminal line discipline modules such as old line discipline and the new line discipline modules supported by 4.3BSD.

4.3BSD considers a terminal device in one of three modes.

- **Cooked mode** provides all the processing steps. The input is collected into lines and all special character processing is done. This is normal mode for interactive use.

- **Raw mode** lets the process receive every charter as in is input, with no interpretation done by the system. Raw mode is used for example by full screen editors such as vi and also by programs that use a serial line something other than interactive use.

- **Cbreak mode** is somewhere between cooked mode and raw mode. The cbreak mode provides character at a time input to the process reading from the terminal, instead of collecting the input into lines. The signal generating keys are still processed; however the editing features are disabled.

## 15.6  CONTROL TERMINAL

In 4.3BSD we have the child process from the fork dissociate from its control terminal before it opens the pseudo-terminal slave device. When the slave is opened it becomes the control terminal. Since we only want the new shell process that the child process execs to disassociate from its control terminal- we do not want the recording process that is reading from your actual terminal to do this- we must do this in child process and not in the parent. This is precisely why the opening of a pseudo-terminal pair into pieces. We do not want to open the slave device until we are in the child process.

## 15.7  RLOGIN OVERVIEW

The terminal line discipline on the local system is placed into the raw mode with echoing disabled by the rlogin client process, so that all keystrokes are passed to the remote system. The raw mode is required to run programs such as the vi editor on the remote system. In the normal UNIX fashion characters that are entered on the local are echoed by the remote system. If the remote system is in a cooked mode then the echoing is done by the terminal line disciplines on the remote system. If the remote system is in a raw mode then the echoing is done by that remote process itself. The figure 15.3 shows the 4.3BSD rlogin processes.
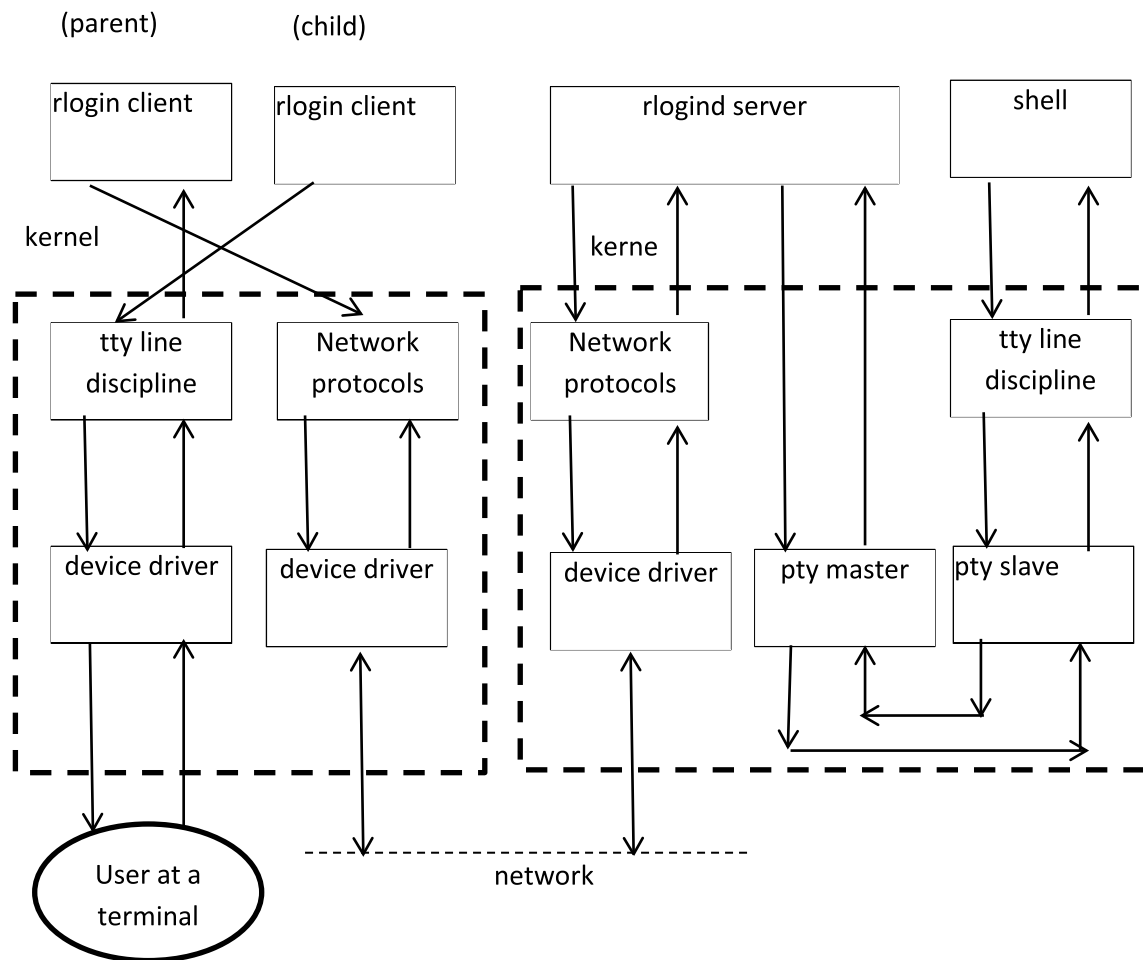
Figure 15.3 : 4.3BSD rlogin processes

The r login f acility p rovides a r emote-echoed, l ocally flow-controlled virtual terminal with proper flushing of output. It is widely used between U NIX hos ts be cause i t pr ovides t ransport of m ore of t he U NIX terminal e nvironment s emantics t han do es t he T elnet pr otocol, and because on m any U NIX hos ts i t c an be c onfigured not t o r equire us er entry of passwords when connections originate from trusted hosts.

Apart from this, rlogin suffers most of the same security disadvantages as telnet, s uch a s t he fact t hat a ll c ommunication, i ncluding pa sswords, i s transmitted in c lear-text. T he t rusted hos ts feature b ypasses pa ssword authentication w hen a n rlogin/rhosts-file i s s pecified. T his p oses a great security risk as the files themselves are not very well secured, and in many cases, can be found on t he host's NFS s hare. Because of these p roblems, rlogin i s not i n m uch u se t oday a nd ha s m ostly been r eplaced b y t he superior SSH protocol.

## 15.8  RPC TRANSPARENCY ISSUES

The system needs to provide a transparent interface for the client, so that there is no distinction between making a remote procedure call and making a local function call.

The client and server stubs hide the network code, but there are other issues that need to be addressed:

**Parameter Passing** – can't pass parameters by reference, since the subroutine and the calling program don't share the same address space. Sun RPC allows only a single argument and a single result. A structure is required for multiple values.

*Binding* – the client needs some way to determine which host is a server. Choices are to require that the client knows which host to contact, or uses super server that keeps track of the addresses of each server, or use a centralized database where each host indicates which servers it is willing to run. Sun RPC takes the following approach. The port mapper on the remote host is contacted for the port number. The port mapper also accepts the broadcast requests. If a matching server is found, the request is passed to the server. The port number is then returned with the results, so the client can be connected directly to the server on future calls.

*Transport Protocol* – Sun RPC supports TCP and UDP. TCP is a byte-stream protocol, so there are no message delimiters. To solve this, a 32-bit integer giving the number of bytes is placed at the beginning of each record. With UDP on older systems, at most 8192 bytes can be sent for the arguments or results of one call. The maximum can never exceed the size of a UDP datagram, which is 64 K – headers.

**Exception Handling** – not only could the typical errors, such as segmentation fault, occur in the remote procedure, but also network problems are also possible. A timeout is usually used to detect server crashes.

The client might also wish to terminate the server. With Sun RPC, the client cannot send an interrupt to the server. Both UDP and TCP handle

timeouts a nd r etransmissions. UDP w ill te rminate a fter s ome n umber o f unsuccessful attempts.

**Call Semantics** – because of ne twork pr oblems, t he r equest t o s tart a remote pr ocedure m ight be s ent m ultiple t imes. P rocedures t hat c an be executed mu ltiple time s w ithout a p roblem a re c alled i dempotent. Examples i nclude c omputing a s quare r oot or c hecking an a ccount balance.

There are three different forms of RPC semantics

1.  Exactly o nce-Means t hat t he r emote p rocedure w as ex ecuted o ne time pe riod. This t ype o f ope ration i s ha rd t o achieve, ow ning t o the possibility of server crashes.

2.  At m ost onc e –Means t hat t he r emote p rocedure w as ei ther n ot executed at al l o r i t w as ex ecuted o ne t ime at most. I f a n ormal return i s m ade t o t he ca ller, w e k now t he r emote p rocedure w as executed o ne t ime. But i f an e rror return i s made, w e are not certain if remote procedure was executed once or not at all.

3.  At l east o nce-Means t hat t he r emote p rocedure w as ex ecuted at least one time, but pe rhaps m ore. This i s t ypically for i dempotent procedures-the client keeps transmitting its request until it receives a valid response. But if the client has to send its request more than once t o r eceive a v alid r esponse, t here i s a p ossibility t hat t he remote procedure was executed more than once.

**Data Representation** – Need a s tandard representation, so the client and server can execute on different architectures. Sun RPC uses the XDR data representation standard.

**Performance** – there c an b e co nsiderable o verhead f or calling an R PC. For e xample, t he ove rhead m ight be 100 t imes the ove rhead of a l ocal procedure cal l. S un R PC u ses s everal m echanisms, s uch as passing pointers, t o m inimize c oping d ata. T he pu rpose of R PC i s t o s implify network programming, not just to replace LPC with RPC.

**Security** – May need to restrict who can execute a program on the server. In t he l ocal cas e t he cal ler o f a f unction can b e s ure t hat i t i s cal ling an authorised pr ovider of t he s ervice, a nd t he pr ocedure can b e s ure i t i s called b y an au thorised user, because they are l inked t ogether at co mpile time. With a remote call, neither party can be sure.

To assure clients and server that they are talking to authorised servers and clients, Sun R PC includes an authentication mechanism. The client sends its c redentials a nd a v erifier to th e s erver with its R PC c all, th e s erver

returns its own verifier with the results. The standard authentication methods provided by the library are Null, UNIX, Short and DES, but it is easy to add new methods. Using the authentication mechanisms is not transparent; it requires some extra programming on the client and server sides.

## 15.9 SUMMARY

Remote Login is comparatively complicated networking example, which we have discussed. The most complicated part of remote login is terminal handling. Also, Users want remote login to be as simple as local login. In this chapter 4.3 BSD rlogin client and server was described in a step wise mode. First a recording process was developed to understand the terminal line disciples and pseudo terminals.

## 15.10  TERMINAL QUESTIONS

1. Why RPC not pass parameters by reference?

2. Explain how Sun RPC maintains at-most-once semantics?

3. What are different terminal modes?

4. What are the different transparency issues with RPC?

5. Write short note on (a) Pseudo-Terminal (b) rlogin

6. What are the different forms of RPC semantics?