

1/1/19

Topics:

- 1) TCP } socket programming
- 2) UDP }
- 3) Fundamental of distributed computing
- 4) Synchronization
- 5) DC paradigm
- 6) Distributed object
- 7) Advanced DC paradigm
- 8) SOA (Service Oriented Architecture)
 - Value
 - characteristic
 - SOA concept
 - ↳ XML std
 - ↳ Web service
 - ↳ Service design principles
 - ↳ SOAP
 - ↳ WSDL (Web Service Description Language)
 - ↳ Basic SOA architecture (BPEL, ESB)
 - ↳ Restful web service, microservice
 - ↳ Data intensive computing
 - ↳ Message driven programming (JMS)

⇒ Inter process communication (IPC)

- pipe, FIFO, Shared memory
 - ↑ ↑ ↑ faster
 - comm. between comm. between (4)
 - two related processes two unrelated processes
 - located by : processes
 - fork() - (K)
 - (K)



U - user mode
K - kernel mode

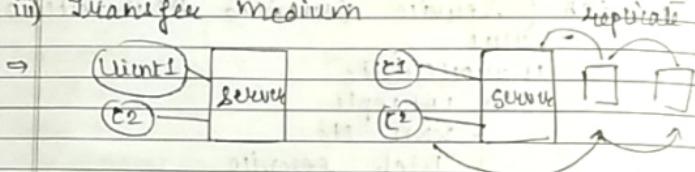
task/ Problem

↓
distributed into processes.

(same or different machines)

Type DC :

- i) Task
- ii) fixed medium
- iii) Transfer medium



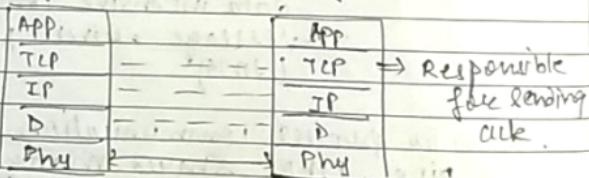
single point
failure)

forward request

Once the problem of
server is fixed, the instance
are destroyed, talk on dynamic
basis

Basic ex. of DC : socket prog
Other ex. : cloud computing, grid computing

• TCP / UDP (transport layer protocol)



Ack. must come from Transport layer.

Postal Service : When receiver calls for white book it's not ack.

when the receiver signs for the package & this signed receipt is sent to the sender ; it is called CLK.

SMS is also connection-less. We don't need to establish the connection for delivery we sending the message. WhatsApp messenger is connection oriented because it is end-to-end.

8/7/19 Client and server are user processes. TCP & IP are part of protocol stack within the kernel.

TCP/IP layers:

- i) APPLⁿ
- ii) Transport
- iii) Network
- iv) Host to Network (LL & Physical)

Read & write commands read/write from the TCP buffer (Transport layer) to the APPLⁿ (User process) or vice versa. Acknowledgement and retransmission are a part of TCP/IP stack which resides in kernel space. It has nothing to do with user space.

Responsibilities of layers:

- i) Physical
 - transmit message in form of electric signals.
 - amplification of signals.

ii) Host to Network

- Allows upper layers to access the media
- Controls how data is placed over the network (MAC)

iii) Network

- sending packets across one or more networks
- best effort efficient delivery
- Non-reliable
- Reliability must be implemented by upper layer

Protocol

- APP can bypass TCP or UDP
 - directly use IP (called raw socket)
Ex - ping
- APP can even bypass IP
 - directly read/write DLL frames
Ex - tcpdump

How to use Transport protocols?

- directly by APP
- indirectly by services such as FTP, SMTP.

A17119

POSIX std allows us to use the transport layer protocols

IETF is responsible for providing IPv4 & IPv6 to Unix

Port Number - used to identify process
(2 bytes)

Special addresses:

127.0.0.1 - Loopback address

0.0.0.0 - Unspecified address.

Multihomed hosts (Ex - Router)

A machine residing in two different networks. (2 NIC / 2 ethernet cables)

UDP - One socket can be used to send datagrams to different servers.

TCP waits for 1 round-trip time before resending the packet.

difference of time
when packet sent and
acknowledgment received.
(ICMP packet) (ping cmd)

17/19

TCP Connection Establishment

Three way handshake occurs

- connection is always initiated by the client

i) Server must be prepared to accept an incoming connection.

ii) Client issues an active open by calling connect.

iii) The server (TCP) must acknowledge (ACK) the client's (TCP) SYN and server (TCP) will also send its SYN.

server - process is sending data
 server (TCP) - TCP (buffer) waiting on
 the server

- Server is listening mode \rightarrow passive open
- To handle multiple clients - flock system call is used
- The no. of clients connected to a server are maintained in the backlog queue (on the server)

TCP Options :

- MSS (Max Segment Size)
 - client $\xrightarrow{\text{SYN}}$ server
 - $\xleftarrow{\text{SYN K, ACK SH}}$
 - $\xrightarrow{\text{ACK KH}}$

is for the connection.

May have to fragment if conn. pro. is slow

SYN J - max. data that server can send to client at a time.

SYN K - max. data that client can send to server (buffer empty)

- Window Scale Option (is for the machine)
 - 2 GB is the max. transfer rate (Max. Window) and not the connection speed.
 - Used so that the sender cannot overflow the receiver's buffer.

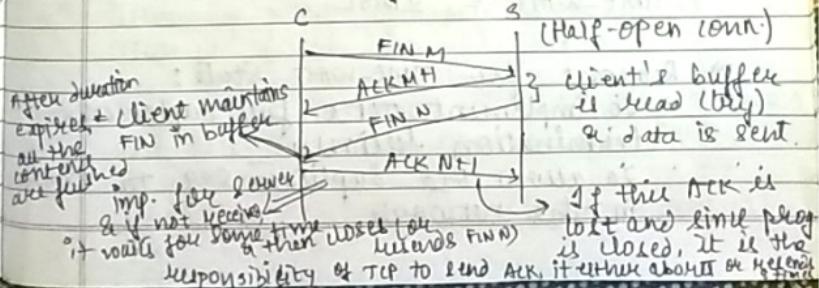
Timestamp Option

- Needed for high-speed conn.
- One packet is sent and reaches a router

- beyond which network is down. The router sends an ICMP and the sender in turn retransmit the packet through a different route. If both the packets make it to the destination at the same time, the packet with the latest timestamp is considered.
- AND used to distinguish two packets having same sequence no.
 - can also be used to calculate time to live value.

⇒ 99% conn are established by client but only 50% conn are terminated by the client.
if server does not have the right to close the conn, its active conn. queue may become full

close system call only means that the app is not going to send anymore data but, it can receive FIN segment if sent. In this time, the buffer is read and required data is sent by the receiver of FIN)



8/7/19

- ⇒ 11 states possible in TCP.
- ⇒ If ACK of client's seq. is sent along with server's reply, it is called piggybacking.
- ⇒ The end that performs active close enters in TIME-WAIT state. It maintains the contents of info for some time. It keeps this in the buffer in case the ACK for FINN is lost (App. layer buffer have been cleared).

If it does not wait then a new conn. requested by the client can be destroyed when FINN is received from the server (which was not meant for this new conn. at all).

- ⇒ MSL (Max. Segment lifetime) - time that any IP datagram can live on Internet.

TTL (Time to live): cannot exceed 255 hop count. TTL is the 8 bit limit for IP packet.

$$\text{TIME-WAIT} = 2 \text{MSL}$$

- ⇒ Reasons for TIME-WAIT State:
 - To implement TCP's full duplex conn. termination reliably.
 - To allow old duplicate seg. to expire in the network.

27/19
Port numbered.

Socket (IP address + port no.)

Ranges:

- Well known: 0 - 1023 (controlled by IANA)
- Registered: 1024 - 49151 (registered with IANA)
- Dynamic: 49152 - 65535 (ephemeral port)

Ex: while using Skype, the port no. is not controlled by IANA but, it is registered with it

⇒ In Unix, ports less than 1024 can only be assigned by super user process

[Kernel can choose IP add. & port no. from 1024 - 5000]

⇒ [*:21, *:*] → remote (client)

↑
ex - server
kernel can choose
any IP add. but,
port no. is fixed

We can only know
about this when
a request comes
in.

(Used when server is multi-homed)

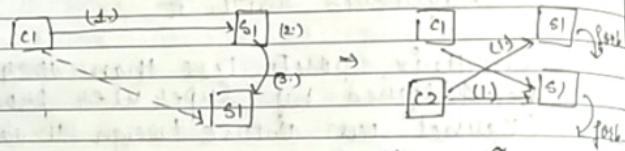
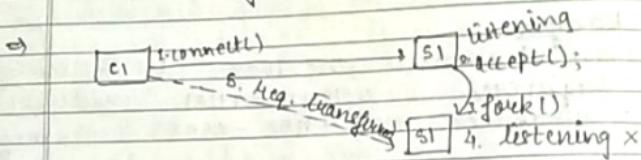
⇒ listening port of any server is always fixed.

⇒ Client starts on host with IP: 206.168.112.219
and does an active open to the Server's IP: 12.106.32.254
Client TCP (kernel) assigns ephemeral
port no. 1500
(server port no.: 21)

⇒ multiplexing, demultiplexing → which layers?

Page No. _____
Date _____

If the machine is not in any network
Kernel assigns IP addrs. all [127.0.0.1]



Hence, generates
4 child processes.

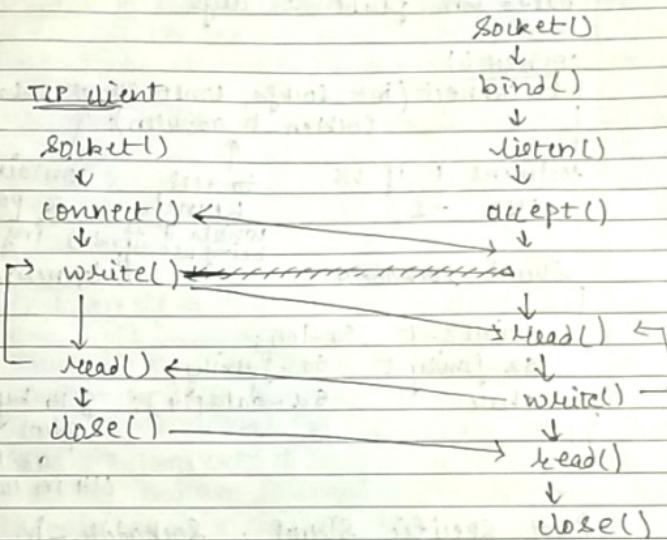
But, since UDP is not reliable, we
can handle many clients on a single
socket if they.

⇒ Seg. arrives with port no. 1500, it is
delivered to first child.

Seg. arrives with port no. 1501, it is
delivered to second child.

Any seg. destined for coming from
port no 21 is delivered to the
listening socket (parent).

TCP Sequence



Socket func.

int socket(int family , int type , int protocol)
return: non-negative descriptor if OK , -1 otherwise

family: AF_INET for . IPv4

Address family

type:	SOCK_STREAM	for	TCP
	SOCK_DGRAM	for	UDP

protocol: IPPROTO_TIP - TCP

IPPROTO_UDP - UDP

If we don't want to write this explicitly
write 0. Proto. is automatically chosen

based on struct & args.

connect()

int connect(int sockfd, const struct sockaddr *servaddr, size_t servlen + addrlen);

returns 0 if OK,
else, -1

to tell kernel,
which IP is being used, & port no.
contains info

struct sockaddr

// generic struct

{

uint8-t sa_len;

sa_family-t sa_family;

char

sa_data[14]; // 14 bytes of
phot. specific
addrs

};

14 in case of IPv6

Phot. specific struct : sockaddr_in
(for IPv4)

-ing → points
IPv4

192.168.31.2 - Not in numeric form, it is
a string. Thus, we need to
use structure.

15/7/19

address family in the struct is used
to specify family of service.

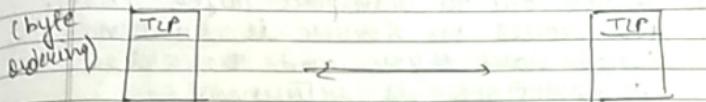
in-addr stores the 32-bit IP address of
server.

family, port no & IP addrs, are mandatory.
Other variables may be may not be
assigned.

bind() - Binds IP addrs & port no. to the process

int bind(int sockfd, struct sockaddr*, socklen_t len)

⇒ Network & should be Big Endian (well defined)



Client
mc itself knows about its endianness. Thus, if it is little endian it calls htonl (host to network) sys. call to convert it into big endian format. If mc is big endian, this step is skipped.

Server
mc, server knows about its own endianness. It calls the sys. call htonl (network to host) if it is little endian.

htons - for 16-bit
htonl - for 32-bit

presentation format - 192.168.0.1

dotted decimal (presentation) → binary (network)

Server: ntohs

const char *inet_ntop (int family, const void *addr, length)

Client: htons

Responses to connect()

- i) ETIMEOUT (Every time out)
Client's SYN did not reach server.
Client sends 8 times. (fixed)
after every (6-2 min)
- ii) Error due to Transport layer (RST)
No process on server is running
called hard error and the error
EINVALREFUSED is returned.

If server is busy, the client's packet
is simply dropped and no acknowledgement
is sent

HOSTUNREACH is generated when no
response is received and neither
ICMP is received.

= listen()

int listen (int sockfd, int backlog)

↑
maintains no. of
conn. at a given time

Backlog is sum of 2 queues:

- conn. which have completed 3-way handshake
- SYN received but, 3-way handshake incomplete

(these queues are maintained by the
kernel).

accept()

client's IP & port no. are read from TCP buffer and ntohs or ntohs il called and result is put into appln layer variable.
(struct sockaddr)

getsockname()

To get IP addre & port no. of socket.
(on machine itself)

getpeername()

IP addre & port no. of the remote machine (ex- client's IP & port no. when called on server, after conn is established)

echoTCPserver.c

```
#include <sys/types.h>
#include <sys/socket.h>
#include <stdio.h>
#include <netinet/in.h>
```

```
#define MAXLINESIZE 100
#define SERV_PORT 5555
```

```
int listenfd, clientfd;
char buffer[MAXLINESIZE + 1];
struct sockaddr_in servaddr;
int nBytesRead = 0;
```

```
void process(clientfd);
```

```
int main() {
```

```

if !(listened = socket (AF_INET, SOCK_STREAM)) {
    fprintf (stderr, "cannot create socket\n");
    exit (-1);
}

```

```

break ( & servaddr, sizeof (servaddr));
servaddr.sin_family = AF_INET;
servaddr.sin_port = htons (SERV_PORT);
servaddr.sin_addr.s_addr = htonl (INADDR_ANY);

```

```

bind (listened, (struct sockaddr *) & servaddr,
      sizeof (servaddr));

```

```

listen (listened, 5);

```

```

for (;;) {
    clientfd = accept (listened, (struct sockaddr *) & clientaddr,
                        NULL);
    if (fork () == 0) {
        close (listened);
        processclient (clientfd);
        close (clientfd);
        exit (0);
    }
    close (clientfd);
}

```

7/7/19

echoTCPClient.c

→ bzero function assigns '0' value to all its data members.

bzero (&recvaddr, sizeof (recvaddr));

↑
pointer (address) ↑
size of the
to the struct

→ write (connected, sendBuffer, strlen (sendBuffer) + 1);

(even if only if sizeof is used then
5 bytes are sent like 10b bytes are sent)

In Server prog:

- Addresses are assigned at transport layer and socket exist at data link layer. Hence, assigning the family while assigning address is required.
- Use (listen()) inside if(fork() == 0) i.e. to avoid creating multiple child processes (from children of parent). Only the parent process should be able to create a child.
- Use (clientd) outside the if(fork() == 0) is written to let the parent process know about the termination of the child process denoted by clientd.
- When one of the client is terminated abnormally (window is closed instead of Ctrl+C), the server is not notified immediately and the child becomes a zombie process (since, the descriptor connected / allocated with the child process (of server) does not exist anymore, but, the parent process does not know about it and assumes that its child is still communicating; it is said to be a zombie process).

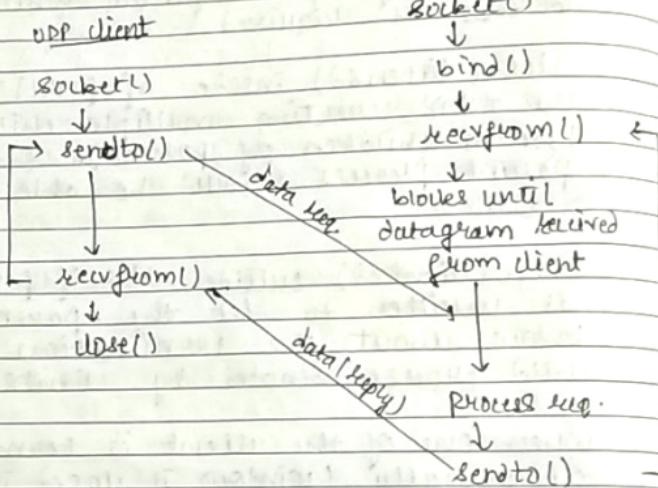
Qn 7th UDP Client / Server

```
server
int n;
socklen + len;
char msg[MAXLINE];
for (i;)
    len = uilen;
    n = recvfrom (sockfd, msg, MAXLINE, 0,
                  peeraddr, &ulen);
    sendto (sockfd, msg, n, 0, peeraddr, len);
```

if client & server are on different machines
we need to use port no. on one of the
machines (preferably server)
if htons → what adds. kernel a.
to the add. of the
port

```
client
int n;
char sendline [MAXLINE], receive [MAXLINE];
while (fgets (sendline, MAXLINE, fp)
      & sendto (sockfd, sendline, strlen (sendline),
                peeraddr, servlen) == -1)
    perror ("error sending message");
```

UDP Server



→ `recvfrom()`

`ssize_t recvfrom(int sockfd, void *buff,
size_t nbytes, int flags,
struct sockaddr *from,
socklen_t *addrlen);`

Stored in Kernel space
and to bring it
into user's memory

`INADDR_ANY` - gives liberty to kernel
to choose an address. If
the addre. is not assigned,
kernel assigns 127.0.0.1 to it

Server

```

int n;
socklen_t len;
char mesg[MAXLINE];
for(;;)
{
    len = ulen;
    n = recvfrom(SOCKfd, mesg, MAXLINE, 0,
                  pclient, &len);
    sendto(SOCKfd, mesg, n, 0, pclient, len);
}

```

- If client & server are on different machines, we need to use pton func. on one of the machines (preferably server)
 If htons → what addrs. kernel assigns is the addrs. of the other party.

client

```

int n;
char sendline[MAXLINE], recvline[MAXLINE+1];
while ((getchar(sendline, MAXLINE, fp) != NULL)
{
    sendto(SOCKfd, sendline, strlen(sendline), 0,
           pclient, &len);
    n = recvfrom(SOCKfd, recvline, MAXLINE, 0,
                 NULL, NULL);
}

```

↙ /
 Server is the sender (we know) and we already know its IP (we sent the data to it) so no need to receive IP again.

23/7/19

Since, there is only one server, and it is communicating with a client, any other client's request is simply dropped (or put into the buffer if only maintained)

Multiplexing & demultiplexing occurs at transport layer only. In case of TCP, it is due to the mechanism that the server is able to forward client requests to appropriate child (server) process.

To know about socket options after it is created:

```
int getsockopt(int sockfd, int level,  
              int optname, ↑  
              void *optval,  
              socklen_t *optlen);
```

13 data link
levels + transport
network

In case of TCP, it can be called only after it is created but, before connect is called.

For server, socket options can be set only before listen is called.

getpeername() is of no use in UDP and is called after connect on client and after accept on server.

23/11/11

Page No. _____
Date _____

optval - whether socket is broadcasting
(value = 1) or not.

Binary options (true/false)

Value options (specific values / structures)

Options:

SO_BROADCAST (binary opt.) (can be get & set)

SO_ERROR (cannot be set)

SO_KEEPALIVE (for TCP only) [if there is no activity for 9 hours, server sends a mes. to that particular client. If ACK is not received, server closes the connection]

SO_RCVBUF, SO_SNDBUF - size of send & receive buffer

SO_LINGER - delays comm. protocol for some time. (for TCP only) when close is called

Options for network level

IP : LEVEL = IPPROTO_IP

TCP: LEVEL = IPPROTO_TCP

Options: SO_KEEPALIVE, NODELAY

→ If dest is also broadcast address,
SO_BROADCAST has to be set to 1; else
EACCES is returned.

Possible peer responses for KEEPALIVE:

• ACK (everything OK)

• RST (peer crashed & rebooted): ECONNRESET

• NO response to keep alive probe

RST

- Socket's pending timer is set to ECONNRESET
- Socket is closed. (the socket of the sender who sent ACK KEEPALIVE probe)

NO response

- 8 additional probes (each 75 sec. apart) are sent
- If any of these probes, if an ICMP is received (from a router), the router keeps the packet alive and probes 8 more times. If still no response, then EHOSTUNREACH is returned.

Changing the waiting (keep alive) time is not advisable as these parameters are set based on per-kernel basis. On modification, it can cause problems on every host.

a keep alive probe
when the ACK does not get any response, the same sender sends FIN because it may happen that the peer is rebooting and when it regains its control, it knows that the conn was terminated. Otherwise, it stays under the impression that it is still connected to that particular host.

- ⇒ When a process completes, all the sockets are automatically closed but, if some data is present in the buffer, it can be lost. Thus, linger is used.
It waits for 1 linger seconds before closing the socket completely (after closed is called)
If any data is found, the process is put to sleep until:
- either all data is send & ack. by peer
TCP
- or linger time expire.

If socket is nonblocking (asynchronous) it will not wait for linger seconds and closes immediately.

The problem that if some data is left in the buffer even after linger seconds, it is discarded; is handled by SHUTDOWN.

Shutdown wait for FIN (only sent when every data has been read & written according to requirement).

SHT-RD

SHT-WR

SHT-RDWR

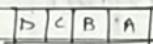
UDP does not send ack. but, if we want we can create an application level function which sends an intimation

(1 or 'Y' or 'Yes') to the sender as soon as data is received (even if it further processes the data later)

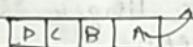
Transport (segment)

↓ → in between exists as MTU which depends on the DLL (frame) network config.
 (MTU = Max. Transfer Unit)

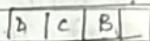
⇒



queue, if new packet arrives, it is dropped



A is head



Even now, if a new packet arrives, it is dropped.
 Hence, we need to use an implement a circular queue only.