

August 21, 2023

## Lecture 1

- Application layer → presentation → session → application
- Transport layer
- Network layer → routing
- Data Link layer
- Physical layer

Switching devices of Datalink Layer

- Hub
- Bridge
- Switch

Hardware:

- End host (source & dest.)
- Switches / Host

Software:

- Protocols of above 5 layers

Application Layer Protocols:

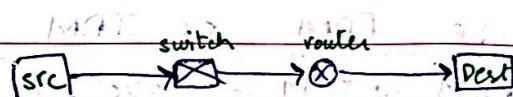
- HTTP
- FTP
- SMTP

Transport Layer Protocols:

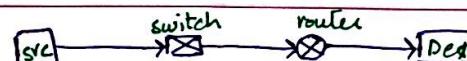
- TCP (connection oriented)
- UDP (connectionless protocol)

Network Layer Protocols:

- IP



message reaches dest at any cost  
TCP is used in emails



complete message may not reach dest  
UDP is used in video streaming

Data Link Layer Protocols:

- Ethernet
- ARP

KAGHAZ  
www.kaghazpk

Wednesday, August 23, 2023

ESSS, IIS Tangail

Lecture 2

Date: 23/08/23

Page No. 1

Endhost :- devices which are used by users.PhysicalMobile Phonesanalog signals 

call specific numbers

which change accordingly

Satellite Phonesradio signals 

specific extensions

available without charges

DatalinkBridge/SwitchNetworkRouter

• 802.11 a RFC document

• 802.11 n available for 5G

Edge router :- which connects one network to another

↳ edge / internet service gateway / gateway

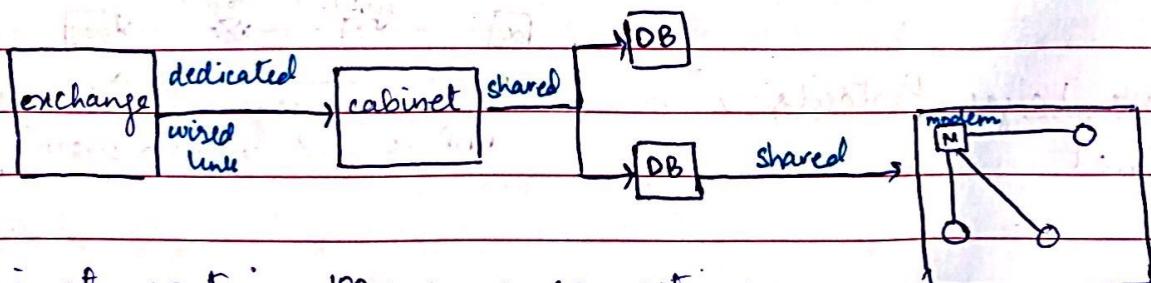
Cable Based Network

use FDM or TDM

FDM:- frequency distributed multiplexing

FDM vs TDM :- have diff bandwidth

TDM:- time distributed multiplexing



cabinet contain 100+ home connections

DB contain 10 home connections

DSL :-KAGHAZ  
www.kaghaz.pl

Monday, August 28, 2023  
Lecture 3

unicast:

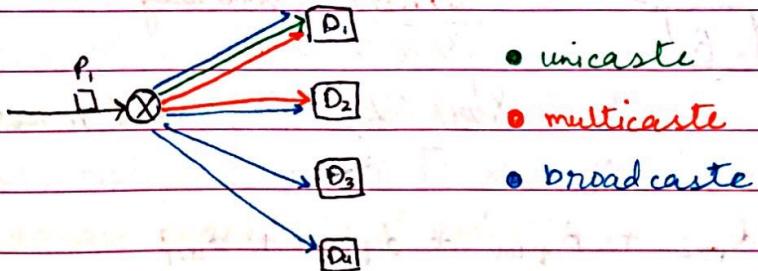
If a packet is forwarded by the switching device to a single digital device

multicast:

If a switching device forwards a packet to more than one digital device

broadcast:

If a switching device forwards a packet to all devices in the range to that switching device but is upto the digital device to accept or reject the packet

Transmission delay

the time a packet takes to completely push onto the line

one hop numerical example

$$\frac{L}{R} = \frac{10 \times 10^3 \text{ bits}}{100 \times 10^6 \text{ bps}}$$

$$= 0.1 \times 10^{-3} \text{ s}$$

$$= 0.1 \text{ ms} \quad \text{Transmission delay}$$

radio range

88 - 108 MHz



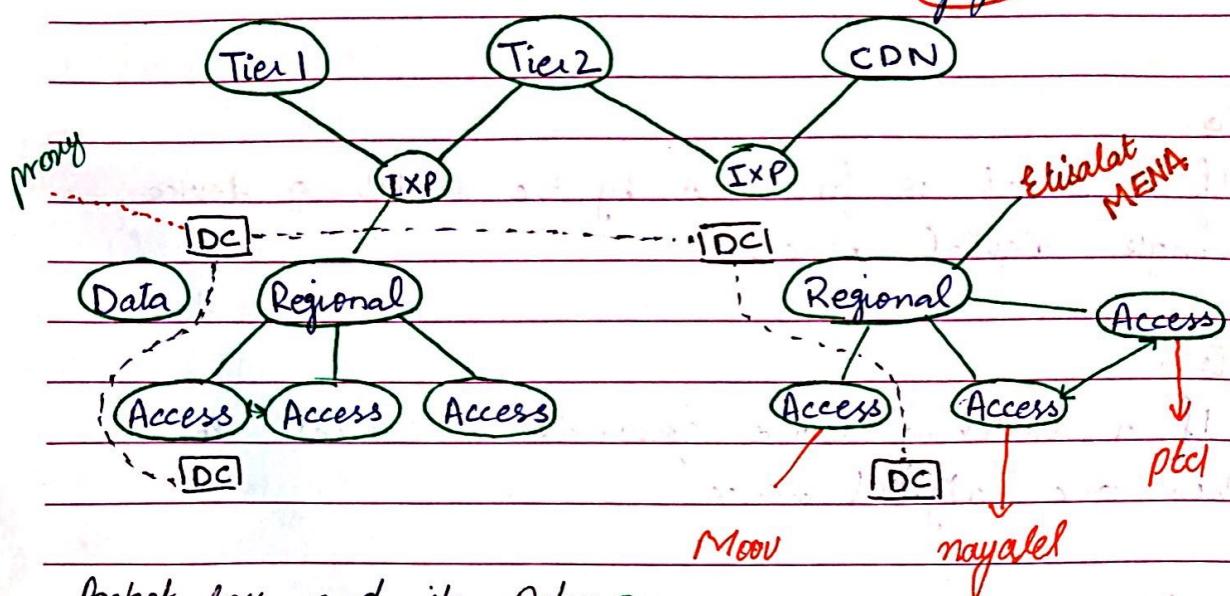
KAGHAZ  
www.kaghazpk

day / date:

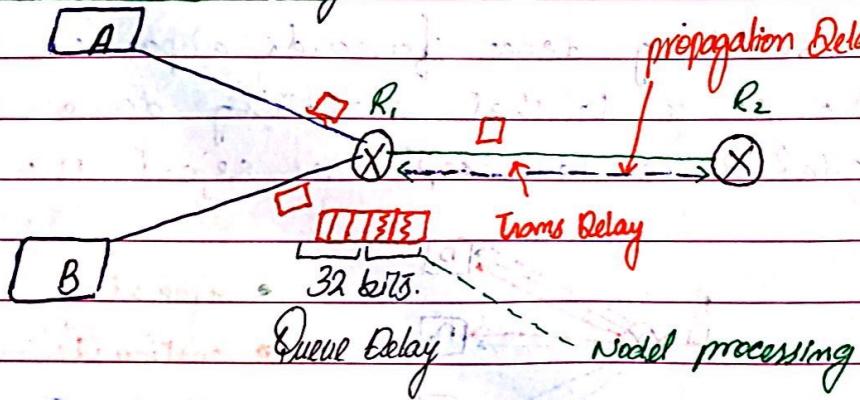
Wednesday, August 30, 2023

Lecture 4

google



Packet loss and its Delay :-



$$d_{\text{node}} = d_{\text{proc}} + d_{\text{queue}} + d_{\text{trans}} + d_{\text{prop}}$$

Packet Queuing Delay :

L<sub>q</sub> : arrival rate of bits  
R service rate of bits

# CHAPTER 1

Nauman Moazzam

## Lecture 1.1

Network Edge :- comprises of end systems / hosts

End Systems :- either request services or provide services

• client and server

• server is not part of the network core

### Server

- provide access to network resources

- always present on host

e.g I request some service on google. Here google is the host from which we are (client) want some kind of service. Here the server is present at google which would provide the service which was requested by the client

- servers have permanent IP address

↳ for websites or services which are accessed 24/7 around the world if the IP address is dynamic then routing table would have to be updated constantly.

\* static IP address makes it easier to find you via DNS

DNS - turns domain names into IP address

### Client

- may request services

- have dynamic IP addresses

whenever a device is connected to the internet an IP is assigned to that device. When we shut down the device the IP is returned back to the pool of addresses & a new IP is assigned when connected again.

- do not communicate with each other

## Access Network & Physical Media

↳ how host connect to any network or internet  
connect end systems to edge router / first router

Edge Router :- last internally managed routing device that connects one network to another

telephony or video conferencing apps are mostly peer-to-peer

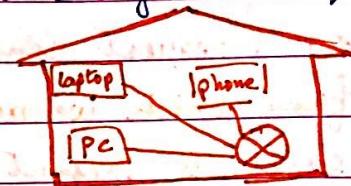
First Router :- the first device which is

connected to the device or server in a network

local ISP routers

Access Network: is a network which physically connects and end system to the immediate router (edge router)

e.g. in a home network where internet is connected by stormfiber, the edge router is not present in the house, the edge



router is present at the office of stormfiber &

the edge router is connected to the devices via fiber cables

## Transmission Modes

→ unicaste :- one to one communication

→ multicast :- one to a selected group communication

→ broadcast :- one to all members of the group communication

## Access Network Technologies

### 1. DSL (digital subscriber line)

• use the telephone wires to transmit voice & data signal

• the data over DSL line goes to Internet

• voice over DSL line goes to telephone network using splitter

### 2. DSLAM (digital subscriber line access multiplexer)

multiplex & demultiplex signals through DSL

internet signals are forwarded to ISP & voice signal

dedicated access network



KAGHAZ  
www.kaghazpk

## 2. Cable Modem / Network

- similar to DSL
  - same network is used to provide data signals & TV signals
  - splitter is used at home to split the signals accordingly
  - & at the cable network end CMTS is used to demultiplex the signal to cable network or ISP
  - the difference b/w DSL and Cable network is that DSL is a dedicated access network & Cable is a shared access network
- shared network is a broadcast medium

## 3. Ethernet (enterprise access network)

- commonly installed wired LAN technology

• used to connect devices together on a Network to transfer broadband data.

- upto 10 Gbps transmission rates

## 4. Wireless Access Network

- Shared

- connects end systems to router via access points

## 5. Physical Media

### Guided Media

signals propagate in solid media

e.g. copper, fiber, coax

### Unguided Media

signals propagate freely

e.g. radio



# Nauman Moazzam

## Lecture 1.2

### Physical Media

- coaxial cable
- fiber optic
- radio

### Approaches to Sharing

#### Reservation

- we reserve a line for permanent use through physical medium

#### circuit switching

e.g. landline

#### On demand

- whenever we need to send / receive a line is created after delivery link is returned back

#### packet switching

e.g. Internet

### Packet Switching

- data is broken in packets
- if packets are created on demand
- if I need to send some data only then I will use the link otherwise link will be free from my side
- it is not necessary for each packet to take the same route to reach its destination
- each packet is considered as individual data
- each packet while have its own src address & dest address
- then the network edge device will reassemble the data & then used

→ addm. admission control : per packet

$$\text{Packet Transmission} = \frac{\text{L bits}}{\text{R bits/sec}} \cdot \text{no. of bits of Packet - length}$$

data network  $\Rightarrow$  routers

telephone network  $\Rightarrow$  switches

day / date:

## Circuit Switching

circuit: connection

- when two device want to communicate a circuit is established

- now even if no data is send or received the circuit is reserved by the device / user

e.g. if I call someone, a circuit is established b/w the two devices. Now if we talk for 5 mins & for the next 15 mins both muted their devices & no one talks. No data is transmitted in those 15 mins but a circuit is reserved even if no data is transmitted since the call is not terminated.

- since the circuit is reserved b/w devices, the resources are dedicated to them & no other device can share these resources.

- doesn't route around trouble

## FDM vs TDM

**FDM:** Frequency division multiplexing

The frequency (bandwidth) is divided equally among the users sharing a link & each user use that allocated line to share data. If a user is not sending data then at that time the frequency is being wasted.

The allocated frequency is not shared.

**TDM:** Time division multiplexing

Instead of frequency, time is divided among users. If there are 4 users then each gets 15 secs of 1 min total in a cycle until time is completed. In that allocated time the user will use full frequency to send data. If there is no data for the user then that time cannot be used by other users.



KAGHAZ  
www.kaghazpk

## Numerical

How long does it take to send a file of 80 Kbytes from host A to host B over a circuit-switched network?

All links are 1.536 Mbps. Each link uses TDM with 24 slots/sec. Time to establish end-to-end circuit is 500 msec.

$$80 \text{ Kbytes} = 640,000 \text{ bits}$$

$$1 \text{ Kbyte} = 8000 \text{ bits}$$

$$1 \text{ Kbyte} = 1024 \text{ bytes}$$

Link transmits 1.536 Mbps per sec

There are 24 slots per sec

Each circuit has  $\frac{1.536 \times 10^6}{24} = 64000 \text{ bps}$  (One slot has)

Total data : 640,000 bits

It takes  $\frac{640,000}{64000} = 10 \text{ sec}$  (To transmit)

To establish link :  $500 \text{ msec} = 0.5 \text{ sec}$

$$1 \text{ msec} = 0.001 \text{ sec}$$

Total time :  $10 \text{ sec} + 0.5 \text{ sec}$

$$= 10.5 \text{ seconds}$$

Packet Switching: Store and Forward

entire packet must arrive at router before it can be transmitted to the next line.

Packet Loss and Delay

- delay occurs when packet arrival rate to link exceeds output link capacity
  - router temporarily stores / buffers packets in its queue  
↳ queuing
  - If the queue fills up & packets keep arriving then the router drops those packets which leads to loss.
- happens in case of congestion in links

4 types of Delays

→ queuing : first in first out

→ processing : router extracts the src & dest IP address & index the if

→ transmission : either given or negligible in numerical

→ propagation : time taken from src to dest affected by speed of light

to get end-to-end delay of a packet we need to calculate all the delays & add them

$$d_{\text{total}} = d_{\text{proc}} + d_{\text{queue}} + d_{\text{trans}} + d_{\text{prop}}$$

$$d_{\text{trans}} = \frac{L}{R}$$

$d_{\text{prop}} \rightarrow$  d length of physical link  
s propagation speed ( $\approx 2 \times 10^8$  m/s)

Transmission delay is the amount of time required for the router to push out the packet. Propagation delay is the time it takes to travel from 1 router to the next.

## Traffic Intensity

La

R

L: packet length

a: arrival rate

R: link bandwidth

• if  $La/R \approx 0$  avg queuing delay is small• if  $La/R \approx 1$  avg queuing delay is large• if  $La/R > 1$  more work is arriving than it can be processed.

## Throughput

rate at which bits are transferred b/w sender

&amp; receiver

e.g. while downloading a video a packet travels through diff links with diff bandwidths. Since each packet may take a diff route, each packet will take diff time to reach the dest. The avg time taken by these bits is throughput.

- The throughput is generally the rate of the link which is minimum.

avg throughput for diff. links =  $\frac{1}{\max(\text{latency})} \times \min(\text{bandwidth})$

all diff. bandwidths exist for different diff. protocols

so select diff protocol. During this the avg

## CHAPTER 2

Nauman Moazzam

Lecture 2.1

### Application Layer

- specifies the shared protocols & interfaces methods used by hosts in a communication network
- helps to identify communication partners & synchronise communication
- ensures an application can effectively communicate with other applications on different computer systems & networks.
- application layer do not exist on core network devices e.g. routers & switches meaning a developer doesn't have to write code for routers / switch

core network device

- network layer
- data link layer
- physical layer

end user device

- application layer
- transport layer
- network layer
- data link layer
- physical layer

### Application Architecture

- client - server
- peer - to - peer → self scalability

## Process Communicating

processes communicate within a host using inter-process communication | OS on servers  
 ↳ set of rules defined by OS in mostly Linux

process in different hosts communicate by exchanging messages

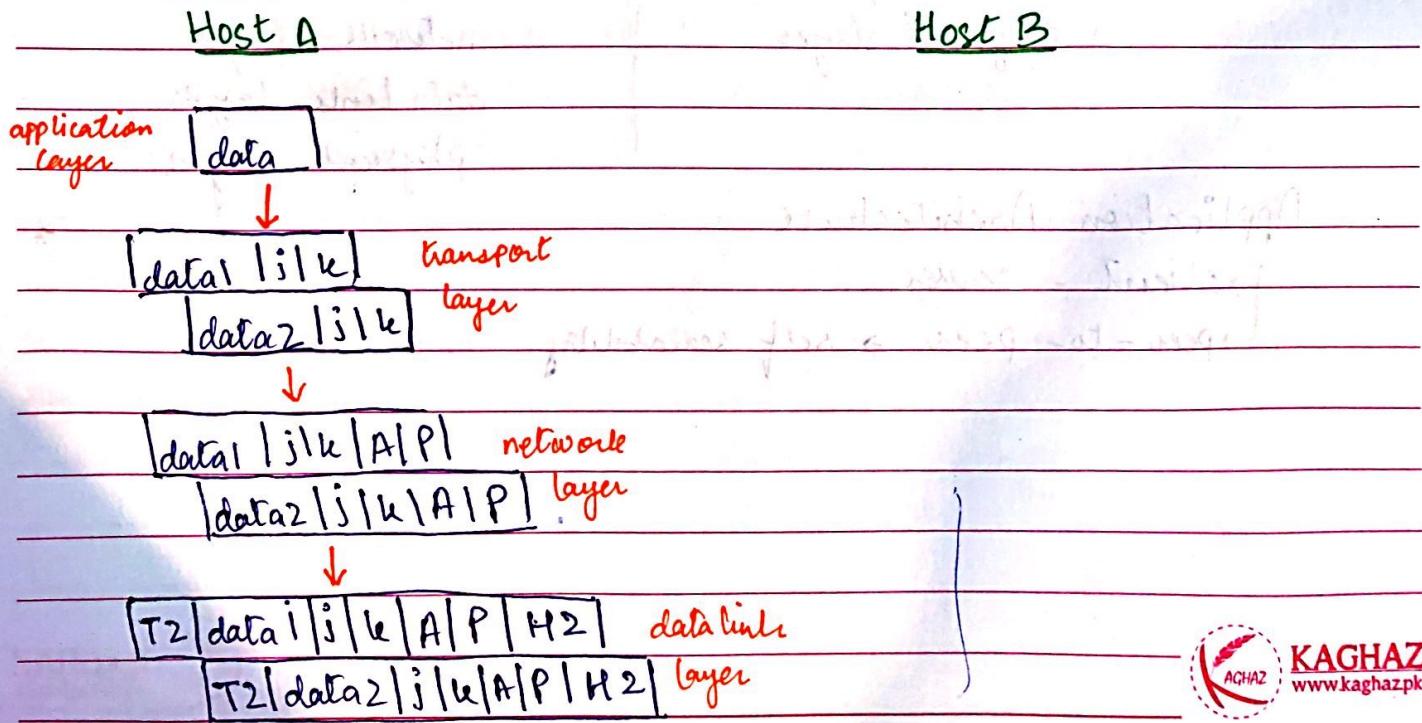
client process: process that initiates communication

server process: process that awaits to be contacted

each process opens a socket in the transport layer  
 & send its data in that socket

socket: communication link b/w 2 processes running on the network

combination of port number & IP address



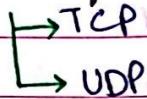
- host A send data to application layer
- application layer breaks the data into packets & send it to transport layer through sockets
- transport layer attach 1 port address j & k to the data packets  
 j :- address of sending process  
 k :- address of receiving process
- then send it to network layer
- network layer append the source & destination IP address to the data packets
- then send to data link layer which attach the header & trailer to the packets
- then the physical layer converts the data into bits & send to the internet
- port number helps the end host to identify which ps is the correct process to send the data.

### Transport Services

- data integrity  
 ↳ e.g. apps for file transfer require 100% reliable data transfer
- timing  
 ↳ e.g. internet, telephony require low delay to be "effective"
- throughput  
 ↳ e.g. multimedia
- security  
 ↳ e.g. bank transfer



## Transfer Layer Protocols



UDP exists because some apps still do not require so many services which TCP provide.

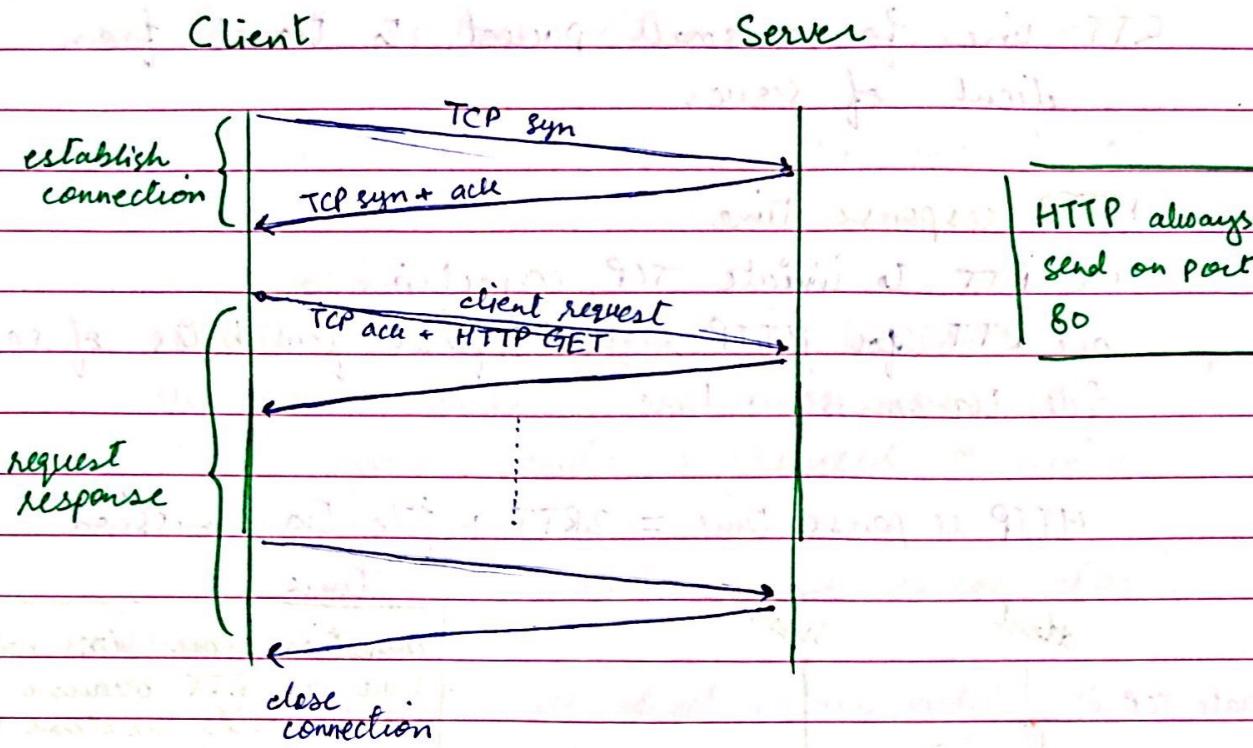
UDP is a light-weight protocol which do not require a large header.

## Lecture 2.2

- differentiating aspect of Internet from other competitor networks was on-demand

### HTTP

- web's application layer protocol
- when we enter some kind of URL in the browser, it send a message to the server known as HTTP request
- This request is received on port number 80
- the server in return sends objects
- HTTP uses TCP protocol
- HTTP is stateless, meaning that the server maintains no prior information about the client searches if information needs to be saved the cookies is used



### Non-Persistence HTTP

at most only 1 object is sent over a TCP connection

- TCP initiate

- client initiate TCP connection on port 80

- server at host accept connection on port 80 & notify client

- client send request message into socket

- server forms response message to the request & send the message into its socket

- client receives response message  
server closes TCP connection

If a message request contains 10 jpeg images  
then these steps will be repeated for all 10 images

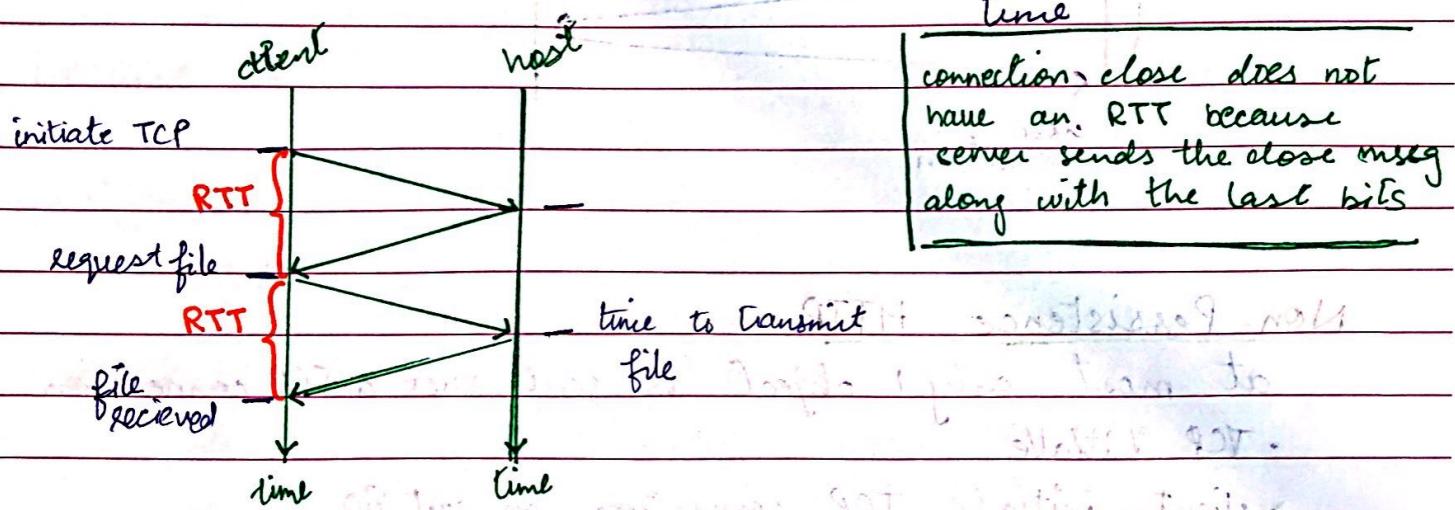
RTT = time for a small packet to travel from client to server

HTTP response time :-

one RTT to initiate TCP connection

one RTT for HTTP request & first few bytes of response  
file transmission time

$$\text{HTTP response time} = 2\text{RTT} + \text{file transmission time}$$



### Persistent HTTP

- server sends the base HTML file to client but keeps the connection open
- if a request message contains 10 more reference objects server will use the same TCP connection to send those files/objects
- all this takes as little as 1 RTT



persistent pipeline often

persistent HTTP often use the concept of pipeline

Non-pipeline:-

clients sends 1 request & in response

the server sends 1 response message.

client again sends a request & client sends a response.

this continues until the connection closer.

Pipeline:-

client requests all objects in one batch & the server will receive all the requests and

then send the response objects.

## HTTP request message format

- first line is request line

<method> <URL> <version> <br> <lf>

method include: POST, GET, PUT, HEAD

version of HTTP: 1.0, 1.1, protocol version

cr: carriage return \r

lf: line-feed \n

- the next couple of lines are header lines

- a line with only cr & lf indicate the end of header lines



## Form Input

when we type something in the webpage

→ POST method:

some webpages include some kind of form

e.g. while shopping online, that input is

acquired by the server using POST method

which is optional in the body of request

message

→ GET method:

when we type something in the search bar

what we type is appended in the URL

## Method Types

### HTTP / 1.0

GET

POST

HEAD :- often used

for debugging, do not return any object

### HTTP / 1.1

GET

POST

HEAD

PUT :- client upload a file in server

DELETE :- client delete a file in server

## HTTP response message

→ first line contain the status line

<method>/<version><status code><cr><lf>

some examples of status code include

(2xx) 200 OK : request succeeded

(3xx) 301 Moved Permanently : object moved

(4xx) 400 Bad Request : message not understood by server.

(4xx) 404 Not Found : requested doc not found on server

(5xx) 505 HTTP Version Not Supported

→ next couple of lines are header lines

DATE :- when the client made the request to server

Last-modified :- when the requested object was updated

→ at the end is the data objects which were requested

## Cookies :-

- used to maintain the state of 'stateless' HTTP

- client stores small state on behalf of the server

- when accept cookies our state is saved in a small database on the server corresponding to the id assigned to us by cookies

when we again request the same HTTP the cookies id is sent to the server along with our request message

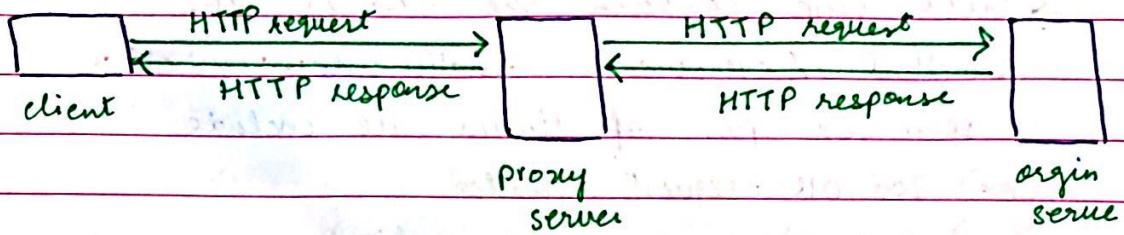
e.g. in daraz, we first accept cookies, our information is saved in the database. After sometime if again add something to the card & checkout our information is saved because of the cookies.

cookies is maintained it is deleted  
or the cookies expire



## Web Caches

compliant HTTP protocol to reduce traffic



- whenever a client performs a request if the objects are stored in the cache server. if another client performs the same request then the proxy server do not forward the request to the origin server, but it returns the objects requested by the second client from as they are stored in the cache of proxy server
- In this way the traffic is reduced.
- proxy server acts as both server & client
- typically cache is installed by ISP

### Example

avg object size : 100k bits

avg request rate from browser to origin server : 15

avg data rate to browser : 1.50 Mbps

internet delay RTT from institutional router to origin server : 2 sec

access link rate : 1.54 Mbps

LAN utilization :  $\frac{1.50 \text{ Mbps}}{10 \text{ Mbps}} = 15\%$

or

$$\text{La/R} = \frac{(100 \text{ k bits} \times 15)}{10 \text{ Mbps}} = 0.15 \text{ (negligible)}$$

(negligible)



Access link utilization:  $1.50 \text{ Mbps} / 1.54 \text{ Mbps} = 97\%$

or

$$La/R = (100 \text{ kbytes} + 15) / 1.54 \text{ Mbps} = 0.97 \text{ (Huge)}$$

Total delay = internet delay + access delay + LAN delay  
 $= 2 \text{ sec} + \text{minutes} + \mu\text{sec}$

97% utilization is equal to 0.97 traffic intensity

Now suppose that,

local web cache is installed

cache hit rate is 0.4

40% requests satisfied at cache

60% requests satisfied at origin

Access Link utilization:  $(0.6 * 1.50) / 1.54 = 0.58 \approx 58\%$

Total delay = internet delay + access delay + LAN delay  
 $= 0.6(2 + \frac{\mu\text{sec}}{1000}) + 0.4(\mu\text{sec})$   
 $= 0.6(2.01) + 0.4(\sim \mu\text{sec})$   
 $= \sim 1.2 \text{ sec}$

## Lecture 2-4

| email is asynchronous

### Electronic Mail

→ user agents

present at user end & systems

→ mail servers

contain dedicated mailbox & message queue

→ SMTP (app layer protocol)

simple mail transfer protocol, push protocol

- the sending user agent push the mail in the outgoing queue of mail server
- mail server use SMTP to push the mail to the user mail box to the receiving user agents mail server
- hence SMTP is a push server
- now the mail is in the mailbox of receivers mail server
- the receiver pulls the mail from the mail box at their convenience , this pull protocol MAP

pull protocol means that if the receiver do not have internet connection in the receivers mail box got a mail then when the receiver will get internet connection they the mail will be downloaded , means that it will be pulled at the receivers convenience.

### SMTP : RFC

- uses TCP
- error message is displayed
- sender & receiver maintain a direct link
- 3 phases of transfer
  - SMTP handshaking initiate TCP
  - SMTP transfer of message persistent TCP
  - SMTP closure

in-band :-

out-band :-

SMTP is out-band

7-bit ASCII response / command



KAGHAZ  
www.kaghaz.pk

## Comparison with HTTP

HTTP

pull protocol

ASCII command/response

object encapsulated in its own response msg

SMTP

push protocol

ASCII command/response

multiple objects sent in multipart msg

## DNS

domain name system

- acts as an identifier for IP address

- user type a domain name in web browser; DNS server translate requests for names into IP addresses; control which server an end user will reach
- FTP invoke an IP address against a domain name
- a single web server may have multiple DNS servers, different DNS queries must be distributed among these servers

in the same way a web services have multiple web servers & every web server have an IP address.

DNS distribute the load among multiple web servers

e.g millions of people access google.com at a time

DNS distribute the load by assigning IP address of web servers to those DNS queries

- some web services have a canonical name & an alias name; we type the alias name & DNS converts it to its canonical name e.g CNN.com



KAGHAZ  
www.kaghazpk

DNS servers are in 3 layers:

- Root DNS server layer
- TLD Top level DNS server

| 13 DNS root servers |

authoritative server layer  
(server of an particular organisation)

- . when we type something in browser
- . DNS service translate domain name to IP address
- . HTTP protocol invoke DNS protocol
- . DNS protocol sends the DNS query to Root DNS server
- . Root DNS transfer the query to a correct TLD server  
e.g. we type FAST.edu in browser the root DNS will send the query to edu DNS server at TLD layer

Root DNS will transfer the query corresponding to the Domain name

- . Each domain name have a TLD server layer
- . TLD server will then send the query to correct organisation server
- . Local DNS server works similarly to proxy server web cache when a DNS query goes to DNS server from an organisation the response is stored against the domain

### Iterated Query

in response to a DNS query, it provides reference to another DNS server



KAGHAZ  
www.kaghazpk

day / date:

## Recursive Query

when one server communicates with several other DNS servers to hunt down an IP address & return it to the client.

## Caching) DNS Information

### DNS records

(name, value, type, ttl)

- canonical name of mail server
- host name :- can canonical name (real name of server)
- domain name :- alias name
- host name of authoritative name server

type = A :- authoritative DNS

type = NS :- TLD server

type = CNAME mapping b/w alias & canonical name

type = MX similar to CNAME but for mail servers

# CHAPTER 3

Nauman Moazzam

Lecture 3.1

- transport layer provides services to the application layer
- provide logical communication b/w apps in diff devices
- run in end systems
- when data packets are in transport layer they are called segments

↳ data packet + Transport layer header



Port addresses

Transport layer:- process to process delivery

Network layer:- responsible for node to node delivery

Link layer:- responsible for delivery b/w neighbouring nodes

e.g. host to first router

a single node contain many processes

e.g. server & client are using html app so the process responsible for that html app ; communication b/w those processes is done alone by transport layer

## Functions of Transport Layer

→ provide process to process communication

↳ accomplished by port number

→ provide end-to-end error checking, error control & flow & congestion control

both TCP & UDP

only TCP

only TCP

TCP is a reliable protocol because it provides error control and flow & congestion control whereas as UDP is unreliable since it doesn't provide these

TCP & UDP both tell if error occurs but only TCP

tells how to resolve the error



KAGHAZ  
www.kaghaz.pk

## Multiplexing and Demultiplexing

- multiplexing takes place on sender's side
- demultiplexing takes place on receiver's side
- on port number in transport layer

### Port

- Transport layer identifier
- to choose among the multiple processes running on same host
- is a 16 bit integer  
0 - 65,535
- on a single client (laptop) there are multiple processes running like HTTP, spotify, skype, etc. Every process will have a port number which will help to identify which process send a message

IANA had divided port numbers into 3 ranges:

→ Well known: 0 - 1023; assigned to processes e.g. HTTP - 80

→ Registered: 1024 - 49,151; registered not assigned

→ Dynamic: 49,152 - 65,353; can be used by any process

**Multiplexing** :- gather data from multiple sockets; envelop data with a header

**Demultiplexing** :- receive segments; deliver data to correct app layer processes by using info from header & identifying correct sockets; remove header



## How / Demultiplexing Works

- Suppose A is the sender & B is the receiver
- both have like & WhatsApp installed
- A wants to send WhatsApp message to B
- A must mention IP address of A & B & and the port number of WhatsApp
- If A wants to send message on like as well
- A must mention IP address of & B and the port number of like while sending the message
- The message from boths apps will be wrapped up along with appropriate headers

header :- source IP address, dest IP address, source portNo ,  
dest port no

- a single message will be send to B
- This is called multiplexing
- At B, the received message will be unwrapped.
- messages from WhatsApp & like are sent to appropriate apps by looking at dest port no.
- This is called demultiplexing

## Connection Connectionless Demultiplexing

- used by UDP
- client select an ephemeral port no randomly e.g 9157
- transport layer will attach the source port no : 9157 &  
the dest port no (well known) : 6428 to the packet
- now at the server side transport layer will strip  
the header & from the dest portno:6428 send the packet  
through correct socket to the appropriate app



- UDP socket is identified by 2-tuple
  - dest IP address
  - dest port no.
- , at host transport layer
  - UDP segment arrives
  - checks port no dest port no
  - directs UDP segment to appropriate socket
- UDP violates the layering principle
  - network layer headers include source & dest IP address
  - transport layer uses the information of network layer
  - several violations exists
- in case of same dest port no & diff source port no
  - the data packets will be delivered to the same socket

### Connection-oriented Demultiplexing

used by TCP

TCP socket is identified by 4-tuple

- source IP address
- source port no
- dest IP address
- dest port no.

- now in case of same dest port no, transport layer can perform segregation on the basis of source IP address & source port no.
- , every socket will create a diff socket in server



User Datagram Protocol

UDP

barebone protocol  
connectionless

TCP is connection oriented so it have extra cost & extra delay

↳ no handshaking

apps which require no delay but can't afford unreliable data transfer prefer UDP

↳ e.g. streaming multimedia, DNS, even if we use UDP does not mean that reliability can not be provided

↳ we can add reliability but within the app each UDP segment is handled independently

↳ UDP does not ensure ordered delivery

UDP segment

|                        |          | 32 bits       |            |  |
|------------------------|----------|---------------|------------|--|
|                        |          | Source Port # | dest port# |  |
| length always in bytes | ← length |               | checksum   | header length =<br>$4 \times 16 \text{ bits} = 64 \text{ bits}$<br>8 bytes |
|                        |          | application   |            |  |
|                        |          | data          |            |  |

length :- total length of segment including header & data

data length = length - header length

checksum used for error detection

## UDP checksum

each segment is treated as a sum of 16 bit words

sender :- add all the words & take its one's complement  
to store in the checksum

receiver :- takes one's complement of the data received  
& match with the checksum value  
if matched == no error  
else error occurred.

## Example

data of 32 bits (4 words) will be divided into 2 16 bit words

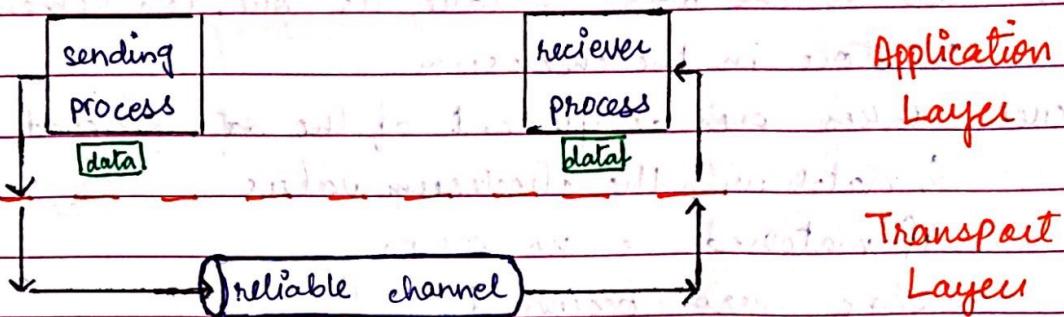
$$\begin{array}{r}
 \textcircled{1} \quad \textcircled{1} \quad \textcircled{1} \quad \textcircled{1} \\
 1 \ 1 \ 1 \ 0 \ 0 \ 1 \ 1 \ 0 \ 0 \ 1 \ 1 \ 0 \ 0 \ 1 \ 1 \ 0 \\
 + \\
 \hline
 \end{array}$$

wrapping at bit 11

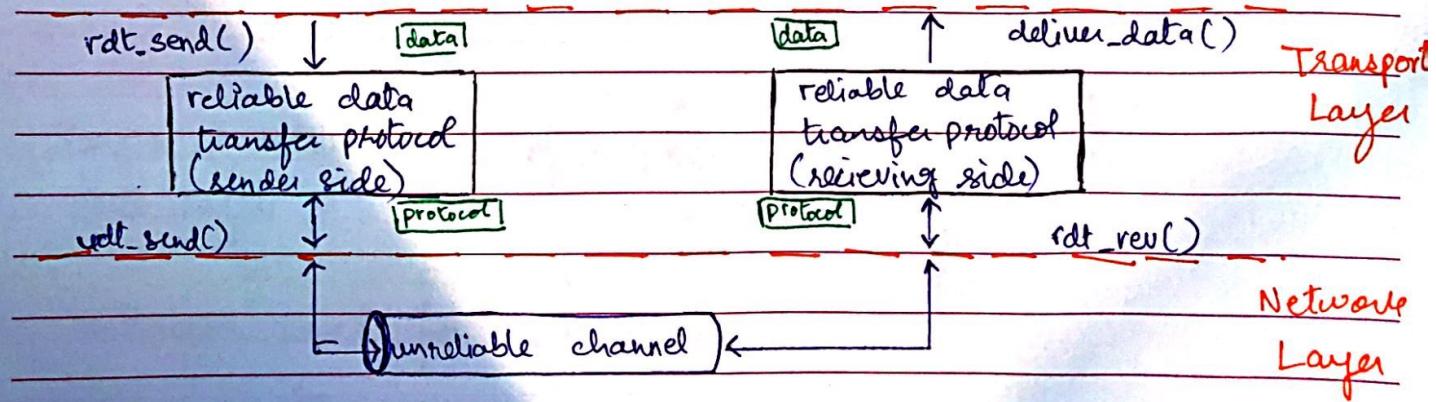
$$\begin{array}{r}
 \boxed{1} \ 0 \ 1 \ 1 \ 1 \ 0 \ 1 \ 1 \ 1 \ 0 \ 1 \ 1 \ 1 \ 0 \ 1 \ 1 \ 0 \ 1 \\
 \hline
 \text{Sum} \quad 1 \ 0 \ 1 \ 1 \ 0 \ 1 \ 1 \ 1 \ 0 \ 1 \ 1 \ 1 \ 0 \ 0
 \end{array}$$

Checksum 0 1 0 0 0 1 0 0 0 1 0 0 0 0 1 1

## Principles of Reliable Data Transfer



- in an ideal world the data will always be reliable  
there would be no data loss or corruption
- sending process will give the data to the transport layer who have to do nothing & it will give the data to the underlying channel
- transport layer will not need a reliable protocol or function because the channel is reliable.
- ! it will call an unreliable function & transfer the data unreliable because the channel is reliable
- receiver will also know that the data it received is reliable so it will not call checksum.



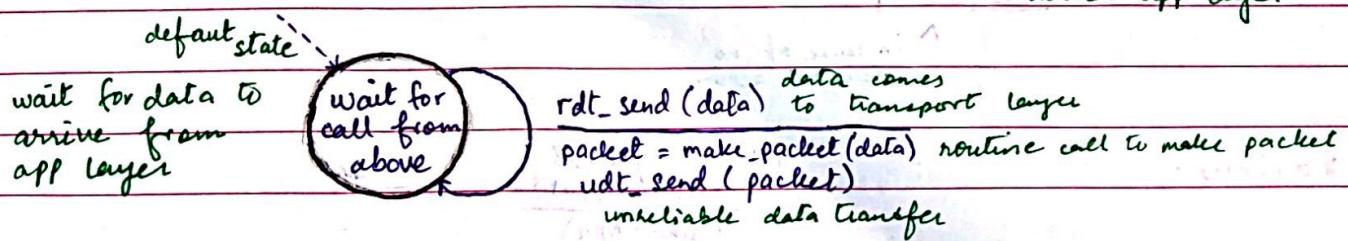
- application layer send data
- it calls a reliable data transfer function in transport layer
- transport layer also call a function in network layer which is unreliable because network layer have unreliable channel
- reliability is ensured at transport layer
- at the receiver side, transport layer receives the data through a reliable function because it has to ensure that data is sent to application layer reliably

rdt 1.0

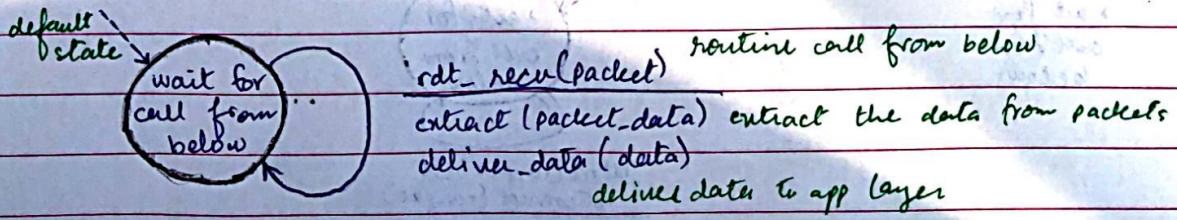
underlined channel is reliable

- no error
- no loss

sender &amp; receiver have separate FSM



### Sender



### Receiver

sender or receiver don't change state because they know that data will be transferred reliably, hence don't wait for any signal / acknowledgement

rdt 2.0

underlined channel is unreliable

↳ bit error

3 new functionalities to recover from error

↳ error detection

↳ checksum

↳ feedback

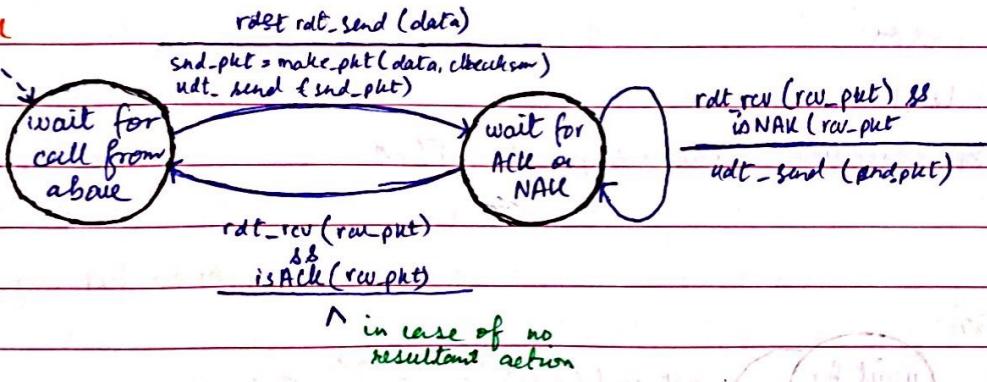
↳ ACK : Positive acknowledgment

↳ NAK : negative acknowledgment

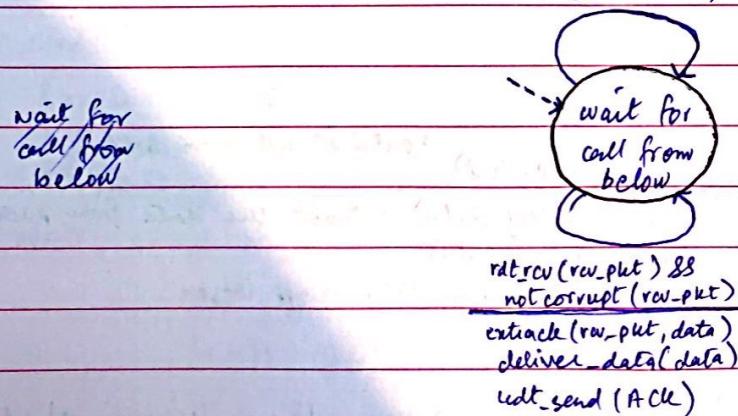
↳ retransmission

↳ packet in error is retransmitted by sender

### Sender



### Receiver



in case of no error

at sender data arrives from above  
packets are formed & checksum is attached  
sends data & changes state

at receiver data arrives from below

check the checksum; data not corrupted

data/packets are extracted

send ACK to sender

at sender, waiting for ACK or NAK from receiver

accepts ACK

changes state to default

in case of corrupted data

at sender side data arrives from the main path

packets are formed; checksum is attached

sends data & changes state

at receiver side data arrives

check the checksum; data is corrupted

send NAK to sender

at sender, waiting for ACK or NAK

accepts NAK

sends packet again

at receiver data arrives; check the checksum

if not corrupted, extract the data & sends ACK to sender



## Flaw in rdt 2.0

- ACK or NAK are bits
- they can also be corrupted
- if receiver sends ACK and on the way it is corrupted  $\rightarrow$  sender receives NAK and resends the packet
- receiver will get duplicated data
- to solve duplications sender adds sequence number
- in this way receiver can identify duplicate packets and discard them

## rdt 2.1

the flaw in 2.0 is solved in rdt 2.1  
add sequence number

- at sender, waits for data at seq#0
- data arrives and made into packets  
 $\text{snprintf} = \text{make\_pkt}(0, \text{data}, \text{checksum})$
- send data & change state
- $\rightarrow$  if sender receives NAK0 (data at receiver was corrupted)
  - resend the 0 packet again
  - if sender receives ACK0 from receiver  
*packet reached receiver without error*
  - changes state for 1 (next data with seq#1)
- $\rightarrow$  if receiver sends NAK  $\rightarrow$  becomes corrupted & sender gets ACK ; sender receives ACK 1
- sender was waiting for ACK0 but received ACK1  
hence it knows that ACK was corrupted
- retransmits the packet

Why sequence # 0 & 1 only?

## rdt 2.2

- Same functionality as 2.1
- remove NAK
- when receiver needs to send negative ACK, receiver will send ACK of the previous sequence no.
- receiver is in wait state 1 but received corrupt data for seq # 1 then receiver will continue to send ACK 0 until correct data for 1 is received
- when sender gets ACK 0 when it sends the packet with seq# 1 sender will know that receiver is not getting correct data for seq# 1

## rdt 3.0

- underlying channel with errors and loss
- in this case checksums, seq#, ACK, retransmission will be of no help
- in this case when sender sends a packet it will start a timer, if sender receives any ACK - then well in good but if sender does not receive any ACK before the timer expires it will send the packet again

• rdt 3.0 gives poor utilization

stop n wait protocol

- limiting the use of resources because the sender can not send another packet until it receives the acknowledgement of previous packet

$$U_{\text{sender}} = \frac{L/R}{RTT + L/R}$$

## Go-back-N

- window size = N
- consecutive transmitted; unACKed packets
- receiver sends cumulative ACK of last correctly received packet
- when sender receives ACK(n) window is forwarded to n+1
- sender has a timer for oldest unACKed packet
- when <sup>timer</sup> packet of a packet expires resend all unACKed packets
- seq # range depends on sequence number field in header.

bits allocated to the sequence # in header is the limit for the range of seq #

$$0, 2^k - 1$$

k :- no of bits allocated for the sequence # in header

e.g. if  $k=3$

$$2^k - 1 = 2^3 - 1$$

$$= 7$$

so 0-7 in binary

the window size would also be limited by this

$$0-7$$

$$N=8$$

if  $N=1$  then Go-back-N would become stop

wait since pipeline main file his packet

before sending next packet

base : sequence # of oldest in flight packet

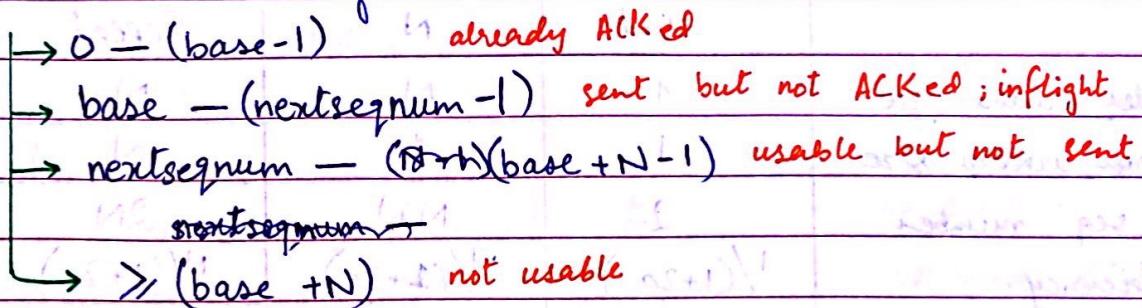
inflight packet : sent but not acked

nextseqnum : next available seq# which can be sent in pipeline



KAGHAZ  
www.kaghazpk

- There are 4 ranges



## Selective Repeat

receiver sends individual ACK of each packet

sender maintains a timer for each packet

receiver maintains a buffer of size N

↳ window size

## Dilemma

- suppose that there are seq # 0, 1, 2, 3
- and window size N is 3
- sender send pkt 0, 1, 2
- receiver receives them & send acknowledgment  
ack 0, ack 1, ack 2, & move the buffer along
- but the ACK's are lost before reaching the sender
- timer for packet 0 expires on sender side  
& it retransmits pkt 0
- at receiver side the buffer moved to 3, 0, 1
- since it does not know that ACK's were lost
- when pkt 0 reaches receiver, it does not know  
that it retransmission or a new new packet

Solution :- reduce window size

ideally window size should be less than or equal  
to half the sequence number space

$$N \leq \text{sequence number}/2$$

|                         | Stop n<br>Wait | Go Back<br>N  | Selective<br>Repeat  |
|-------------------------|----------------|---------------|----------------------|
| sender window size      | 1              | $(1-N)$       | $N$                  |
| receiver window size    | $(1+2a)$       | $(1+2a)(1-a)$ | $N$                  |
| min seq. number         | 2              | $N+1$         | $2N$                 |
| efficiency              | $1/(1+2a)$     | $N/(1+2a)$    | $N/(1+2a)$           |
| type of Ack             | individual     | cumulative    | individual           |
| order at receiving side | -              | inorder only  | out of order as well |
| num of retransmission   | $1$            | $N$           | $1$                  |

$a = \text{ratio of } D_{\text{prop}} \text{ and } D_{\text{trans}}$

- as Go Back N acks the pkt cumulatively, it rejects out of order pkts
- as Selective Repeat supports receiving out of order packets

↳ sorts the window after receiving packets

it uses independent acks



KAGHAZ  
www.kaghaz.pk

## Transmission Control Protocol

- connection oriented protocol

- inorder transmission

- full duplex

↳ bi-directional flow of data

- multicasting is not possible

↳ one sender, one receiver

- data flows in byte stream

not overload

receiver flow control & congestion control are diff

not overload  
the channel links

## TCP Segment

|                  |           | 32 bits       |   |             |   |                 |   |                       |   |                |                  |                           |                            |
|------------------|-----------|---------------|---|-------------|---|-----------------|---|-----------------------|---|----------------|------------------|---------------------------|----------------------------|
| header length    | fixed len | source port # |   | dest port # |   | sequence number |   | acknowledgment number |   | receive window |                  | sequence # of segment     |                            |
|                  |           | not used      | C | E           | U | A               | R | S                     | F | checksum       | Urg data pointer | options (variable length) | App Data (variable length) |
| 6 flags of 1 bit |           |               |   |             |   |                 |   |                       |   |                |                  |                           |                            |

- no. of bits available for seq# is the range of seq, in this case available bits = 32 ( $0, 2^{32}-1$ ) range of seq# all pipeline  $2^{32}-1$  segments ja saletay

- seq# is basically byte num, data stream buffer se byte grab karti aur segment mein dal deti



KAGHAZ  
www.kaghaz.pk

- header length is variable because the field of options in the segment is optional  
there are diff TCP options which when used will vary the TCP header length  
if this is not used then the header length is standard i.e 20 bytes

$$\text{Data length} = \text{Total segment size} - \text{header length}$$

- There are 6 flags in TCP segment of 1 bit
  - C, E → congestion notification
  - if A flag is set → ACK is included in the segment
  - & RST flag are used when connection is established or tear down
  - reset flag, sync flag, finish flag
- receive window includes the max # of bytes a receiver can process  
↳ this is used for flow control  
number of bytes receiver is willing to accept at any given point of time

Sequence Number

- TCP views data as byte stream after having fix.
- TCP app layer see data as it has been sent by sending side pr our sender buffer main ata ja raha hai
- TCP sender grabs chunks of data from the buffer & make packets/segments
- seq # have to be 0 in order for TCP to assume inorder delivery as network doesn't care about the order
- sequence number is the first byte in the first chunk of data which TCP grabbed

500 bytes data in buffer

MSS 1000 bytes

no. of segments =  $\frac{500}{1000} = 500$  segments

assuming, each segment size is 1000 bytes of 1000.

seq# of 1st segment 0 - 999

0 - 999 :- 1<sup>st</sup> segment seq# field / range

1000 - 1999 :- 2<sup>nd</sup> segment seq# range

$2000 + 2999 = 5000 \times (x-1) = 7296$

Acknowledgment Number

ACK for segment 0 - 999 will be 1000

- seq# of next expected byte from the sender
- this means that bytes till 999 have been correctly received and next expected byte is

$$1000 \times 3 = 7296 + 1 = 7297$$

17796 :- TCP sends cumulative ACK

• sender send 3 segments 0-999, 1000-1999, 2000-2999

TCP will send the ACK for the last correctly received byte

• ACK in this case will be 3000



KAGHAZ  
www.kaghaz.pk

## TCP RTT, timeout

- if timeout value < RTT

→ timer expires before ACK reaches sender  
retransmission

- if timeout value > RTT

→ after ACK gets lost, too long times expire

main note: retransmission nahi hogi.

→ slow reaction to ACK / segment lost

- optimal value of RTT

- when connection develops a sample RTT is stored

- this variable is used to set the timer

- during the connection, sample RTT is measured with some intervals & update the variable

- time is computed using the variable

SampleRTT is not measured for retransmissions.

$$\text{EstimatedRTT} = (1 - \alpha) * \text{EstimatedRTT} + \alpha * \text{sampleRTT}$$

where,  $\alpha = 0.02$  to  $0.125$  typically

*depends on network designer*

$$\text{Timeout Interval} = \text{EstimatedRTT} + 4 * \text{DevRTT}$$

$$\text{DevRTT} = (1 - \beta) * \text{DevRTT} + \beta * |\text{SampleRTT} - \text{EstimatedRTT}|$$

where,  $\beta = 0.25$  typically

TCP ACK generated at receiver

1. inorder segments arrive at receiver  
receiver waits up to 500ms for the next expected inorder segment so it can send a cumulative ACK
2. inorder segment arrive, along with another unACKed during the 500ms wait another inorder segment arrives at the receiver, it will immediately send a cumulative ACK
3. out of order segment arrives, higher than expected seq #  
receiver expects seq# 6, but receives seq# 8  
so this means that seq# 7 was lost so receiver will send the ACK for last correctly received segment i.e. 6
4. arrival of segment that fills the gap  
immediately sends the ACK  
receiver receives seq# 7 which fills the gap & send ACK for 7  
sends new ACK only if gap is filled at the lower end

# Chapter # 14

## Network Layer

**Data plane:** local per routing functions primarily forwarding moving a Datagram from an input link to an output link at a router

**Control plane:** network-wide logic, end-to-end getting packets one end of the network to the other end of the network coordination of devices

## Services and Protocols

**Sender Host :** network layer will take a transport layer segment from UDP or TCP and attach its header to it and encapsulate the segment in a datagram

**Receiver Host :** receives the datagram, checks info like the checksum, extracts the payload & multiplexes to the correct transport layer protocol

**Router :-** receives datagram from neighbouring host/router from the network layer & pass it to the appropriate router/host

## Network Layer Function

Local :- decision or action made at an individual router

Forwarding :- router local action

forwarding the datagram from a routers input port to its output port implemented in hardware

Global :- network wide

Routing :- global action

determining the route taken by a packet

from source to destination

implemented in software

Dataplane :- per router, per IP address device. Local function

Forwarding

Controlplane :- does network wide logic

Routing

Router :- the router matches the bits in datagram header with the bits in the forwarding table to determine the appropriate output link the datagram should be forwarded.

how does these local forwarding tables be computed?

→ entered by hand by the network developer

→ traditional routing algo

distributed routing algo is run in all the routers  
routing algo function communicates with the routing algo function in other router to compute values in forwarding table

→ Software Design Networking

remote software controller process, computes & distributes forwarding table used by each router



KAGHAZ  
www.kaghaz.pk

## Services

- guaranteed delivery with specific delay
- in-order delivery
- flow guaranteed in min bandwidth

## Best Effort Services

- simplicity
  - ↳ easy to add a new host, to add a new network
  - ↳ easy to manage
- provisioning of bandwidth
  - ↳ enough capacity allowed performance of real-time facilities
- Replicated
  - ↳ distributed application level infrastructure
  - ↳ services provided from multiple locations.

## Internet Protocol

about

→ datagram format

→ IP address structure

→ packet handling

# IP Datagram Format

|                                | 32 bits           |                   |                 |                      | - total datagram length |
|--------------------------------|-------------------|-------------------|-----------------|----------------------|-------------------------|
|                                | ver               | head len          | type of service | length               |                         |
|                                | 16 bit identifier |                   |                 | flgs fragment offset |                         |
|                                | time to live      | upper layer       | header checksum |                      |                         |
| max length = 64k bytes         |                   | source IP address |                 |                      | 32 bit                  |
| Typically < 1500 bytes or less |                   | dest IP address   |                 |                      | 32 bit                  |
|                                |                   | options (if any)  |                 |                      |                         |
|                                |                   | payload data      |                 |                      |                         |

ver :- IP protocol version number (4 bits)

IPv4 header

head len :- length of header

Data length :- Total Datagram - Header Length

Length

generally 20 bit if options field is not used

Type services :- to distinguish diff types of datagram from each other

ECN :- congestion notification

diff serv :- 6 types of diff bits

time to live :- decremented each time datagram passes through a router

if TTL reaches 0 the datagram must be dropped at the router

upper layer :- Transport layer protocol

6 :- TCP segment

17 :- UDP segment

day / date:

- 16 bit identifier, flags, fragment offset -

used when a large datagram is segments into multiple smaller datagrams

- header checksum - internet checksum computed over the contents of the IP header

checksum need to be computed at every router since the TTL is changed with every hop

20 bytes of TCP

20 bytes of IP

40 bytes + app layer overhead

## IP Addressing

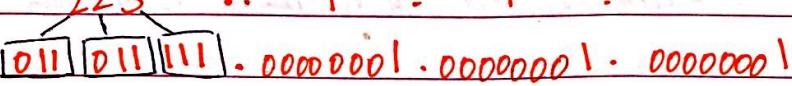
IP address does not identify a host or router per say rather over link layer interface of router or host

interface - connection b/w host/router & physical link  
router have more than 1 interfaces

multiple incoming & outgoing links

what more does router interfaces actually connected

IP address are in dotted-decimal notation

e.g. 223 .. 1 .. 1 ..  


Subnet router to router connection is also a subnet piece of a network that contains all devices that can reach each other without passing through intervening network layer devices routers that are directly connected to each other via some link layer technology

IP Address has 2 part

- Subnet part device in same subset have common high order bits

• Host part IP address will have diff host part which is the lower order bits

223.1.1.1  
all was in network "C" 3 bits subnet ID host ID

for example unhi ke liye FAST ki ID hai 223.1 jo usko I CAN se mil jayega aur local administrator isko subnets main divide karے ga

FAST ka network to identify ho gaya ab subnet specify karne ke liye bhi bits chahiye 1100 milegi host ID se (8 bits) last 2 bits jo host ID hi hain wo main se subnet ID aye gi aur host ID hi bits kum rkh jaye gi agar

Agen 3 subnets banane hain to 3 bits chahiye

$$2^x = 3$$

$$\therefore x = 2$$

host ID se 2 bits le le ge for subnet ID



KAGHAZ  
www.kaghaz.pk

day / date:

$$\text{max IP add} = 2^{32} = 4B$$

IP addresses in Class A, B, C have 2 levels

↳ IP address = net ID, host ID

↳ IP address = net ID, subnet ID, host ID

Delivery of IP address packets

↳ delivery to the site router

↳ delivery to the subnet router

↳ delivery to end host

## Subnet Mask

process of extracting the network address / net ID (if no subnetting is done) or subnet ID (if subnetting is done)

32 bit pattern

"1" in every net ID location & "1" in every in every subnet ID (if any) & "0" in every host ID location

e.g. 141.14.21.2 /  
10001101.00001110.0000010.00010101

subnet mask

## Defaults Masks

with subnetting

Class A 255.0.0.0

Class B 255.255.0.0

Class C 255.255.255.0

Net ID = IP address AND Subnet Mask

## CIDR

Classless InterDomain Routing

format : a.b.c.d/x

→ no. of bits of subnet part

there is no rigid boundary b/w host and network & subnet part

network IP ← → subnet part → host part →  
 200.23.16.0 / 23

host part : 9 bits

total IP add available :  $2^9$

usable :  $2^9 - 2$

host part which is all 0 / all 1

first & last bit is not usable

maanji thi IP address jos main 500 IP address hoon, unkon ne 500 ke closest  $2^9 = 512$  de di

day / date:

## DHCP

Dynamic Host Configuration Protocol  
dynamically obtain IP address from server when it joins the network using DHCP server  
when device leaves the network the IP is released so it can be reused

DHCP uses UDP

client uses port 68 server uses port 67

server uses port 67

ACK

port number used by client

port number used by server

client sends broadcast message to find server

server responds with unicast message

client sends broadcast message to find server

client sends broadcast message to find server

server responds with unicast message

client sends broadcast message to find server

port 68      port 67

port 68      port 67

port 68      port 67

client sends broadcast message to find server

server responds with unicast message

client sends broadcast message to find server

server responds with unicast message



# Chapter 3

## Flow Control

receiver tells the sender that rate at which segments need to be sent

receive window variable (rwnd)

$\text{LastByteRcvd} - \text{LastByteRead} \leq \text{RcvBuffer}$

last byte of  
data stream  
received from  
network

last byte in  
the stream to  
be read by  
app.

size of buffer

$$\text{rwnd} = \text{RcvBuffer} - [\text{LastByteRcvd} - \text{LastByte Read}]$$

$\boxed{\text{LastByteSent} - \text{LastByteAcked} \leq \text{rwnd}}$

variables maintained  
by sender

## Connection Management

TCP connection have 3 way handshake

SYN bit is used to ~~in~~ in 3 way handshake

SYN bit is set to 1

FIN bit is used to close the connection

FIN bit is set to 1

## TCP Tahoe

timeout, duplicates ACK, data loss

cwnd will always go to 1

## TCP Reno

timeout  $\rightarrow$  cwnd goes to 1

duplicate ACK  $\rightarrow$  cwnd will be half

ss thresh will be half in any case.

linear growth :- AIMD / congestion avoidance

exponential growth :- slow start

if cwnd becomes half :- duplicate ACK

fast recovery

Slow start - Congestion Avoidance - AIMD - Reno

Slow start  $\Rightarrow$  linear increase in throughput

Congestion avoidance  $\Rightarrow$  exponential increase

AIMD  $\Rightarrow$  constant throughput

Shannon's theorem

estimated capacity and bottleneck 937

congestion avoidance  $\Rightarrow$  throughput =  $\frac{1}{2}$  \* capacity

Fast Recovery  $\Rightarrow$  throughput = capacity

slow start and congestion avoidance 937

fast recovery  $\Rightarrow$  throughput = capacity

150 . 100 . 80 . 0 / 22

150 . 100 . 01010000 . 00000000

$\begin{array}{c} 001 \\ 010 \\ 011 \end{array}$  host

- convert into bits

$$\cdot 1/x \cdot 32 - x = \text{host ID bits}$$

$$\cdot \text{dekhlo kitne IP add} = y$$

$$2^x = y$$

$x$  = no of bits required for IP address

$$\cdot \text{Subset bits} = \text{host ID bits} - x$$

- update mask

$\rightarrow$  net ID + subnet

$\rightarrow$  mask  $32 - x$

# Chapter # 5

## Network Layer : Control Plane

### Routing Algorithms Classification

#### → Global

- compute paths at a global level using global information of the network
- link state algorithms

#### → Decentralized

- calculations based on their vicinity
- based on the information of their neighbours
- distance vector algorithms

#### → Static

- routes once calculated change very slowly over time

#### → Dynamic

- quick changes in response to link cost changes

How is the cost of a link computed

#### → cost of the link itself

e.g. fiber optic link will have high cost compared to submarine link

#### → congestion

e.g. high congestion on a link would mean higher cost so routers don't use that link



## Distance Vector Algorithm

### Bellman-Ford Equation

$D_x(y)$ : cost of least-cost path from  $x$  to  $y$

$$D_x(y) = \min_{v \in N(x)} \{ c_{xv} + D_v(y) \}$$

cost of neighbour  $v$  to dest  $y$

cost to neighbour  $v$

min taken over all neighbours  $v$  of  $x$

## Hierarchical Routing

- autonomous system (AS)

- an AS is typically under 1<sup>same</sup> administrative control

- routers in an AS will be running the same routing algorithm & will have information of each other

- Intra-AS :- <sup>Protocol</sup> within an AS

AS:- set of Internet routable IP prefixes (of same subnet)

belonging to a network or a collection of networks.

that are managed by a single entity or organization

gateway router:- edge router of an AS which helps in communication b/w two AS

- Inter-AS :- different AS can have different protocols communication b/w AS is done by gateway routers.



- one AS can have multiple gateway routers
- forwarding tables have both Intra-AS & Inter-AS routing protocols

In case of multiple gateway routers which how would a router know to which gateway router to forward the datagram?

there are 2 tasks performed by Inter-AS protocols

1. obtaining reachability info from neighboring AS  
must know which network is accessible through which gateway router

2. propagating this info to all internal routers

every subnet is connected to routers & these routers are connected to gateway routers

end host will a packet to its subnet router, now

how will this router know to which gateway router to forward that packet to reach it's destination.

This information will be provided by the gateway routers. (gateway routers bataye ge ke his network ya

subnet ke sare connected hain)

| 2 communicating AS must be running on the same Inter-AS routing protocol generally BGP.v4 |

• hierarchical routing saves table size

if all Internet networks is considered as a single set of routers, so all routers will be able to communicate with each other. This will result a flood of routing information on links that there will be no space for data. Secondly, every forwarding table will have millions of entries



## Inter-AS

## Intra- AS

Policy

admin controls traffic; policy can also control traffic which is not destined for it e.g. AS1 & AS3 communicate via AS2, so AS2 have control whether to allow it or not.

single admin so no policy decisions needed

.. performance focus (on) performance policy dominates performance since need to protect the network, change, & restrict the traffic

## Hot Potato Routing

- suppose a packet at AS1 needs to go to dest x
- x is reachable through AS2 and AS3
- the router at AS1 must decide which path to take
- sends the packet with the least cost link

## OSPF (open shortest path first)

usually implemented in upper tier ISP's

use Dijkstra's algorithm (link state)

open:- routing protocol specification is publicly available  
each router creates a map of the entire AS not only its neighbor

network administrator configures chooses the individual link cost

choose link cost to be one :- minimum hop routing

cost can be inversely proportional to the links capacity



security :- all messages are authenticated  
simple authentication  
MD5

multi-path :- multiple paths to the dest with same cost  
OSPF allow multiple paths to be used  
load sharing b/w paths

TOS :- for each link, multiple cost metric for different service can be used  
complex cost computation

multicast :- 1 to a group of nodes in AS

### Hierarchical OSPF

- AS can be configured hierarchically into areas
- each area runs its own OSPF routing algo with each other
- each router in an area broadcast its link states to all other routers in that area
  - 1 or more area border routers are responsible for routing packets outside that area
- exactly 1 OSPF area in the AS is configured to be the backbone area
- backbone area contains all the area border routers and may also contain non area border routers
- inter-area routing within AS packet is first routed to an area border router then it is routed to the backbone router and then to the dest area
- flooding of advertisement is contained only in its area
- broadcast link state every 30 mins or whenever a link state changes

checks if links are operational via "Hello messages" send to an attached neighbors

## BGP (Border Gateway Protocol)

- inter-AS routing protocol
- dominantly implemented over Internet
- RFC 4271, 4724

**functions of BGP**

- eBGP :- obtain subnet reachability info from neighboring ASes, may be connected directly or indirectly
- iBGP :- propagate this info to all the routers in the AS and not only the gateway router.
- determine "good" routes  
not only on basis of cost but also policy

**BGP session**

- BGP connection exist b/w 2 routers (a pair) **BGP peers**

semi-permanent TCP connection  
use port 179

exchange BGP messages

- BGP session b/w 2 AS's is called external BGP (eBGP)
- BGP session b/w routers in the same AS is called internal BGP (iBGP)
- BGP sessions do not always correspond to the physical links

not necessary that there is a BGP session when there is a direct link ; might exist when there is no direct link

- an AS is identified globally by a unique Autonomous System number and is listed in RFC 1930
- a stub AS do not need an AS number
  - ↳ an AS on which traffic originates or terminates.

- attributes of BGP

AS-Path :- identifies the ASes through which a route has passed ; used to detect & prevent looping advertisement

Next-Hop :- indicates specific internal-AS router to next hop AS.

router interface that begins AS-Path (**Next-Hop value**)

used by routers to properly configure forwarding table  
a little b/w inter & intra-AS protocol

## BGP Route Selection

more than 1 route to destination AS

1. local preference value (**Assigned by network administration**)

based on policy decision / cost

routes with the highest preference values are selected

if only 1 route is selected routing will be done on that, next steps won't be performed.

2. shortest AS-Path

3. closest Next-Hop router

hot-potato routing

least cost path

at this path generally a single path is filtered out

4. additional criteria

more BGP attributes

## SDN Control Plane

network's forwarding devices - packet switches

forwarding decisions are made on the basis of transport layer, network layer & link layer packet header fields



• logically centralized controller computes & distributes to the forwarding tables

- controller interacts with control agent (CA) of each router

- interacts using well-defined protocol to configure & manage routers flowtable

- Control Agent communicates with controller & do what it commands

CA's don't directly interact with each other

don't take part in computing forwarding table

• Characteristics of SDN

Flow-based forwarding

→ packet forwarding rules are specified in switch's flow table

SDN control plane compute, manage & install flow table entries

Segregation of data plane & control plane

dataplane consist of network's switches

control plane consist of servers & software that determine and manage switches flowtables.

→ Network control function

external to dataplane switches

controller maintains accurate network state info

provides this to the network-control applications

provide the means through which these apps can monitor, program & control the underlying network

• Programmable Network

network control application might determine end-end paths

specify & control dataplane in network device

- SDN control plane is divided into 2 components
  - SDN controller
  - SDN network-control applications

### SDN controller

- control the operation of a remote SDN enabled switch/host
- a protocol is needed to transfer info b/w controller & device
- communication b/w controller & device is known as controller's "southbound" interface
- controller require up-to-date info about the state of network hosts, links, switches & other devices
  - ↳ for control decisions
- controller might maintain a copy of flowtables
- controller interacts with network-control applications using "northbound" interface
  - ↳ allows network-control apps to read/write network state & flow tables

Routing

Access Control

Load Balancer

### Northbound API

Interface for network control apps

Network wide distributed, robust state management

### SDN Controller

Statistics

... Flowtables

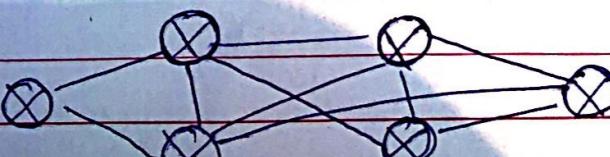
Link state info

... Host info

... Switch info

Communication to/from controlled device

### Southbound API



KAGHAZ  
www.kaghaz.pk

## Internet Control Message Protocol (ICMP)

- used mostly to signal error conditions b/w hosts & routers

- Example: ping / traceroute

- communicate network level info  
error reporting

unreachable host, protocol, network, port, etc

echo request

ping

- ICMP messages are carried directly inside payload

↳ UDP / TCP message are carried as payload in IP datagram

- upper layer protocol

- protocol number: 1

↳ used for demultiplexing up from IP

- message consists of

1 byte type field

1 byte code field

2 byte checksum field

upto 8 byte header field

of IP datagram that caused the message  
to be issued

e.g first 8 bytes of a datagram whose TTL exceeded

- Type: 11, code: 0 TTL expired

router received a datagram, decremented TTL,

& now TTL is 0



## Traceroute

- works by sending a set of 3 udp datagrams to dest
- 1st set is sent with an IP TTL field value set to 1
- 2nd set is sent with a TTL value of 2
- 3rd is sent with a TTL value of 3

⇒ when a router forwards datagram, always decrement the TTL value by 1

- when TTL value decremented to 0, datagram is dropped at that router
- The router may send an ICMP message back to the source indicating that TTL value has expired
- IP add of the message containing that ICMP TTL expired message is the IP address of the router where the packet was dropped
- if a packet have TTL value = n  
reply back from the router is the router which is n hops away
- Source also records the amount of time when it sent an IP datagram to the time when the corresponding ICMP message is received  
**measurement of the RTT from the host to that router**
- when segment finally reach the destination dest host may return ICMP "port unreachable" message type:3 , code:3  
the source knows the segment / packet reached the end of its path

- ICMP is a tool for network management



## Chapter 6

### Link Layer

- responsibility of transferring datagram from one node to physically adjacent node over a link
- encapsulate datagram in a frame & transfer it from 1 node to another adjacent node over link
- each link may use different protocols and each protocol provide different services

**Network layer - end to end**

**Link Layer : adjacent node**

### Link Layer Services

- encapsulate datagram into frame & add header AND trailer
- access channel can be a point-to-point link or a shared medium e.g. Ethernet
- reliable delivery between 2 adjacent nodes
- transport layer also provide rdt but between to processes
- don't use rdt on low-bit-error links but like fiber
- why both link layer & end-to-end reliability?
- efficiency is improved
- burden off TCP - TCP doesn't have to perform retransmissions.
- bit error corrects locally so TCP find the packet without bit error
- but in an end-to-end transmission some link might not provide reliability so rdt is ensured at TCP
- provide flow control
- error detection
- error correction
- half duplex & full duplex services

- most of its portion
- link layer is mostly implemented in the NIC card or the ethernet card or LAN card
- many of the link layer services are implemented in a chip at the heart of the NIC card
- controlling functionalities are implemented in the hardware
- the part of the link layer which implemented on the software runs on the hosts CPU

Most of the linklayer is implemented in a chip in the NIC card (hardware) rest of the link layer is implemented on the software which runs on the hosts CPU

## Error Detection

EDC :- error detection & correction

D :- Data

- Datagram is forwarded to the link layer, which attaches EDC bit to the payload, then this payload is forwarded to the link which may be prone to error or not

At the receiver side D' & EDC' is received which runs runs an error detection algorithm, if error occurs datagram is dropped else forwarded to the upcoming layer.

- more EDC bits, higher error detection reliability
- less EDC bits, lower error detection reliability
- There are primarily 3 error detection techniques
  - parity checking
  - internet checksum
  - CRC



## Parity Check

- single EDC parity bit is added
- even parity
- EDC bit is 1 :- even number of 1's in the data + parity bit
- EDC bit is 0 :- odd number of 1's in the data + parity bit
- parity bit is chosen so the total number of 1's in data + parity bit is even
- this is error prone ; is not suitable for bursts of error
- to overcome this we use 2D bit parity
- in this we convert a stream of bits in rows & columns
- each row & each col will have its own parity bit

## Internet Checksum

- same as we did in Transport layer
- add 2 bit words & take 1's complement of the sum

## Cyclic Redundancy Check (CRC)

- forward error correction :- when receiver detects error & resolve it
- backward error correction :- when receiver detects error but don't correct it, instead sends NAK to sender, sender retransmits
- D : data bits      R : CRC bits      G : pre agreed generator code
- receiver compute the generator code to check for errors
- D+R should be exactly divisible by G  
if remainder is 0 then no error
- G is an algebraic polynomial which needs to be converted into bits

## Multiple Access Protocols (MAC)

- There are 2 types of links

Point - to - Point

have 1 sender at one side of the link and 1 host at the other side

Broadcast

can be wired or unwired  
shared channel

- problem: simultaneous transmission by node over shared channel

solution: multiple access protol, algorithm to determine how nodes share a channel

### Ideal MAC characteristics

- one node should be able to utilise complete capacity of the link / rate =  $R$  bps
- if  $M$  nodes want to transmit, the rate would be divided  $R/M$  bps
- decentralized, no one node to decide

### Categories of MAC

channel partitioning

e.g TDM and FDM

random access

allow collisions but can recover from them  
taking turns



KAGHAZ  
www.kaghaz.pk

## Random Access Protocols

- transmitting at full rate  $R$  bps
- no prior agreement between nodes sharing a link
- when a node needs to transmit, it does so without checking if link is free which leads to collisions
- but there is mechanism for detection & recovery from collisions
- when collision is detected each node waits a random amount of time before transmitting again collision probability decrease

### 1- Slotted Aloha

- time is divided into slots
- 1 slot's length is atleast 1 frame wide
- 1 frame should be transmitted completely in 1 slot
- if a slot have data to transmit, it will for the beginning of a slot (always)
- Analogy : each node flip a coin ; tails :- won't transmit  
heads :- transmit in the beginning of next slot
- head corresponds to probability  $p$  (transmit)
- tail corresponds to probability  $(p-1)$  (do not transmit)
- each node transmits at full bandwidth
- max efficiency is 37%. (at best performance par bhi low efficiency)  
 $\hookrightarrow$  only 37% of the slots do useful work
- Probability a given node has success :-  $p(1-p)^{N-1}$
- Probability the remaining node do not transmit :-  $(1-p)^{N-1}$
- Probability that any of  $N$  nodes has a success :-  $Np(1-p)^{N-1}$   
 $\hookrightarrow$  efficiency
- max efficiency :-  $\lim_{N \rightarrow \infty} Np^*(1-p^*)^{N-1} \approx 1/e = 0.37$

### 2- Pure Aloha

- fully de-centralized, no synchronization, high collision rate
- node immediately transmits a packet
- Probability that any of  $N$  a given node has success :-  $p(1-p)^{2(N-1)}$

### 3- CSMA (carrier sense multiple access)

- before transmission node listen before transmitting if senses channel is idle it transmits else defer packet.
- collision can still occur nonetheless due to propagation delay
- node may not hear <sup>another</sup> node have transmitted, so it senses that channel is free and starts its own transmission
- propagation delay is introduced due to the distance b/w 2 nodes
- even after collision nodes transmit the whole packet/frame  
↳ no collision detection

### 4- CSMA/CD

- it further perform collision detection
- a transmitting node continues hearing on the channel and if detects another transmission, ceases its own transmission
- waits a random amount of time before retransmission
- ethernet uses CSMA/CD, and it is implemented in NIC
- the random time is chosen using binary (exponential) backoff
- node choose K randomly, from  $(0, 1, 2, \dots, 2^{n-1})$  where  $n$  is the number of collisions experienced previously by the node for that frame.
- the amount of time node waits is  $K \cdot 512$  bit times  
if  $K=1$ , it waits 512 bit times i.e. 512 ms for a 100 Mbps channel
- more collisions experienced, larger the interval from which K is chosen
- efficiency =  $\frac{1}{1 + S_{\text{prop}} / t_{\text{trans}}}$
- less  $t_{\text{prop}}$ , greater efficiency since colliding nodes will abort immediately efficiency = 1 (100%)
- greater  $t_{\text{prop}}$ , less efficiency



## Taking Turns

- using channel partition protocols and MAC protocols
- mechanism is used so nodes take turns to transmit
- 2 mechanisms
  - ↳ polling
  - ↳ token passing

### Polling

- one node on the channel is master & the rest are slaves
  - master controls the turns, assigns turns
  - usually polls using round robin fashion
  - goes around asking each node if it wants to transmit, and if a node does want to it is assigned a number a certain number of frames, after the limit moves to the next node.
  - if a master detects transmission is finished transmitting then move to the next node
- (disadvantage)*
- highly centralized
  - polling delay
- | example : bluetooth |*

### Token Passing

- there is a token exchanged b/w nodes, whichever node has the token it will transmit and only hold onto the token if it has any frame(s) to transmit using full capacity
- if the node which has token crashes then the whole system crashes since token will also be gone
- in this case, token recovery process is also needed

a router have 2 interfaces  
with diff IP add & MAC add

day / date:

- ARP is a simple request-response protocol
- Host A broadcasts a request packet containing IP address of B & MAC address FF.F.F.F.F.F.F
- all hosts/routers on the LAN will accept the frame & extract the datagram
- B will match the IP address with its own IP & send a response packet only to A containing its MAC address
- A will save the MAC address of B in its ARP table
- the rest of the nodes will discard the frame
- ARP table entry have 3 things IP add : MAC add : TTL  
TTL is normally 20 min after which the entry will be discarded from the table
- ARP request is broadcast
- ARP response is unicast

## Ethernet

- bus topology in mid 90's , broadcast LAN
- hub-based star topology in late 80's  
hub is a physical layer device which acts on bits not frames  
boosts the bits energy strength & transmit it onto other interfaces
- switched ethernet prevails today  
active switch in centre rather than a hub

## Frame Structure

| preamble | dest address | source address | type | data (payload) | CRC |
|----------|--------------|----------------|------|----------------|-----|
|----------|--------------|----------------|------|----------------|-----|

- data contains IP datagram ;  $46 \text{ bytes} \leq \text{payload} \leq 1500 \text{ bytes}$
- if a payload  $< 46 \text{ bytes}$  then ethernet protocol stuff bytes in the payload to make it a minimum of 46 bytes
- if payload  $> 1500 \text{ bytes}$  network layer perform fragmentation

- dest address is a 6 byte field, contains MAC add of dest
- source address is 6 byte field, contains MAC add of source
- type is a 2byte field, used to multiplex network layer protocol  
a host may support multiple network-layer protocols for diff applications  
so the adapter at any host needs to know to which network layer protocol to should it pass the frame, this is identified in type field
- CRC is a 4 byte field
- Preamble is a 8 byte field, 56 bits of alternating 1's & 0's  
first 7 bytes 10101010 (always)  
last byte is SFD (8 bits) which is also alternating except last 2 bits  
10101011
- 1st 7 bytes are used to wakeup the receiver's adapter & to synchronize with sender's clock
- clocks are out of synchronisation because there will always be a drift b/w target rate & transmitted rate, this drift is not known by the receiver or any other adapter on the LAN
- the last 11 of the SFD tells the receiver that preamble byte is about to end and actual /important stuff is coming

## Ethernet Switches

- layer 2 switches, network-to-link layer device
- do not need to configured unlike routers
- perform 2 functions - filtering & forward  
determines whether to forward the datagram or filter/drop  
if forward then onto which interface to forward
- each switch table entry have 3 things  
MAC address  
interface which leads to that MAC address  
TTL / time at which the entry was made

- suppose a frame arrives at a switch at interface  $x$
- there are 3 possible cases
  - there is no entry for the dest MAC address in the table  
in this case the switch broadcasts the frame to all interfaces except  $x$
  - 2- there is an entry in the table associating dest add to interface  $x$ , so there is no need to forward the frame, and filter function is performed ( $\text{src add} == \text{dest add}$ )
  - 3- there is an entry associating dest add with interface other than  $y$ , selective forwarding is performed for that interface only.

### Self Learning

- switching table is initially empty
- the frames received by the switch at any interface source MAC add, interface, time it arrives is stored in the table
- in this way it populates the table
- switch deletes an entry after sometime if no frames are received with that address as source this time period is known as aging time
- there are no collisions
- full duplex
- allow simultaneous switching (transmission) without collisions as long as dest is different
- transmission with same dest add can not happen simultaneously

- suppose a frame arrives at a switch at interface  $x$
- there are 3 possible cases
  - 1- there is no entry for the dest MAC address in the table  
in this case the switch broadcasts the frame to all interfaces except  $x$
  - 2- there is an entry in the table associating dest add to interface  $x$ , so there is no need to forward the frame, and filter function is performed ( $\text{src add} == \text{dest add}$ )
  - 3- there is an entry associating dest add with interface other than  $y$ , selective forwarding is performed for that interface only.

### Self Learning

- switching table is initially empty
- the frames received by the switch at any interface source MAC add, interface, time it arrives is stored in the table
- in this way it populates the table
- switch deletes an entry after sometime if no frames are received with that address as source this time period is known as aging time
- there are no collisions
- full duplex
- allow simultaneous switching (transmission) without collisions as long as dest is different
- transmission with same dest add can not happen simultaneously